



Bergische Universität  
Wuppertal

# Umfangreiche C-XSC-Langzahlpakete für beliebig genaue reelle und komplexe Intervallrechnung

Frithjof Blomquist, Werner Hofschuster, Walter Krämer

Preprint

**BUW-WRSWT 2012/2**

Nachfolger des Preprints BUW-WRSWT 2011/1

Wissenschaftliches Rechnen/  
Softwaretechnologie



Diese Arbeit ist entstanden im Rahmen des Projektes  
*C-XSC Schnittstelle zur MPFR- und MPFI-Bibliothek*

## Impressum

Herausgeber: Prof. Dr. W. Krämer, Dr. W. Hofschuster Wissenschaftliches Rechnen/Softwaretechnologie Fachbereich C (Mathematik und Naturwissenschaften) Bergische Universität Wuppertal Gaußstr. 20 42097 Wuppertal, Germany
--

## Internet-Zugang

Die Berichte sind in elektronischer Form erhältlich über die World Wide Web Seiten

<http://www2.math.uni-wuppertal.de/wrswt/literatur.html>

## Autoren-Kontaktadresse

Frithjof Blomquist  
Adlerweg 6  
D-66436 Püttlingen  
E-mail: [blomquist@math.uni-wuppertal.de](mailto:blomquist@math.uni-wuppertal.de)

Werner Hofschuster  
Bergische Universität Wuppertal  
Gaußstr. 20  
D-42097 Wuppertal  
E-mail: [hofschuster@math.uni-wuppertal.de](mailto:hofschuster@math.uni-wuppertal.de)

Walter Krämer  
Bergische Universität Wuppertal  
Gaußstr. 20  
D-42097 Wuppertal  
E-mail: [kraemer@math.uni-wuppertal.de](mailto:kraemer@math.uni-wuppertal.de)

# Inhaltsverzeichnis

<b>1. Einleitung</b>	<b>11</b>
<b>2. Installation</b>	<b>13</b>
2.1. Installation der MPFR- und MPFI-Bibliotheken	13
<b>3. MpfrClass-Interface zur Anbindung der MPFR-Bibliothek an C-XSC</b>	<b>15</b>
3.1. Grundlegendes	15
3.1.1. Allgemein	15
3.1.2. Aufbau	15
3.1.3. Präzision	15
3.1.4. Variablentyp PrecisionType	15
3.1.5. Variablentyp RoundingMode	16
3.1.6. Zahlenformat	17
3.2. Konstruktoren / Destruktor	18
3.2.1. Konstruktoren	18
3.2.2. Destruktor	18
3.3. Eingabe / Ausgabe	19
3.4. Rundungsmodi und Precision Handling	20
3.5. Anwendungsprogramm	21
3.6. Typ-Umwandlungen	22
3.6.1. <code>real</code> , <code>double</code> , ... → MPFR	22
3.6.2. MPFR → <code>real</code> , <code>double</code> , ...	22
3.6.3. MPFR → <code>mpfr_t</code>	22
3.6.4. MPFR ↔ <code>mpfr_t</code>	22
3.6.5. <code>mpfr_t</code> → MPFR	22
3.6.6. <code>string</code> → MPFR	23
3.6.7. MPFR → <code>string</code>	23
3.6.8. MPFR → MPFR	23
3.6.9. Verschiedenes	24
3.7. Zuweisungs-Operatoren	25
3.8. Abfragen	25
3.9. Vergleiche	26
3.9.1. Vergleichsfunktionen	26
3.9.2. Vergleichsoperatoren <code>=</code> , <code>≠</code> , <code>&gt;</code> , <code>≥</code> , <code>&lt;</code> , <code>≤</code>	26
3.10. Arithmetische Operatoren	28
3.10.1. Addition	28
3.10.2. Subtraktion	28
3.10.3. Multiplikation	29
3.10.4. Division	29
3.10.5. Vom Current-Rundungsmodus abweichende Rundungen	30
3.10.6. Arithmetik-Funktionen zum Auf- und Abrunden	31
3.11. Mathematische Funktionen	32
3.11.1. Standard-Implementierung	32
3.11.2. Davon abweichende Funktionen und Konstanten	33
3.11.3. Elementarfunktionen	37

3.11.4. Skalarprodukt aus zwei Teilprodukten . . . . .	40
3.11.4.1. Beispiel . . . . .	42
3.11.5. Funktionen der Mathematischen Physik . . . . .	44
<b>4. MpfiClass-Interface zur Anbindung der MPFI-Bibliothek an C-XSC</b>	<b>47</b>
4.1. MPFI-Bibliothek . . . . .	47
4.1.1. Entwickler . . . . .	47
4.1.2. Allgemein . . . . .	47
4.1.3. Installation . . . . .	47
4.2. Grundlegendes . . . . .	48
4.2.1. Allgemein . . . . .	48
4.2.2. Aufbau . . . . .	48
4.2.3. Präzision . . . . .	48
4.2.4. Variablentyp PrecisionType . . . . .	48
4.3. Konstruktoren / Destruktoren . . . . .	49
4.3.1. Konstruktoren . . . . .	49
4.3.2. Destruktor . . . . .	49
4.4. Leeres Intervall . . . . .	49
4.5. Zuweisungs-Operatoren . . . . .	50
4.6. Arithmetische Operatoren . . . . .	51
4.6.1. Addition . . . . .	51
4.6.2. Subtraktion . . . . .	52
4.6.3. Multiplikation . . . . .	52
4.6.4. Division . . . . .	53
4.7. Eingabe / Ausgabe . . . . .	54
4.8. Base und Precision Handling . . . . .	55
4.9. Typ-Umwandlungen . . . . .	56
4.9.1. <code>real, double, ...</code> → MPFI . . . . .	56
4.9.2. MPFI → <code>interval</code> . . . . .	56
4.9.3. MPFI → <code>mpfi_t</code> . . . . .	56
4.9.4. MPFI ↔ <code>mpfi_t</code> . . . . .	56
4.9.5. <code>mpfi_t</code> → MPFI . . . . .	57
4.9.6. <code>string</code> → MPFI . . . . .	57
4.9.7. MPFI → <code>string</code> . . . . .	57
4.9.8. MPFI → MPFI . . . . .	57
4.10. Abfragen . . . . .	58
4.11. Vergleiche . . . . .	59
4.11.1. Vergleichsfunktionen . . . . .	59
4.11.2. Vergleichsoperatoren <code>=, ≠, &gt;, ≥, &lt;, ≤</code> . . . . .	59
4.12. Durchschnitt . . . . .	62
4.13. Konvexe Hülle . . . . .	63
4.14. Mathematische Intervall-Funktionen . . . . .	64
4.14.1. Standard-Implementierung . . . . .	64
4.14.2. Davon abweichende Funktionen und Konstanten . . . . .	65
4.14.3. Elementarfunktionen . . . . .	68
4.14.4. Funktionen der Mathematischen Physik . . . . .	71
4.14.5. Skalarprodukt aus zwei Intervallprodukten . . . . .	72
4.15. Optimale Intervall-Einschließungen . . . . .	74
<b>5. MpfcClass-Interface für komplexe Langzahlrechnungen in C-XSC</b>	<b>75</b>
5.1. Grundlegendes . . . . .	75
5.1.1. Allgemein . . . . .	75

5.1.2.	Aufbau	75
5.1.3.	Präzision	75
5.1.4.	Variablentyp PrecisionType	75
5.1.5.	Variablentyp RoundingMode	75
5.2.	Rundungsmodi und Precision Handling	77
5.3.	Konstruktoren / Destruktor	78
5.3.1.	Konstruktoren	78
5.3.2.	Destruktor	79
5.4.	Zuweisungs-Operatoren	80
5.5.	Eingabe / Ausgabe	81
5.6.	Typ-Umwandlungen	82
5.6.1.	real, double, complex, ... → MPFC	82
5.6.2.	MPFC → complex	82
5.6.3.	MPFC → mpfr_t	82
5.6.4.	mpfr_t → MPFC	83
5.6.5.	MPFC → string	83
5.6.6.	string → MPFC	84
5.6.7.	MPFC → MPFC	84
5.7.	Abfragen	85
5.8.	Vergleiche	86
5.8.1.	Vergleichsoperatoren =, !=	86
5.9.	Arithmetische Operatoren	87
5.9.1.	Addition	87
5.9.2.	Subtraktion	88
5.9.3.	Multiplikation	89
5.9.4.	Division	89
5.10.	Mathematische Funktionen	90
5.10.1.	Standard-Implementierung	90
5.10.2.	Davon abweichende Funktionen	90
5.10.3.	Elementarfunktionen	92
<b>6.</b>	<b>MpfcClass-Interface für komplexe Langzahl-Intervallrechnungen in C-XSC</b>	<b>95</b>
6.1.	Grundlegendes	95
6.1.1.	Allgemein	95
6.1.2.	Aufbau	95
6.1.3.	Präzision	95
6.1.4.	Variablentyp PrecisionType	95
6.2.	Rundungs und Precision Handling	96
6.3.	Konstruktoren / Destruktor	97
6.3.1.	Konstruktoren	97
6.3.2.	Destruktor	98
6.4.	Zuweisungs-Operatoren	99
6.5.	Eingabe / Ausgabe	100
6.6.	Typ-Umwandlungen	101
6.6.1.	real, interval, ... → MPFCI	101
6.6.2.	MPFCI → cinterval	101
6.6.3.	MPFCI → mpfi_t	101
6.6.4.	mpfi_t → MPFCI	102
6.6.5.	string → MPFCI	102
6.6.6.	MPFCI → string	102
6.6.7.	MPFCI → MPFCI	103
6.7.	Abfragen	103

6.8.	Vergleiche . . . . .	104
6.8.1.	Vergleichsoperatoren =, !=, <, <=, . . . . .	104
6.9.	Durchschnitt . . . . .	107
6.10.	Konvexe Hülle . . . . .	108
6.11.	Arithmetische Operatoren . . . . .	109
6.11.1.	Addition . . . . .	109
6.11.2.	Subtraktion . . . . .	110
6.11.3.	Multiplikation . . . . .	111
6.11.4.	Division . . . . .	112
6.12.	Mathematische Funktionen . . . . .	113
6.12.1.	Standard-Implementierung . . . . .	113
6.12.2.	Davon abweichende Funktionen . . . . .	113
6.12.3.	Elementarfunktionen . . . . .	115
6.12.4.	Optimale, komplexe Intervall-Einschließungen . . . . .	119
<b>7.</b>	<b>Anwendungen</b>	<b>121</b>
7.1.	Erste Nullstelle von $J_0(x)$ . . . . .	121
7.2.	Einschließung reeller arithmetischer Ausdrücke . . . . .	124
7.3.	Automatische Differentiation . . . . .	125
7.4.	Globale Optimierung . . . . .	131
7.4.1.	Optimale Einschließung des Wertebereichs analytischer Funktionen . . . . .	135
7.5.	Reelle, eindimensionale Taylor Arithmetik . . . . .	139
7.5.1.	Elementarfunktionen . . . . .	140
7.5.2.	Konstruktoren / Destruktoren . . . . .	141
7.5.2.1.	Konstruktoren . . . . .	141
7.5.2.2.	Destruktor . . . . .	141
7.5.3.	Zuweisungs-Operatoren . . . . .	141
7.5.4.	Arithmetische Operatoren . . . . .	142
7.5.5.	Unabhängige Variablen und Konstanten . . . . .	142
7.5.6.	Precision Handling . . . . .	142
7.5.7.	Zugriff auf die Vektorkomponenten . . . . .	143
7.5.8.	Ausgabe . . . . .	143
7.5.9.	Programme . . . . .	144
7.6.	Komplexe, eindimensionale Taylor Arithmetik . . . . .	147
7.6.1.	Elementarfunktionen . . . . .	148
7.6.2.	Konstruktoren / Destruktoren . . . . .	149
7.6.2.1.	Konstruktoren . . . . .	149
7.6.2.2.	Destruktor . . . . .	149
7.6.3.	Zuweisungs-Operatoren . . . . .	150
7.6.4.	Arithmetische Operatoren . . . . .	150
7.6.5.	Unabhängige Variablen und Konstanten . . . . .	150
7.6.6.	Precision Handling . . . . .	151
7.6.7.	Zugriff auf die Vektorkomponenten . . . . .	151
7.6.8.	Ausgabe . . . . .	152
7.6.9.	Programme . . . . .	152
7.7.	Nullstellen komplexer Ausdrücke . . . . .	155
7.8.	Einschließung aller Nullstellen nichtlinearer Funktionen . . . . .	158
<b>A.</b>	<b>Neue Funktionen vom Typ MpfrClass</b>	<b>179</b>
A.1.	Grundregeln für garantierte Rundungen . . . . .	179
A.1.1.	Unitäre Operatoren . . . . .	180
A.1.2.	Addition . . . . .	180

A.1.3.	Subtraktion . . . . .	180
A.1.4.	Multiplikation . . . . .	180
A.1.5.	Division . . . . .	181
A.2.	$x^2 - y^2, x^2 + y^2$ . . . . .	184
A.3.	$x/(x^2 + y^2)$ . . . . .	185
A.4.	$(x^2 - y^2)/(x^2 + y^2)^2$ . . . . .	186
A.4.1.	Numerische Beispiele . . . . .	186
A.5.	$2xy/(x^2 + y^2)^2$ . . . . .	187
A.5.1.	Numerische Beispiele . . . . .	187
A.6.	$2xy/(4x^2y^2 + (1 + x^2 - y^2)^2)$ . . . . .	188
A.6.1.	Numerische Beispiele . . . . .	188
A.7.	$(1 + x^2 - y^2)/(4x^2y^2 + (1 + x^2 - y^2)^2)$ . . . . .	189
A.7.1.	Numerische Beispiele . . . . .	190
A.8.	$x^2 + a \cdot x + b$ . . . . .	191
A.8.1.	Numerische Beispiele . . . . .	191
A.9.	$2\sin(x) \cosh(y)/(\cosh(2y) - \cos(2x))$ . . . . .	192
A.9.1.	Optimale Einschließung . . . . .	192
A.9.2.	Numerische Beispiele . . . . .	194
A.10.2	$\cos(x) \sinh(y)/(\cos(2x) - \cosh(2y))$ . . . . .	195
A.10.1.	Optimale Einschließung . . . . .	195
A.11.	$\sqrt{x^2 - 1}$ . . . . .	197
A.12.	$x/\sqrt{1 + x^2}$ . . . . .	198
A.13.	$x/\sqrt{1 - x^2}$ . . . . .	198
A.14.	$x/\sqrt{x^2 - 1}$ . . . . .	198
A.15.	$\ln(\sin(x))$ . . . . .	199
A.16.	$\ln(\sqrt{x^2 + y^2})$ . . . . .	201
A.17.	$\ln(\sqrt{(1 + x)^2 + y^2})$ . . . . .	201
A.18.	$\operatorname{arcosh}(1 + x)$ . . . . .	202
A.19.	$\Gamma'(x)$ . . . . .	203
A.20.1	$1/\Gamma(x)$ . . . . .	203
A.21.	$(1/\Gamma(x))'$ . . . . .	204
<b>B.</b>	<b>Neue Funktionen vom Typ MpfClass</b>	<b>205</b>
B.1.	$\ln(\sqrt{(1 + x)^2 + y^2})$ . . . . .	205
B.2.	$x/(x^2 + y^2)$ . . . . .	208
B.3.	$(x^2 - y^2)/(x^2 + y^2)^2$ . . . . .	210
B.3.1.	Maximumbestimmung . . . . .	213
B.3.2.	Minimumbestimmung . . . . .	214
B.3.3.	Numerische Beispiele . . . . .	214
B.4.	$2xy/(x^2 + y^2)^2$ . . . . .	216
B.4.1.	Maximumbestimmung . . . . .	223
B.4.2.	Minimumbestimmung . . . . .	224
B.4.3.	Numerische Beispiele . . . . .	225
B.5.	$x^2 + a \cdot x + b$ . . . . .	226
B.5.1.	Numerische Beispiele . . . . .	226
B.6.	$(1 + x^2 - y^2)/(4x^2y^2 + (1 + x^2 - y^2)^2)$ . . . . .	227
B.6.1.	Maximumbestimmung . . . . .	228
B.6.2.	Minimumbestimmung . . . . .	231
B.6.3.	Numerische Beispiele . . . . .	232
B.7.	$2xy/(4x^2y^2 + (1 + x^2 - y^2)^2)$ . . . . .	233
B.7.1.	Maximumbestimmung . . . . .	236
B.7.2.	Minimumbestimmung . . . . .	239

B.7.3. Numerische Beispiele . . . . .	240
B.8. $x/\sqrt{1+x^2}$ . . . . .	241
B.9. $x/\sqrt{1-x^2}$ . . . . .	241
B.10. $x/\sqrt{x^2-1}$ . . . . .	241
<b>C. Elementarfunktionen für komplexe Punkt- und Intervallargumente</b>	<b>243</b>
C.1. Elementarfunktionen für komplexe Punktargumente . . . . .	248
C.1.1. Exponentialfunktion, Realteil . . . . .	248
C.1.2. $\sin(z)$ . . . . .	249
C.1.3. $\cos(z)$ . . . . .	249
C.1.4. $\tan(z)$ . . . . .	249
C.1.5. $\cot(z)$ . . . . .	249
C.1.6. $\arg(z)$ . . . . .	250
C.1.7. $\operatorname{argpl}(z)$ . . . . .	251
C.1.8. $ z $ . . . . .	252
C.1.9. $\log(z)$ . . . . .	252
C.1.10. $1/z$ . . . . .	253
C.1.10.1. Realteil . . . . .	253
C.1.10.2. Imaginärteil . . . . .	253
C.1.11. $1/z^2$ . . . . .	254
C.1.11.1. Realteil . . . . .	254
C.1.11.2. Imaginärteil . . . . .	254
C.1.11.3. Numerische Ergebnisse . . . . .	254
C.1.12. $1/(1+z^2)$ . . . . .	255
C.1.12.1. Realteil . . . . .	255
C.1.12.2. Imaginärteil . . . . .	255
C.1.12.3. Numerische Ergebnisse . . . . .	255
C.1.13. $z^2$ . . . . .	256
C.1.14. $\sqrt{z}$ . . . . .	256
C.1.15. $z^2 + a \cdot z + b$ . . . . .	257
C.1.16. $1/\sqrt{z}$ . . . . .	258
C.1.17. $\sinh(z)$ . . . . .	260
C.1.18. $\cosh(z)$ . . . . .	260
C.1.19. $\tanh(z)$ . . . . .	260
C.1.20. $\operatorname{coth}(z)$ . . . . .	260
C.1.21. $\arcsin(z)$ . . . . .	261
C.1.21.1. Realteil . . . . .	261
C.1.21.2. Imaginärteil . . . . .	263
C.1.22. $\arccos(z)$ . . . . .	265
C.1.23. $1/(1-z^2)$ . . . . .	265
C.1.24. $\log(1+z)$ . . . . .	266
C.1.24.1. Realteil . . . . .	267
C.1.24.2. Imaginärteil . . . . .	270
C.1.24.3. Numerische Ergebnisse . . . . .	270
C.1.25. $e^z - 1$ . . . . .	271
C.1.25.1. Realteil . . . . .	272
C.1.25.2. Imaginärteil . . . . .	272
C.1.25.3. Numerische Ergebnisse . . . . .	272
C.1.26. $\sqrt{z+1} - 1$ . . . . .	274
C.1.26.1. Numerische Ergebnisse . . . . .	277
C.1.27. $\sqrt{1+z^2}$ . . . . .	279
C.1.27.1. Realteil . . . . .	279



C.1.27.2. Imaginärteil . . . . .	282
C.1.28. $\sqrt{1-z^2}$ . . . . .	283
C.1.29. $\sqrt{z^2-1}$ . . . . .	284
C.1.29.1. Realteil . . . . .	285
C.1.29.2. Imaginärteil . . . . .	285
C.1.30. $\sqrt{-z^2-1}$ . . . . .	286
C.2. Elementarfunktionen für komplexe Intervallargumente . . . . .	287
C.2.1. $z^2$ . . . . .	287
C.2.2. $1/z$ . . . . .	287
C.2.3. $1/z^2$ . . . . .	288
C.2.3.1. Numerische Ergebnisse . . . . .	288
C.2.4. $1/(1+z^2)$ . . . . .	289
C.2.4.1. Numerische Ergebnisse . . . . .	289
C.2.5. $1/(1-z^2)$ . . . . .	290
C.2.5.1. Numerische Beispiele . . . . .	290
C.2.6. $z^2 + a \cdot z + b$ . . . . .	291
C.2.6.1. Numerische Beispiele . . . . .	291
C.2.7. $\log(z)$ . . . . .	293
C.2.7.1. Analytische Funktion . . . . .	294
C.2.7.2. Nicht-analytische Funktion . . . . .	295
C.2.8. $\log(1+z)$ . . . . .	296
C.2.8.1. Analytische Funktion . . . . .	297
C.2.8.2. Nicht-analytische Funktion . . . . .	299
C.2.9. $\sqrt{z}$ . . . . .	300
C.2.10. $\sqrt{z}$ , Beide Quadratwurzeln . . . . .	301
C.2.11. $\sqrt[n]{z}$ . . . . .	302
C.2.12. $\sqrt[n]{z}$ Alle Wurzeln . . . . .	304
C.2.13. $\sqrt{z+1}-1$ . . . . .	306
C.2.13.1. Realteil . . . . .	307
C.2.13.2. Imaginärteil . . . . .	308
C.2.13.3. Numerische Ergebnisse . . . . .	308
C.2.14. $\sqrt{1+z^2}$ . . . . .	310
C.2.14.1. Realteil . . . . .	311
C.2.14.2. Imaginärteil . . . . .	312
C.2.14.3. Numerische Ergebnisse . . . . .	313
C.2.15. $\sqrt{1-z^2}$ . . . . .	314
C.2.16. $\sqrt{z^2-1}$ . . . . .	315
C.2.16.1. Realteil . . . . .	316
C.2.16.2. Imaginärteil . . . . .	317
C.2.16.3. Numerische Ergebnisse . . . . .	318
C.2.17. $\sqrt{-z^2-1}$ . . . . .	319
C.2.17.1. Numerische Ergebnisse . . . . .	320
C.2.18. $1/\sqrt{z}$ . . . . .	321
C.2.18.1. Realteil . . . . .	321
C.2.18.2. Imaginärteil . . . . .	324
C.2.18.3. Numerische Beispiele . . . . .	327
C.2.19. $\cot(z)$ . . . . .	328
C.2.20. $\arcsin(z)$ . . . . .	329
C.2.20.1. Algorithmus . . . . .	330
C.2.21. $\arccos(z)$ . . . . .	339
C.2.21.1. Algorithmus . . . . .	340

C.2.22. $\arctan(z)$ . . . . .	354
C.2.22.1. Algorithmus . . . . .	355
C.2.23. $\operatorname{arccot}(z)$ . . . . .	362
C.2.24. $\operatorname{arsinh}(z)$ . . . . .	363
C.2.25. $\operatorname{arcosh}(z)$ . . . . .	364
C.2.26. $\operatorname{artanh}(z)$ . . . . .	365
C.2.27. $\operatorname{arcoth}(z)$ . . . . .	366
C.2.28. $z^p$ , $p \in \mathbb{P}:\text{Mpf}i\text{Class}$ . . . . .	367
C.2.29. $e^z - 1$ . . . . .	370
C.2.29.1. Realteil . . . . .	370
C.2.29.2. Imaginärteil . . . . .	371
C.2.29.3. Numerische Ergebnisse . . . . .	371
<b>D. Einschließung reeller und komplexer Ausdrücke</b>	<b>373</b>
D.1. Einschließung reeller arithmetischer Ausdrücke . . . . .	373
D.1.1. Problemstellung und Definitionen . . . . .	373
D.1.2. Verbesserung einer Einschließung . . . . .	374
D.1.3. Optimale Einschließungen . . . . .	374
D.2. Einschließung komplexer arithmetischer Ausdrücke . . . . .	377
D.2.1. Problemstellung und Definitionen . . . . .	377
D.2.2. Verbesserung einer Einschließung . . . . .	378
D.2.3. Optimale Einschließungen . . . . .	378
<b>E. Tabellen mathematischer Funktionen</b>	<b>383</b>
E.1. Reelle Punkt-Funktionen . . . . .	383
E.2. Reelle Punkt-Funktionen der Mathematischen Physik . . . . .	385
E.3. Reelle Intervall-Funktionen . . . . .	386
E.4. Reelle Intervall-Funktionen der Mathematischen Physik . . . . .	388
E.5. Reelle Funktionen zur Automatischen Differentiation . . . . .	389
E.6. Reelle Funktionen zur Taylorarithmetik . . . . .	391
E.7. Komplexe Punkt-Funktionen . . . . .	393
E.8. Komplexe Intervall-Funktionen . . . . .	395
E.9. Komplexe Funktionen zur Taylorarithmetik . . . . .	397
<b>F. Laufzeitvergleiche</b>	<b>399</b>
<b>Literaturverzeichnis</b>	<b>401</b>
<b>Stichwortverzeichnis</b>	<b>405</b>

# 1. Einleitung

C-XSC [34, 36, 37, 65] ist eine auf **C++**basierende Programmierumgebung, mit deren selbst-verifizierenden numerischen Algorithmen die exakten mathematischen Lösungen von linearen oder nichtlinearen Gleichungssystemen, Integralen, Ableitungen von praktisch beliebig komplizierten Ausdrücken, Differential- und Integralgleichungen, oder Nullstellen von differenzierbaren Funktionen, ... durch Intervalle oder Intervallvektoren Funktionsschläuche, etc. nahezu optimal eingeschlossen werden können. In C-XSC sind optimale/semimorphe (siehe [46]) reelle und komplexe Punkt- und Intervall-Arithmetiken auch für Vektoren und Matrizen über Gleitkommazahlen implementiert.

Alle Rechnungen basieren auf dem IEEE-double-Standard, so dass eine Mantisse von ca. 16 Dezimalstellen und ein Zehnerexponentenbereich von  $-324$  bis  $+308$  zur Verfügung stehen. Rechnungen in höherer Präzision können in C-XSC 2.5.0 mit Hilfe der beiden Datentypen `l_interval` und `lx_interval` durchgeführt werden, die beide auf einem Staggered-Correction-Format (Staggered-Format) basieren, so dass jeweils eine maximale Präzision von etwa  $324 + 308 = 632$  Dezimalstellen benutzt werden kann.

Im Gegensatz zum `l_interval`-Typ, der den obigen Exponentenbereich von  $-324$  bis  $+308$  besitzt, steht mit dem `lx_interval`-Typ der sehr große Exponentenbereich von  $-2711437152599603$  bis  $+2711437152599603$  zur Verfügung, so dass hier Überlauf- oder Unterlaufprobleme praktisch keine Rolle mehr spielen. Auch die Auswertung der komplexwertigen Elementarfunktionen ist in beiden Staggered-Formaten möglich.

Der Vorteil des Staggered-Formats besteht nun darin, dass unabhängig von der gewählten Präzision die vier Grundoperationen  $\{+, -, *, /\}$  hochgenau und beliebig lange Skalarprodukte mit Hilfe des langen Akkumulators sogar maximalgenau ausgewertet werden können [3, 46]. Da dieser Akkumulator derzeit hardwaremäßig nicht unterstützt wird, muss er softwaremäßig simuliert werden, was zu langen Laufzeiten führt. Ein weiterer Nachteil der Staggered-Arithmetik besteht darin, dass bei zu breiten Eingangsintervallen die Ergebnisse nur noch in der Präzision des *double*-Formats, d.h. also mit nur etwa 16 Dezimalstellen, berechnet werden können.

Die geschilderten Nachteile des Staggered-Formats lassen sich jedoch vermeiden, wenn man die in Frankreich ab 2002 entwickelten MPFR- [21] und MPFI-Bibliotheken [60] benutzt, die von Laurent Fousse, Guillaume Hanrot, Vincent Lefèvre, Patrick Pélissier, Paul Zimmermann bzw. Nathalie Revol, Fabrice Rouillier, Sylvain Chevillard, Hong Diep NGUYEN und Christoph Lauter entwickelt wurden. Die MPFR-Bibliothek bietet eine Langzahl-Punktarithmetik, wobei die Grundoperationen und die zahlreichen Elementarfunktionen mit den folgenden Rundungsoptionen maximalgenau ausgewertet werden können:

- `MPFR_RNDN`      Rundung zur nächsten Rasterzahl des Langzahlformats
- `MPFR_RNDD`      Abrundung zur nächstkleineren Rasterzahl des Langzahlformats
- `MPFR_RNDU`      Aufrundung zur nächstgrößeren Rasterzahl des Langzahlformats
- `MPFR_RNDZ`      Rundung zur Null im gewählten Langzahlformat
- `MPFR_RNDA`      Rundung weg von der Null im gewählten Langzahlformat

Mit Hilfe der beiden Rundungsoptionen `MPFR_RNDD` und `MPFR_RNDU` können damit z.B. optimale Einschließungen exakter Funktionswerte der implementierten MPFR-Elementarfunktionen sehr effektiv berechnet werden.

In der MPFI-Bibliothek sind neben den Intervall-Grundoperationen  $\{\oplus, \ominus, \odot, \oslash\}$  fast alle Elementarfunktionen implementiert, mit denen zu beliebig breiten Eingangsintervallen maximalgenaue Einschließungen ihrer Funktionswerte berechnet werden können. Dabei werden Ober- und Unterschranken dieser Einschließungen in der vorgegebenen Präzision maximalgenau berechnet. Die aktuelle MPFI-Version ist 1.5. Um sie nutzen zu können, müssen die GMP-Bibliothek (Version 4.1 oder höher) und die MPFR-Bibliothek (Version 3.1.0 oder höher) vorhanden sein. Da die Bibliothek auf der MPFR-Bibliothek basiert, profitiert MPFI von den oben beschriebenen korrekten Rundungen der MPFR-Bibliothek. Die Bibliothek befolgt den IEEE-754-Standard der Gleitkomma-Arithmetik.

Da die MPFR- und MPFI-Bibliotheken auf einer Integer-Arithmetik basieren, können die entsprechenden Hardware-Ressourcen direkt angesprochen werden, was im Vergleich zur Staggered-Arithmetik [42, 49, 14, 15] die Laufzeit erheblich reduziert.

Unter C-XSC lassen sich die beiden Bibliotheken mit Hilfe zweier Interfaces ( $\mathbf{C}^{++}$ -Wrapper-Klassen) benutzen, die beide von Hans-Stephan Brand im Jahre 2010 im Rahmen einer Bachelor-Arbeit [16] entworfen und prototypmäßig realisiert wurden. Die MPFR- und MPFI-Bibliotheken kommen damit unter C-XSC mit den dort definierten Sprachelementen sehr einfach zur Anwendung. Insbesondere die MPFR-Bibliothek liefert zusätzliche Funktionen, wie z.B.  $\Gamma(x)$ ,  $\text{digamma}(x)$ ,  $\text{erf}(x)$ ,  $\text{erfc}(x)$ ,  $\zeta(x)$ ,  $\zeta(n)$ ,  $J_n(x)$ ,  $Y_n(x)$ , wobei  $J_n(x)$  und  $Y_n(x)$  die Besselfunktionen der ersten und zweiten Art sind, mit  $n = 0, 1, 2, \dots$ . Die streng monotonen Funktionen  $\text{erf}(x)$ ,  $\text{erfc}(x)$ ,  $\text{digamma}(x)$  und  $\Gamma'(x)$  lassen sich mit den beschriebenen Rundungsoptionen in C-XSC zusätzlich auch als Intervallfunktionen sehr einfach implementieren. Für komplexe Punkt- und Intervall-Rechnungen werden für C-XSC entsprechende  $\mathbf{C}^{++}$ -Wrapper-Klassen entwickelt.

Die  $\mathbf{C}^{++}$ -Wrapper-Klassen werden so implementiert, dass sie nicht zusammen mit den in C-XSC schon integrierten Staggered-Formaten genutzt werden können. Für Vergleichsrechnungen mag dies ein gewisser Nachteil sein, dafür bleibt der Programmieraufwand aber überschaubar.

## 2. Installation

Die Installation der MPFR- und MPFI-Bibliotheken ist nur unter einem Linux/Unix-System möglich. Dazu müssen unter **OpenSUSE 11.4** zunächst mit Yast2 u.a. installiert sein:

- `gcc44-c 4.4.1-20090817-2.3.4` Der GNU C++-Compiler
- `libgmp3 4.3.1-2.2` Shared library for the GNU MP Library
- `libgmpxx4 4.3.1-2.2` C++bindings for the GNU MP Library
- `make 3.81-130.2` GNU make
- `texlive 2008-13.18.1` und `texinfo 4.13a-3.2`

Weitere Informationen zur GNU MP Library findet man unter

<http://gmplib.org/>

### 2.1. Installation der MPFR- und MPFI-Bibliotheken

Die Installation der aktuellen GMP-, MPFR- und MPFI-Bibliotheken erfolgt unter **OpenSUSE 11.4** durch folgende Schritte:

1. Die Dateien `mpfr-3.1.1.tar.gz`, `mpfi-1.5.1.tar.gz`, `gmp-5.0.5.tar.xz` aus dem Netz laden und ins Heimatverzeichnis kopieren. Die letzte Datei findet man zum Download unter <http://gmplib.org/>
2. Mit Yast2 `gmp-devel` deinstallieren und `m4` neu installieren. Auf meinem Rechner waren neben `make 3.82-140.1` und `gcc 4.5.1` zusätzlich noch installiert:

```
libgmp10 5.0.1-4.5-x86_64    libgmp10-32bit 5.0.1-4.5-x86_64
libgmpxx4 5.0.1-4.5-x86_64
```

Ob die drei letzten Dateien auch wirklich gebraucht werden, wurde nicht getestet.

3. Die GMP-, MPFR- und MPFI-Bibliotheken deinstallieren und anschließend die entsprechenden Verzeichnisse löschen. Am Beispiel der MPFR-Bibliothek sind dazu folgende Schritte nötig
  - `cd ~/mpfr-3.1.0`
  - `make clean`
  - `gmake uninstall` (als root!)
  - Das Verzeichnis `~/mpfr-3.1.0` löschen

Die GMP- und MPFI-Bibliotheken sind analog zu deinstallieren und die Ordner entsprechend zu löschen.

4. Die aktuelle gmp-Version neu installieren mit:

- Wechseln ins Heimatverzeichnis, wo sich `gmp-5.0.5.tar.xz` befindet.
- `gmp-5.0.5.tar.xz` entpacken, wobei das Unterverzeichnis `~/gmp-5.0.5` entsteht.
- `cd ~/gmp-5.0.5`
- `./configure`
- `make`
- `make check` **dringend** empfohlen!
- `make install` (als root)

5. Die neueste Version MPFR-3.1.1 installieren mit:

- Wechseln in Heimatverzeichnis, wo sich `mpfr-3.1.1.tar.gz` befindet.
- `mpfr-3.1.1.tar.gz` entpacken, wobei das Verzeichnis `~/mpfr-3.1.1` entsteht.
- `cd ~/mpfr-3.1.1`
- `./configure`
- `make`
- `make check` **dringend** empfohlen!
- `make install` (als root)

Die Installation der neuesten Version MPFI-1.5.1 erfolgt ganz analog!

#### Anmerkungen:

1. Nach der beschriebenen Neuinstallation von `mpfr-3.1.1` konnte die Übersetzung eines Programms problemlos und ohne Warnungen durchgeführt werden. Beim Programmstart erfolgte jedoch die Fehlermeldung

```
error while loading shared libraries: libmpfi.so.0:  
cannot open shared object file: No such file or directory
```

Diese Fehlermeldung wurde jedoch beseitigt durch den root-Programmaufruf

```
ldconfig <enter>
```

Das Programm `ldconfig` aktualisiert die Links zu allen Bibliotheken und muss nach der **manuellen** Installation von Bibliotheken als root ausgeführt werden.

## 3. MpfrClass-Interface zur Anbindung der MPFR-Bibliothek an C-XSC

### 3.1. Grundlegendes

Das MpfrClass-Interface ist eine in `mpfrclass.hpp` und `mpfrclass.cpp` implementierte C++-Wrapper-Klasse `MpfrClass` für die C-Bibliothek MPFR, deren C-Funktionen über die implementierten Operatoren und Funktionen aufgerufen werden. Alle Funktionen mit der möglichen Übergabe eines `RoundingModes` oder eines `PrecisionTypes` besitzen als Standard die Werte von `CurrRndMode` bzw. `CurrPrecision`, die beide beliebig gesetzt werden können. Dies gilt auch für fast alle Konstruktoren.

#### 3.1.1. Allgemein

Um in C-XSC das Interface verwenden zu können, muss der Header `mpfrclass.hpp` eingebunden werden. In ihm sind die benötigten Header-Dateien der MPFR-Bibliothek enthalten. Die `MpfrClass`-Klasse liegt im Namensraum "MPFR".

#### 3.1.2. Aufbau

Die Klasse besteht intern aus einer "mpfr\_t"-Variablen. Diese dient zum Speichern des Wertes. Zusätzlich gibt es static Elemente, um den Standard-Rundungsmodus, die Standard-Precision und die aktuelle Basis zu speichern.

#### 3.1.3. Präzision

Die Präzision gibt die Anzahl der binären Mantissenstellen einer `MpfrClass`-Variablen `x` an. Die Präzision von `x` kann mit `x.SetPrecision(prec)` auf den Wert `prec ≥ 2` gesetzt werden. Die maximale Präzision ist durch `MPFR_PREC_MAX = 9223372036854775807` vorgegeben. Man sollte diesen großen Wert jedoch auch nicht annähernd verwenden, da Laufzeit und Speicherkapazität überfordert wären. Außerdem kann bei internen Rechnungen, z.B. zur Garantie optimaler Rundungen, die Current-Präzision etwa verdoppelt werden, so dass bei der Präzision nach oben genug Spielraum vorhanden sein sollte. Die Current-Precision ist die Präzision, mit der z.B. alle arithmetischen Grundoperationen durchgeführt werden. Sie kann global mit `SetCurrPrecision(prec)` gesetzt werden. Wenn dies nicht geschieht, so wird mit der Default-Precision von 53 Bits gerechnet. Unabhängig davon kann die Präzision für jede `MpfrClass`-Variable `x` auch einzeln festgelegt werden.

#### 3.1.4. Variablentyp `PrecisionType`

Mithilfe des Variablentyps `PrecisionType` (Name der Variablen meist `prec`) kann der Präzisionswert einer `MpfrClass`-Variablen eingestellt werden. Der Variablentyp ist ein typedef für eine "mp\_prec\_t"-Variable.

### 3.1.5. Variablentyp `RoundingMode`

Mit dem Variablentyp “`RoundingMode`” (Name der Variablen meist `rnd`) wird der gewünschte Rundungsmodus eingestellt. Der Variablentyp ist ein typedef für eine `mpfr_rnd_t`-Variable.

- `RoundNearest` – Der Wert wird zur nächsten Rasterzahl gerundet (0)
- `RoundUp` – Der Wert wird aufgerundet (2)
- `RoundDown` – Der Wert wird abgerundet (3)
- `RoundToZero` – Der Wert wird in Richtung Null gerundet (1)
- `RoundFromZero` – Rundung weg von der Null (4)

Der Rundungsmodus kann global mit `SetCurrRndMode(rnd)` gesetzt werden, sonst wird mit dem Standard-Rundungsmodus `RoundNearest` gerundet. Soll abweichend von diesem globalen Rundungsmodus, z.B. bei einer Multiplikation, anders gerundet werden, so kann dies mit den entsprechenden C-Funktionen und ihrem eigenen Rundungsparameter `rnd` erfolgen, siehe etwa das Beispiel auf Seite 30.



### 3.1.6. Zahlenformat

Im Zusammenhang mit der Current-Präzision, die mit `SetCurrPrecision(prec)` auf  $\text{prec} \geq 2$  gesetzt werden kann, soll das dazugehörige MPFR-Format im Vergleich zum bekannten IEEE-Format genauer untersucht werden.

Wir betrachten zunächst das **IEEE-Format**:

- Kleinste positive (denormalisierte) Zahl:  $\text{minreal} := 2^{-1074}$ .
- $\text{succ}(\text{minreal}) = 2 \cdot \text{minreal} = 2^{-1073}$ .
- Kleinste positive normalisierte Zahl:  $\text{MinReal} := 2^{-1022} = \text{minreal} \cdot 2^{52}$ .
- Größte positive normalisierte Zahl:  $\text{MaxReal} < 2^{+1024}$ .
- $\text{pred}(\text{minreal}) = 0$ ,  $\text{succ}(\text{minreal}) = 2 \cdot \text{minreal} = 2^{-1073}$ .
- $\text{pred}(\text{MinReal}) = \text{MinReal} - \text{minreal} = (2^{52} - 1) \cdot \text{minreal}$ .
- $\text{succ}(\text{MinReal}) = \text{MinReal} + \text{minreal} = (2^{52} + 1) \cdot \text{minreal}$ .
- $\text{expo}(\text{MinReal}) = \text{expo}(\text{succ}(\text{MinReal})) = -1021$ .
- $\text{expo}(\text{minreal}) = -1073 < \text{expo}(\text{succ}(\text{minreal})) = -1072$ .

Wir betrachten jetzt das **MPFR-Format** mit der Current-Präzision  $\text{prec} \geq 2$ :

- Kleinste positive Zahl:  $\text{minfloat}(\text{prec}) := 2^{-1073741824}$  ist  $\text{prec}$ -unabhängig!
- $\text{mant}(\text{minfloat}(\text{prec})) = 0.5$ ,  $\text{expo}(\text{minfloat}(\text{prec})) = -1073741823$ .
- $\text{pred}(\text{minfloat}()) = 0$ ,  $\text{succ}(\text{minfloat}()) < 2 \cdot \text{minfloat}()$ .
- $\text{expo}(\text{succ}(\text{minfloat}())) = \text{expo}(\text{minfloat}()) = -1073741823$ .
- Ausgehend von  $\text{minfloat}() = 0.5 \cdot 2^{-1073741823}$  ist die kleinste positive Maschinenzahl  $r$  mit dem nächst-größeren Exponenten  $-1073741822$  gegeben durch:  $r = 2 \cdot \text{minfloat}()$ .
- Größte positive normalisierte Zahl:  $2^{+1073741822} < \text{MaxFloat}(\text{prec}) < 2^{+1073741823}$ .  
 $\text{expo}(\text{MaxFloat}(\text{prec})) = +1073741823$ .  $\text{mant}(\text{MaxFloat}(\text{prec})) < 1$ .

Vergleicht man mit den Werten des IEEE-Formats, so erkennt man, dass das MPFR-Format keinen denormalisierten Zahlenbereich besitzt. Für **jeden** Exponenten  $\text{ex} = \text{expo}(x)$  einer Maschinenzahl  $x$ , die nicht unendlich und kein NaN ist, stehen für die Mantisse von  $x$  genau  $\text{prec} \geq 2$  Bits zur Verfügung, so dass es im Gegensatz zum IEEE-Format genau  $2^{\text{prec}}$  Maschinenzahlen für jeden Exponenten gibt.

Bei möglichen Fehlerbetrachtungen muss also nicht mehr zwischen einem normalisierten und denormalisierten Bereich unterschieden werden. Man muss also nur noch darauf achten, dass Zwischenergebnisse möglichst nicht in den Unterlaufbereich fallen.

Für den Vorgänger einer Maschinenzahl  $x = m \cdot 2^{\text{ex}} \geq 0$  gilt mit  $\text{ex} = \text{expo}(x)$  und  $m = \text{mant}(x)$ :

$$(3.1) \quad \text{pred}(x) = \begin{cases} -\text{minfloat}(), & \text{falls } x = 0, \\ 0, & \text{falls } x = \text{minfloat}(), \\ 2^{\text{ex}-1} \cdot (1 - 2^{-\text{prec}}), & \text{falls } m = 0.5, x > \text{minfloat}(), \\ 2^{\text{ex}} \cdot (m - 2^{-\text{prec}}), & \text{falls } 0.5 < m < 1, x > \text{minfloat}(). \end{cases}$$

Für den Nachfolger einer Maschinenzahl  $x = m \cdot 2^{\text{ex}} \geq 0$  gilt die einfachere Darstellung:

$$(3.2) \quad \text{succ}(x) = \begin{cases} \text{minfloat}(), & \text{falls } x = 0, \\ 2^{\text{ex}} \cdot (m + 2^{-\text{prec}}), & \text{falls } 0.5 \leq m < 1. \end{cases}$$

Für  $x \leq 0$  gilt:  $\text{succ}(x) = -\text{pred}(-x)$  und  $\text{pred}(x) = -\text{succ}(-x)$ .

## 3.2. Konstruktoren / Destruktor

### 3.2.1. Konstruktoren

```
MpfrClass ();
```

Der Default-Konstruktor legt ein neues Element mit der Current-Precision an.  
Der Aufruf `MpfrClass y;` initialisiert den Wert: `y = NaN;`

```
MpfrClass (const MpfrClass& op, RoundingMode rnd, PrecisionType prec);  
MpfrClass (const mpfr_t& op, RoundingMode rnd, PrecisionType prec);  
MpfrClass (int op, RoundingMode rnd, PrecisionType prec);  
MpfrClass (const double& op, RoundingMode rnd, PrecisionType prec);  
MpfrClass (const cxsc::real& op, RoundingMode rnd, PrecisionType prec);
```

Mögliche Konstruktor-Aufrufe sind:

1. `MpfrClass y(op);`
2. `MpfrClass y(op, RoundNearest);`
3. `MpfrClass y(op, RoundDown, 3);`

Zu 1. `op` wird mit dem Current-Rundungsmodus auf die Current-Precision gerundet.

Zu 2. `op` wird mit `RoundNearest` auf die Current-Precision gerundet.

Zu 3. `op` wird auf die (sehr kleine) Precision `prec = 3` abgerundet.

Der Aufruf `MpfrClass y(op, 3);` führt zu einer Fehlermeldung, da in der Parameterliste der Rundungsmodus fehlt. Die oberen fünf Konstruktoren erlauben also eine sehr flexible Initialisierung von `MpfrClass`-Objekten.

Mit den Deklarationen

```
int k; double dbl; real r;
```

liefern die folgenden Konstruktor-Aufrufe

- `MpfrClass y(k, RoundNearest, 32); MpfrClass y(dbl, RoundNearest, 53);`
- `MpfrClass y(r, RoundNearest, 53);`

`MpfrClass`-Objekte `y` mit den **ungerundeten** Werten von `k`, `dbl`, `r`.

Der folgende Aufruf: `MpfrClass y(x, RoundNearest, x.GetPrecision());` mit `x` vom Typ `MpfrClass` liefert ein `MpfrClass`-Objekt `y`, mit `y = x`, (Copy-Konstruktor).

```
MpfrClass (const std::string& s, RoundingMode rnd, PrecisionType prec);
```

Der Aufruf `MpfrClass y(s);` rundet `s` mittels der voreingestellten Current-Precision und mit dem Current-Rundungsmodus in das Klassenobjekt `y`. Der Aufruf `MpfrClass y(s, RoundNearest, 140);` rundet `s` mit der Präzision von 140 Bits zur nächsten Rasterzahl dieses Formats.

**Achtung:** Keine Leerzeichen am String-Ende!

### 3.2.2. Destruktor

```
~MpfrClass ();
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

### 3.3. Eingabe / Ausgabe

```
friend std::ostream& operator << (std::ostream& os, const MpfrClass& x);
```

Ermöglicht die Ausgabe einer MpfrClass-Variablen `x` über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(x.GetPrecision()/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Der Wert der Variablen `x` wird dabei in der mit `SetBase(k)` voreingestellten Basis ausgegeben. Für die meist dezimale Ausgabe ist `k = 10` zu wählen.

```
std::ostream& operator << (std::ostream& os, mpfr_t& x);
```

Ermöglicht die Ausgabe einer Variablen `x` vom Typ `mpfr_t` über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(mpfr_get_prec(x)/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Der Wert `x` wird dabei in der mit `SetBase(k)` voreingestellten Basis ausgegeben. Für die meist dezimale Ausgabe ist `k = 10` zu wählen.

```
friend std::istream& operator >> (std::istream& is, MpfrClass& x);
```

Ermöglicht das Einlesen einer MpfrClass-Variablen über den Standardeingabestrom "cin". Die eingegebene Zahl ist auf keine Stellenanzahl begrenzt. Die Zahl muss in der eingestellten Basis eingegeben werden, sonst entsteht eine Fehlermeldung. Die Rundung erfolgt mit dem voreingestellten Current-Rundungsmodus, siehe das Programm auf Seite 21.

### 3.4. Rundungsmodi und Precision Handling

```
PrecisionType GetPrecision () const;
```

Diese Memberfunktion gibt die Präzision des aktuellen Objekts in Bits zurück. Als Beispiel entsprechen dabei 302 Bits  $302/\log_2(10) \approx 91$  Dezimalstellen.

```
void SetPrecision (PrecisionType prec);
```

Diese Memberfunktion setzt die Präzision des aktuellen Objektes auf `prec`. Sein Wert bleibt nicht erhalten.

```
void RoundPrecision (PrecisionType prec, RoundingMode rnd);
```

Diese Memberfunktion rundet das aktuelle Objekt auf die neue Precision `prec`, sein Wert bleibt dabei i.a. nicht erhalten. Sollte die Präzision des Objektes größer sein als `prec`, wird das Objekt mit Hilfe des eingestellten Rundungsmodus so gerundet, dass es in das Format der Präzision `prec` passt. Ist die Präzision kleiner als `prec`, werden die restlichen binären Stellen mit Nullen aufgefüllt.

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision auf `prec`. Wird die Current-Precision nicht gesetzt, so wird mit der Default-Precision von 53 Bits gerechnet.

```
static const int GetBase ();
```

Gibt die aktuelle Basis zurück, diese hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen.

```
static void SetBase (int b);
```

Setzt die aktuelle Basis auf `b`. Dies muss ein Wert zwischen 2 und 36 sein. Die Basis hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die internen Rechnungen erfolgen stets im Binärsystem!

```
static const RoundingMode GetCurrRndMode ();
```

Gibt den aktuellen Rundungsmodus zurück mit den Werten: (0, 1, 2, 3, 4).

```
static void SetCurrRndMode (RoundingMode rnd);
```

Setzt den Current-Rundungsmodus auf `rnd`. Für `rnd` sind fünf verschiedene Modi möglich:

- RoundNearest: Rundung zur nächsten Rasterzahl (0)
- RoundToZero: Rundung in Richtung Null (1)
- RoundFromZero: Rundung weg von der Null (4)
- RoundUp: Aufrunden (2)
- RoundDown: Abrunden (3)

Wird der Rundungsmodus mit `SetCurrRndMode` nicht gesetzt, so wird als Default-Rundungsmodus `RoundNearest` benutzt.

### 3.5. Anwendungsprogramm

Das folgende C-XSC Programm MPFR-01.cpp zeigt einen Konstruktor-Aufruf, den Eingabe- und Ausgabe-Mechanismus und das Rundungs- bzw. Precision-Handling:

```
1 // MPFR-01.cpp
2 #include <iostream>
3 #include <iomanip>
4 #include "mpfrclass.hpp"
5 using namespace MPFR;
6 using namespace std;
7 int main(void)
8 {
9     MpfrClass::SetCurrRndMode (RoundUp);
10    cout << "\nRoundingMode = " << MpfrClass::GetCurrRndMode() << endl;
11    MpfrClass::SetCurrPrecision (40);
12    cout << "Current-Precision = " << MpfrClass::GetCurrPrecision() << endl;
13
14    MpfrClass x(1.2345678, RoundDown, 20); // Konstruktor mit double-Argument
15    cout.precision(x.GetPrecision()/3.32192809); // Dezimale Stellenzahl
16    cout << "x = " << x << endl;
17    cout << "Precision of x = " << x.GetPrecision() << endl << endl;
18
19    x.SetPrecision(70); // Inhalt von x wird dabei geloescht
20    cout << "x = ? " ; cin >> x;
21    cout.precision(x.GetPrecision()/3.32192809);
22    cout << "x = " << x << endl;
23    cout << "Precision of x = " << x.GetPrecision() << endl << endl;
24    return 0;
25 }
```

Bei der interaktiven Eingabe (Zeile 22) von 1.2345678 erhält man die folgende Ausgabe:

```
RoundingMode = 2
Current-Precision = 40
x = 1.23457
Precision of x = 20

x = ? 1.2345678
x = 1.2345678000000000000001
Precision of x = 70
```

- Beim Konstruktoraufruf (14) wird 1.2345678 zunächst vom Compiler, **vermutlich** mit RoundToNearest, in einem internen *double*-Binärformat gespeichert, das dann vom Konstruktor durch Abrunden in das Objekt x mit der Precision von 20 Bit geschrieben wird, was etwa  $20/3.3219 \approx 6$  Dezimalziffern entspricht. Also **Vorsicht**:

Da man nicht genau weiß, wie der Compiler die dezimale Eingabe 1.2345678 intern speichert, weiß man auch nicht genau, welchen Wert das Objekt x durch den Konstruktor-Aufruf erhält. Diese Unsicherheit lässt sich vermeiden, wenn man in C-XSC z.B. eine *real*-Variable anlegt und diese mit cin und entsprechenden Rundungsmanipulatoren (z.B. RndNext) mit 1.2345678 initialisiert und dann diese *real*-Variable an den Konstruktor übergibt.

- Das *double*-Format mit 53 Mantissen-Bits ergibt eine maximale Precision von  $53/3.3219 \approx 16$  Dezimalstellen. Will man den Dezimalwert 1.2345678, z.B. aufgerundet, mit einer Precision von z.B. 70 Bit, d.h. etwa 21 Dezimalstellen, in ein Objekt x speichern, so kann dies nach Zeile 20 mit cin realisiert werden, wobei jedoch zum Aufrunden in Zeile 9 der Current-Rundungsmodus RoundUp zu setzen ist.
- In (15) und (21) wird mit cout.precision(...) das Ausgabeformat festgelegt.

## 3.6. Typ-Umwandlungen

Um ein flexibles Arbeiten zwischen C-XSC und der MPFR-Bibliothek zu ermöglichen, wurden möglichst viele Typ-Umwandlungs-Funktionen bereitgestellt:

### 3.6.1. real, double, ... → MPFR

```
MpfrClass real2Mpfr (const cxsc::real& op);  
MpfrClass double2Mpfr (const double& op);  
MpfrClass int2Mpfr (const int& op);  
MpfrClass mpfr_t2Mpfr (const mpfr_t& op);
```

Obige Funktionen liefern mit dem jeweiligen Eingabewert `op` einen Rückgabewert vom Typ `MpfrClass` in einer Präzision, die gewährleistet, dass der Rückgabewert **genau** dem Wert von `op` entspricht. Die Präzision des Rückgabewertes wird also i.a. nicht mit der Current-Precision übereinstimmen! Die obigen vier Funktionen kommen u.a. bei den Vergleichsoperatoren zur Anwendung.

### 3.6.2. MPFR → real, double, ...

Die folgenden drei Funktionen liefern mit einem Objekt `op` vom Typ `MpfrClass` jeweils einen i.a. **gerundeten** Rückgabewert vom Typ `real`, `double`, `long int`;

```
cxsc::real to_real (const MpfrClass& op, RoundingMode rnd);  
double to_double (const MpfrClass& op, RoundingMode rnd);  
long int to_int (const MpfrClass& op, RoundingMode rnd);
```

Mit z.B. `rnd = RoundUp` wird aufgerundet. Ohne Angabe eines Rundungsmodus wird jedoch nach dem voreingestellten Current-Rundungsmodus gerundet.

### 3.6.3. MPFR → mpfr\_t

```
const mpfr_t& getvalue(const MpfrClass& r)
```

Obige Funktion liefert von einem als `const` deklarierten Objekt `r` eine Referenz auf den Wert seines Attributs `mpfr_rep` vom Typ `mpfr_t`.

**Anwendung:** `const`-Parameter in einer Funktionen-Parameterliste, vgl. z.B. die Funktion `MpfiClass MpfrClass2Mpfi(const MPFR::MpfrClass& v)` in der Datei `mpficlass.cpp`.

### 3.6.4. MPFR ↔ mpfr\_t

```
mpfr_t& MpfrClass::GetValue();
```

Mithilfe der Memberfunktion `GetValue()`, die eine Referenz auf den Wert `mpfr_rep` vom Typ `mpfr_t` des aktuellen Objekts liefert, können sowohl `MpfrClass`-Objekte an eine Funktion übergeben als auch referenzierte Rückgabewerte vom Typ `mpfr_t` von einer Funktion übernommen werden. Mithilfe von `GetValue()` kann man daher Funktionen mit referenzierten `mpfr_t`-Parametern einfach aufrufen, vergleiche dazu das Programm MPFR-02 auf Seite 30.

### 3.6.5. mpfr\_t → MPFR

```
void SetValue(const mpfr_t& t);
```

Mithilfe der Memberfunktion `SetValue()` wird der `mpfr_rep`-Wert des aktuellen Objekts auf `t` gesetzt, wobei `mpfr_rep` die Präzision von `t` übernimmt, d.h. der Wert von `t` wird ohne Rundung übernommen.

### 3.6.6. string → MPFR

```
MpfrClass string2Mpfr(const std::string& op, Roundingmode rnd,  
                    PrecisionType prec);
```

Der Aufruf `string2Mpfr(op);` rundet den String `op` mittels der voreingestellten Current-Precision und mit dem Current-Rundungsmodus in ein Klassenobjekt vom Typ `MpfrClass`. Der Aufruf `string2Mpfr(op, RoundNearest, 140);` rundet `op` in ein Klassenobjekt vom Typ `MpfrClass` der Präzision 140 Bits. Gerundet wird dabei zur nächsten Rasterzahl dieses Formats. Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in eine binäre Zahl konvertiert werden kann, so dass daher Rundungen i.a. nicht zu vermeiden sind. Eine weitere Möglichkeit, einen String in ein `MpfrClass`-Objekt zu verwandeln, besteht in einem entsprechenden Konstruktor-Aufruf, vgl. Seite 18. **Achtung:** Keine Leerzeichen am String-Ende!

### 3.6.7. MPFR → string

```
std::string to_string (const MpfrClass& x, RoundingMode rnd,  
                    PrecisionType prec);  
std::string to_string (const mpfr_t& x,    RoundingMode rnd,  
                    PrecisionType prec);
```

`x` wird mittels `rnd` in einen String `s` mit `prec` Dezimalstellen gerundet, wenn Base gleich 10 ist. Wählt man `prec` hinreichend groß, so stellt der String den Wert von `x` **exakt** dar, weil eine binäre Zahl stets exakt in eine Dezimalzahl umgewandelt werden kann.

Wird `prec` nicht angegeben, so wird `x` mittels `rnd` in einen String gerundet, der bei Base=10 so viele Dezimalstellen besitzt, wie es der Präzision von `x` entspricht. Ist die Präzision von `x` z.B. 302, so wird ein String von  $302/\log_2(10) \approx 91$  Dezimalstellen generiert. Wird neben `prec` auch `rnd` nicht angegeben, so wird mittels des Current-Rundungsmodus in den String mit gleicher Dezimalstellenzahl gerundet.

### 3.6.8. MPFR → MPFR

Beachten Sie bitte, dass bei einer Wertzuweisung an eine `MpfrClass`-Variable mit Hilfe des Operators `=` der linke Operand stets auf die Current-Precision gesetzt wird und dass der rechte `MpfrClass`-Operand dabei **stets** bez. des Current-Rundungsmodus in den linken Operanden gerundet wird, vgl. dazu auch Seite 25. Will man jedoch abweichend von dieser Rundung einen anderen Rundungsmodus benutzen, so kann dies mit folgender Funktion ohne Rückgabewert realisiert werden:

```
void set_Mpfr (MpfrClass& op, const MpfrClass& op1, RoundingMode rnd,  
             PrecisionType prec);
```

Folgende Funktionsaufrufe sind möglich:

1. `set_Mpfr (op, op1, RoundUp, prec);`  
`op` erhält die Präzision `prec` und den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls aufgerundet wird. Setzt man `prec` gleich der Präzision von `op1`, so erhält `op` den **exakten** Wert von `op1`, und zwar unabhängig vom gewählten Rundungsmodus `rnd`.
2. `set_Mpfr (op, op1, RoundDown);`  
`op` wird auf die Current-Precision gesetzt und erhält den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls abgerundet wird.

3. `set_Mpfr (op, op1);`  
`op` wird auf die Current-Precision gesetzt und erhält den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls bez. des Current-Rundungsmodus gerundet wird.

### 3.6.9. Verschiedenes

```
int sign(const MpfrClass& x)
```

Zurückgegeben wird das Vorzeichen von `x`.

```
sign(NaN) = 0;   sign(+inf) = +1;   sign(-inf) = -1;
```

```
long int expo (const MpfrClass& x);
```

Zurückgegeben wird der Exponent `e` von `x`. Ist `m` die Mantisse einer normalen Gleitpunktzahl  $x \neq 0$ , so gilt:  $x = m \cdot 2^e$ ,  $|m| \in [0.5, 1)$ .

```
expo(0)      = -9223372036854775807; expo(NaN)  = -9223372036854775806;
```

```
expo(+inf) = -9223372036854775805; expo(-inf) = -9223372036854775805;
```

```
expo(MaxFloat()) = 1073741823;   expo(minfloat()) = -1073741823;
```

**Beachten** Sie: Die Mantisse `m` wird mithilfe von `mant(x)` bestimmt, vgl. Seite 33. Im Gegensatz zu C-XSC ist bei festem `prec` die Präzision der Mantisse `m` jedoch **unabhängig** vom Wert des Zweier-Exponenten `expo(x)`, wobei für  $x \neq 0$  gilt

$$-1073741823 \leq \text{expo}(x) \leq +1073741823.$$

Danach erhält man z.B. `mant(succ(0.5)) = mant(succ(minfloat(0)))`, wobei `minfloat()` die kleinste positive Zahl vom Typ `MpfrClass` ist, vgl. Seite 34.

**Beachten** Sie auch, dass bei einer Multiplikation von `x` mit  $2^k$ ,  $k \in \mathbb{Z}$ , z.B. mit Hilfe von `times2pown(x,k,rnd)`, der Rückgabewert `x` das **exakte** Produkt  $x \cdot 2^k$  liefert, solange kein Über- oder Unterlauf eintritt.



### 3.7. Zuweisungs-Operatoren

Bei den folgenden fünf Zuweisungs-Operationen erhält der linke Operand als Präzision stets die aktuelle Current-Precision, und der rechte Operand `op` wird nach dem aktuellen Current-Rundungsmodus in den linken Operanden gerundet.

```
MpfrClass& operator = (const MpfrClass& op);  
MpfrClass& operator = (const mpfr_t& op);  
MpfrClass& operator = (const cxsc::real& op);  
MpfrClass& operator = (const double& op);  
MpfrClass& operator = (const int& op);
```

Ist z.B. `op` vom Typ `real` oder `double` und ist die Current-Precision kleiner als 53, so wird `op` i.a. in den linken Operanden gerundet. Nur wenn die Current-Precision größer oder gleich 53 ist, wird der Wert des linken Operanden gleich dem Wert des rechten Operanden. Für andere Typen des rechten Operanden gelten ganz entsprechende Aussagen.

Ist z.B. `op` vom Typ `MpfrClass` und ist seine Präzision größer als die Current-Precision, so wird `op` in den linken Operanden bez. des Current-Rundungsmodus gerundet, d.h. die Werte des linken und rechten Operanden werden dann i.a. **verschieden** sein!

### 3.8. Abfragen

Bei allen folgenden Abfragefunktionen braucht die Präzision von `x` nicht mit der Current-Precision übereinzustimmen.

```
bool isNaN (const MpfrClass& x);  
bool isInf (const MpfrClass& x);  
bool isNumber(const MpfrClass& x);
```

`isNaN` und `isInf` überprüfen, ob `x` gleich NaN bzw.  $\pm\text{Inf}$  ist. `isNumber` überprüft, ob `x` eine normale `MpfrClass` Zahl ungleich NaN und ungleich  $\pm\text{Inf}$  ist.

```
bool isZero (const MpfrClass& x);  
bool isNeg (const MpfrClass& x);  
bool isPos (const MpfrClass& x);  
bool isInteger(const MpfrClass& x);  
bool isEven (const MpfrClass& x);  
bool isOdd (const MpfrClass& x);
```

Die oberen sechs Funktionen sind selbsterklärend, wobei die drei letzten wirklich praktische Bedeutung haben!

Mit `isInteger(x)` wird lediglich geprüft, ob  $x \in \mathbb{Z}$  erfüllt ist, wobei `x` jedoch nicht von Typ `int` oder `long int` sein muss.

## 3.9. Vergleiche

Mit den folgenden Vergleichsfunktionen werden die üblichen Vergleichsoperatoren implementiert, wobei wenigstens ein Operand vom Typ `MpfrClass` sein muss.

### 3.9.1. Vergleichsfunktionen

```
int compare_equal (const MpfrClass& x, const MpfrClass& y);
int compare_less (const MpfrClass& x, const MpfrClass& y);
int compare_lessequal (const MpfrClass& x, const MpfrClass& y);
int compare_greater (const MpfrClass& x, const MpfrClass& y);
int compare_greaterequal (const MpfrClass& x, const MpfrClass& y);
```

Die obigen Funktionen liefern einen Wert ungleich Null, wenn jeweils gilt:

$x = y$ ,  $x < y$ ,  $x \leq y$ ,  $x > y$ ,  $x \geq y$  und den Wert Null sonst. Ist  $x$  oder  $y$  NaN, so wird Null zurückgegeben. Die Präzisionen von  $x$  und  $y$  können verschieden sein und müssen mit der Current-Precision nicht übereinstimmen.

### 3.9.2. Vergleichsoperatoren =, ≠, >, ≥, <, ≤

```
bool operator == (const MpfrClass& op1, const MpfrClass& op2);
bool operator == (const MpfrClass& op1, const double& op2);
bool operator == (const MpfrClass& op1, const int op2);
bool operator == (const MpfrClass& op1, const cxsc::real& op2);
bool operator == (const double& op2, const MpfrClass& op1);
bool operator == (const int op2, const MpfrClass& op1);
bool operator == (const cxsc::real& op2, const MpfrClass& op1);
```

```
bool operator != (const MpfrClass& op1, const MpfrClass& op2);
bool operator != (const MpfrClass& op1, const double& op2);
bool operator != (const MpfrClass& op1, const int op2);
bool operator != (const MpfrClass& op1, const cxsc::real& op2);
bool operator != (const double& op2, const MpfrClass& op1);
bool operator != (const int op2, const MpfrClass& op1);
bool operator != (const cxsc::real& op2, const MpfrClass& op1);
```

```
bool operator < (const MpfrClass& op1, const MpfrClass& op2);
bool operator < (const MpfrClass& op1, const double& op2);
bool operator < (const MpfrClass& op1, const int op2);
bool operator < (const MpfrClass& op1, const cxsc::real& op2);
bool operator < (const double& op2, const MpfrClass& op1);
bool operator < (const int op2, const MpfrClass& op1);
bool operator < (const cxsc::real& op2, const MpfrClass& op1);
```

```
bool operator <= (const MpfrClass& op1, const MpfrClass& op2);
bool operator <= (const MpfrClass& op1, const double& op2);
bool operator <= (const MpfrClass& op1, const int op2);
bool operator <= (const MpfrClass& op1, const cxsc::real& op2);
bool operator <= (const double& op2, const MpfrClass& op1);
bool operator <= (const int op2, const MpfrClass& op1);
bool operator <= (const cxsc::real& op2, const MpfrClass& op1);
```

```
bool operator > (const MpfrClass& op1, const MpfrClass& op2);
bool operator > (const MpfrClass& op1, const double& op2);
bool operator > (const MpfrClass& op1, const int op2);
bool operator > (const MpfrClass& op1, const cxsc::real& op2);
bool operator > (const double& op2, const MpfrClass& op1);
bool operator > (const int op2, const MpfrClass& op1);
bool operator > (const cxsc::real& op2, const MpfrClass& op1);
```

```
bool operator >= (const MpfrClass& op1, const MpfrClass& op2);
bool operator >= (const MpfrClass& op1, const double& op2);
bool operator >= (const MpfrClass& op1, const int op2);
bool operator >= (const MpfrClass& op1, const cxsc::real& op2);
bool operator >= (const double& op2, const MpfrClass& op1);
bool operator >= (const int op2, const MpfrClass& op1);
bool operator >= (const cxsc::real& op2, const MpfrClass& op1);
```

Sind die Operanden `op1` und `op2` beide vom Typ `MpfrClass`, so müssen ihre Präzisionen nicht übereinstimmen und insbesondere auch nicht mit der Current-Precision.

## 3.10. Arithmetische Operatoren

Für alle arithmetischen Operatoren gilt:

Das exakte Ergebnis einer arithmetischen Operation wird unabhängig von der Präzision der Operanden mit dem voreingestellten Current-Rundungsmodus optimal gerundet. Die Ergebnis-Präzision ist dabei stets gleich der voreingestellten Current-Precision.

Die Operatoren  $\odot=$ , mit  $\odot \in \{+, -, \cdot, /\}$ , bedeuten  $u\odot=v \iff u = u \odot v$ . Dabei wird  $u\odot v$  mit dem Current-Rundungsmodus in die Current-Precision gerundet und in  $u$  gespeichert, wobei  $u$  als Präzision die Current-Precision erhält.

### 3.10.1. Addition

```
MpfrClass operator + (const MpfrClass& op1, const MpfrClass& op2);
MpfrClass operator + (const MpfrClass& op1, const mpfr_t& op2);
MpfrClass operator + (const MpfrClass& op1, const double& op2);
MpfrClass operator + (const MpfrClass& op1, const cxsc::real& op2);
MpfrClass operator + (const MpfrClass& op1, const int op2);
```

```
MpfrClass operator + (const mpfr_t& op1, const MpfrClass& op2);
MpfrClass operator + (const double& op1, const MpfrClass& op2);
MpfrClass operator + (const cxsc::real& op1, const MpfrClass& op2);
MpfrClass operator + (const int op1, const MpfrClass& op2);
```

```
MpfrClass& operator += (MpfrClass& op1, const MpfrClass& op2);
MpfrClass& operator += (MpfrClass& op1, const mpfr_t& op2);
MpfrClass& operator += (MpfrClass& op1, const double& op2);
MpfrClass& operator += (MpfrClass& op1, const cxsc::real& op2);
MpfrClass& operator += (MpfrClass& op1, const int op2);
```

### 3.10.2. Subtraktion

Beachten Sie den Hinweis auf Seite 28

```
MpfrClass operator - (const MpfrClass& op1, const MpfrClass& op2);
MpfrClass operator - (const MpfrClass& op1, const mpfr_t& op2);
MpfrClass operator - (const MpfrClass& op1, const double& op2);
MpfrClass operator - (const MpfrClass& op1, const cxsc::real& op2);
MpfrClass operator - (const MpfrClass& op1, const int op2);
```

```
MpfrClass operator - (const mpfr_t& op1, const MpfrClass& op2);
MpfrClass operator - (const double& op1, const MpfrClass& op2);
MpfrClass operator - (const cxsc::real& op1, const MpfrClass& op2);
MpfrClass operator - (const int op1, const MpfrClass& op2);
```

```
MpfrClass& operator -= (MpfrClass& op1, const MpfrClass& op2);
MpfrClass& operator -= (MpfrClass& op1, const mpfr_t& op2);
MpfrClass& operator -= (MpfrClass& op1, const double& op2);
MpfrClass& operator -= (MpfrClass& op1, const cxsc::real& op2);
MpfrClass& operator -= (MpfrClass& op1, const int op2);
```

### 3.10.3. Multiplikation

Beachten Sie den Hinweis auf Seite 28

```
MpfrClass operator * (const MpfrClass& op1, const MpfrClass& op2);
MpfrClass operator * (const MpfrClass& op1, const mpfr_t& op2);
MpfrClass operator * (const MpfrClass& op1, const double& op2);
MpfrClass operator * (const MpfrClass& op1, const cxsc::real& op2);
MpfrClass operator * (const MpfrClass& op1, const int op2);

MpfrClass operator * (const mpfr_t& op1, const MpfrClass& op2);
MpfrClass operator * (const double& op1, const MpfrClass& op2);
MpfrClass operator * (const cxsc::real& op1, const MpfrClass& op2);
MpfrClass operator * (const int op1, const MpfrClass& op2);

MpfrClass& operator *= (MpfrClass& op1, const MpfrClass& op2);
MpfrClass& operator *= (MpfrClass& op1, const mpfr_t& op2);
MpfrClass& operator *= (MpfrClass& op1, const double& op2);
MpfrClass& operator *= (MpfrClass& op1, const cxsc::real& op2);
MpfrClass& operator *= (MpfrClass& op1, const int op2);
```

### 3.10.4. Division

Beachten Sie den Hinweis auf Seite 28

```
MpfrClass operator / (const MpfrClass& op1, const MpfrClass& op2);
MpfrClass operator / (const MpfrClass& op1, const mpfr_t& op2);
MpfrClass operator / (const MpfrClass& op1, const double& op2);
MpfrClass operator / (const MpfrClass& op1, const cxsc::real& op2);
MpfrClass operator / (const MpfrClass& op1, const int op2);

MpfrClass operator / (const mpfr_t& op1, const MpfrClass& op2);
MpfrClass operator / (const double& op1, const MpfrClass& op2);
MpfrClass operator / (const cxsc::real& op1, const MpfrClass& op2);
MpfrClass operator / (const int op1, const MpfrClass& op2);

MpfrClass& operator /= (MpfrClass& op1, const MpfrClass& op2);
MpfrClass& operator /= (MpfrClass& op1, const mpfr_t& op2);
MpfrClass& operator /= (MpfrClass& op1, const double& op2);
MpfrClass& operator /= (MpfrClass& op1, const cxsc::real& op2);
MpfrClass& operator /= (MpfrClass& op1, const int op2);
```

### 3.10.5. Vom Current-Rundungsmodus abweichende Rundungen

Bei allen arithmetischen Operatoren  $\{+, -, *, /\}$  und  $\{+ =, - =, * =, / =\}$  werden die exakten Ergebnisse nach dem mit `SetCurrRndMode(...)` voreingestellten Current-Rundungsmodus in das jeweilige Ergebnisobjekt gerundet, siehe Seite 20.

Wenn jedoch, z.B. bei einer Multiplikation zweier `MpfrClass`-Objekte `a, b`, abweichend vom Current-Rundungsmodus z.B. ein Aufrunden verlangt wird, so kann dies mit Hilfe der MPFR-Funktion

```
int mpfr_mul (mpfr_t ROP, mpfr_t OP1, mpfr_t OP2, mpfr_rnd_t RND)
```

realisiert werden, wenn `RND` durch `RoundUp` ersetzt wird. Weitere Informationen findet man in den Dateien `mpfr.info`, `mpfr.pdf` oder `mpfr.dvi`, wobei die beiden letzten Dateien durch

```
make pdf bzw. make dvi
```

in dem Verzeichnis erzeugt werden können, in dem sich `mpfr.texi` befindet. Das nachfolgende Programm zeigt für das Auf- und Abrunden die entsprechenden Anweisungen:

```
1 // MPFR-02.cpp
2 #include <iostream>
3 #include <iomanip>
4 #include "mpfrclass.hpp"
5
6 using namespace MPFR;
7 using namespace std;
8
9 int main(void)
10 {
11     MpfrClass::SetCurrRndMode (RoundDown);
12     cout << "\nCurrent-RoundingMode = " << MpfrClass::GetCurrRndMode() << endl;
13     MpfrClass::SetCurrPrecision (20);
14     cout << "Current-Precision = " << MpfrClass::GetCurrPrecision() << endl;
15
16     MpfrClass a(1.234567890123456, RoundNearest, 20),
17                b(1.234567890123456, RoundNearest, 20), y;
18
19     mpfr_mul(y.GetValue(), a.GetValue(), b.GetValue(), RoundUp);
20     cout.precision(y.GetPrecision()/3.32193 + 2); // ca. 8 Dezimalstellenzahl
21     cout << "a*b (RoundUp) = " << y << endl;
22     cout << "y.GetPrecision() = " << y.GetPrecision() << endl << endl;
23
24     y = a*b; // Rundung mit RoundDown
25     cout.precision(y.GetPrecision()/3.32193 + 2); // ca. 8 Dezimalstellenzahl
26     cout << "a*b (RoundDown) = " << y << endl;
27     cout << "y.GetPrecision() = " << y.GetPrecision() << endl;
28
29     return 0;
30 }
```

Obiges Programm erzeugt die folgende Ausgabe:

```
Current-RoundingMode = 3
Current-Precision = 20
a*b (RoundUp) = 1.5241584
y.GetPrecision() = 20

a*b (RoundDown) = 1.5241565
y.GetPrecision() = 20
```

Im Vergleich zum exakten Produkt 1.5241578... erkennt man den aufgerundeten Wert und den bez. `Current-RoundingMode = RoundDown (3)` abgerundeten Produktwert.

Hier noch einige Anmerkungen zur MPFR-Funktion `mpfr_mul(...)`

1. Neben den C-XSC Funktionen können die in `mpfr.info` beschriebenen MPFR-Funktionen, wie z.B. `mpfr_mul(...)`, ebenfalls aufgerufen werden. Für die Implementierung weiterer Funktionen ist dies eine große Hilfe.
2. Ergebnisse werden in der Parameterliste der MPFR-Funktionen stets an **erster** Stelle abgelegt.
3. Wird die Ergebnis-Variable `y` vor dem MPFR-Funktionsaufruf durch z.B. `MpfrClass y;` lediglich deklariert, so erhält `y` die voreingestellte Current-Precision und zwar **unabhängig** von den Präzisionen beider Operanden in der Parameterliste.
4. Wird die Ergebnis-Variable `y` jedoch vor dem MPFR-Funktionsaufruf, z.B. durch

```
MpfrClass y(0, RoundNearest, 50);
```

mit der Präzision von 50 Bit initialisiert, so liefert `mpfr_mul(...)` das Ergebnis `y` mit der Präzision von 50 Bit und zwar wieder **unabhängig** von den Präzisionen beider Operanden.

Beachten Sie, dass bei den arithmetischen Operatoren mit dem Ergebnistyp `MpfrClass` die Ergebnis-Präzision stets gleich der Current-Precision ist, wobei gegebenenfalls das exakte Ergebnis bez. des Default-Rundungsmodus in den Ergebnisoperanden gerundet wird.

### 3.10.6. Arithmetik-Funktionen zum Auf- und Abrunden

Um abweichend vom voreingestellten Rundungsmodus bei den arithmetischen Grundoperationen auch das Auf- oder Abrunden zu ermöglichen, werden die folgenden acht Funktionen bereitgestellt, welche die gerundeten Ergebnisse in der CurrentPrecision zurückgeben:

```
MpfrClass addd(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass addu(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass subd(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass subu(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass muld(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass mulu(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass divd(const MpfrClass& op1, const MpfrClass& op2);
MpfrClass divu(const MpfrClass& op1, const MpfrClass& op2);
```

Die Funktion `MpfrClass addd(const MpfrClass& op1, const MpfrClass& op2)` liefert also die optimal abgerundete Summe  $(op1 + op2)$  in der mit `void SetCurrPrecision(prec)` voreingestellten CurrentPrecision.

## 3.11. Mathematische Funktionen

### 3.11.1. Standard-Implementierung

Die MPFR-Bibliothek stellt eine Vielzahl von Elementarfunktionen und einige Funktionen der Mathematischen Physik zur Verfügung, wobei zu einem Argument  $x$  mit beliebiger Präzision zunächst der exakte Funktionswert  $y_0$  berechnet wird. Danach wird dann  $y_0$  in eine Ergebnisvariable  $y$  mit einer möglichen anderen Präzision gerundet, wobei diese Rundung durch einen entsprechenden Parameter `rnd` gesteuert werden kann. Die Deklaration z.B. der Exponentialfunktion ist mit diesen Bezeichnungen gegeben durch:

```
int mpfr_exp (mpfr_t y, mpfr_t x, mpfr_rnd_t rnd)
```

Mit obiger MPFR-Funktion wird die Exponentialfunktion für das C-XSC Interface wie folgt implementiert:

```
MpfrClass exp (const MpfrClass& x, RoundingMode rnd)
{
    MpfrClass y(0); // Die Präzision von y ist jetzt die Current-Precision
    mpfr_exp(y.mpfr_rep, x.mpfr_rep, rnd);
    return y;
}
```

Mit Hilfe des Konstruktoraufrufs wird also zunächst die Präzision der Ergebnisvariablen  $y$  auf die Current-Precision gesetzt, vgl. dazu auch Seite 18. Danach wird mit dem Argument  $x$  und seiner Präzision der exakte Wert  $y_0 = e^x$  berechnet und dann mittels des Rundungsmodus `rnd` nach  $y$  gerundet. Wird `rnd` nicht gesetzt, so wird nach dem Current-Rundungsmodus gerundet und ist dieser nicht gesetzt, so wird zur jeweils nächsten Rasterzahl gerundet, (`RoundNearest`).

Soll nun ganz analog zur Exponentialfunktion im C-XSC Interface eine neue Funktion implementiert werden, die in der MPFR-Bibliothek noch nicht definiert ist, so muss wie folgt vorgegangen werden:

1. Sichere den Wert der ursprünglichen Current-Precision in `prec_old`
2. Setze die neue Current-Precision `prec` auf die Präzision vom Argument  $x$ . Wenn dann gilt `prec < prec_old`, setze `prec = prec_old`.
3. Setze mit `MpfrClass y(0);` die Präzision von  $y$  auf den Wert von `prec`.
4. Berechne in der neuen Current-Precision den mittels `rnd` gerundeten Funktionswert, der in  $y$  abzulegen ist.
5. Runde mittels `rnd` durch `y.RoundPrecision(prec_old, rnd);` auf die ursprüngliche Current-Precision.
6. Durch `SetCurrPrecision(prec_old);` die alte Current-Precision wiederherstellen.
7. Durch `return y;` den gerundeten Funktionswert zurückgeben, fertig!

#### Anmerkung:

- Grundsätzlich wird also der mittels `rnd` gerundete Funktionswert in der ursprünglichen Current-Precision `prec_old` zurückgegeben. Die interne Berechnung erfolgt jedoch in der durch `prec` vorgegebenen Präzision, die mindestens so groß ist wie `prec_old` selbst. Ein Beispiel dazu findet man in `mpfrclass.cpp` bei der Definition der `acoth`-Funktion.



### 3.11.2. Davon abweichende Funktionen und Konstanten

Bei nur wenige mathematische Funktionen ist es sinnvoll, bei ihrer Implementierung von dem auf Seite 32 angegebenen Schema abzuweichen. Die Ausnahmen sind:

```
MpfrClass abs (const MpfrClass& op, RoundingMode rnd, PrecisionType prec);
```

Zurückgegeben wird der i.a. gerundete Absolutbetrag von `op`. Für die Funktion gibt es drei verschiedene Aufrufmöglichkeiten:

1. `abs (op);`  
`|op|` wird bez. des Current-Rundungsmodus in die Current-Precision gerundet und zurückgegeben.
2. `abs (op, RoundUp);`  
`|op|` wird in die Current-Precision aufgerundet und zurückgegeben.
3. `abs (op, RoundDown, prec);`  
`|op|` wird in ein Format mit der Präzision `prec` abgerundet und zurückgegeben. Die Präzision des zurückgegebenen Wertes wird also i.a. von der voreingestellten Current-Precision verschieden sein!  
Will man jedoch `|op|` **rundungsfehlerfrei** mit der gleichen Präzision von `op` zurückgeben, so gelingt dies in allen Fällen, unabhängig von der voreingestellten Current-Precision, mit dem Funktionsaufruf:

```
abs (op, RoundNearest, op.GetPrecision());
```

wobei der Rundungsmodus (hier `RoundNearest`) natürlich beliebig gesetzt werden kann. Stimmt die Präzision von `op` mit der Current-Precision überein, so liefert der Aufruf `abs (op);` ebenfalls den **rundungsfehlerfreien** Wert von `|op|`. In der Praxis wird die rundungsfehlerfreie Rückgabe von `|op|` vermutlich immer im Vordergrund stehen.

Die Funktion `abs` kann also sehr flexibel eingesetzt werden und funktioniert nach 1. und 2. wie bei der Standard-Implementierung von Seite 32. Lediglich der Punkt 3. weicht von dieser Standard-Implementierung ab, um den exakten Wert von `|op|` in jedem Fall garantieren zu können.

**Beachten** Sie außerdem: Soll  $|k|$ ,  $k$  vom Typ `int`, berechnet werden, so muss dies mit `std::abs(k)` erfolgen.

```
MpfrClass mant (const MpfrClass& x);
```

Die Mantisse  $m$  von  $x$  wird rundungsfehlerfrei in der Präzision von  $x$  zurückgegeben. Im Fall einer normalen Gleitpunktzahl  $x \neq 0$  gilt:  $|m| \in [0.5, 1)$ , vgl. auch Seite 24.  
 $x = 0 \rightarrow m = 0$ ;  $x = \text{NaN} \rightarrow m = \text{NaN}$ ;  $x = \pm\text{Inf} \rightarrow m = \pm\text{Inf}$ ;

```
MpfrClass comp (const MpfrClass& x, const long int k);
```

Mit der Mantisse  $x$  und dem Zweierexponenten  $k$  wird  $x \cdot 2^k$  in der Präzision von  $x$  **rundungsfehlerfrei** zurückgegeben. Im Fall  $x = 0$  wird 0 zurückgegeben und in den Fällen  $x = \text{NaN}$  oder  $x = \pm\text{Inf}$  erhält man `NaN`. Ist  $x \neq 0$  eine normale Maschinenzahl und gilt:  $|x| \notin [0.5, 1)$  oder  $|k| > 1073741823$ , so wird ebenfalls `NaN` zurückgegeben.

```
MpfrClass min (const MpfrClass& op1, MpfrClass& op2);  
MpfrClass max (const MpfrClass& op1, MpfrClass& op2);
```

Zurückgegeben wird **rundungsfehlerfrei** das Minimum bzw. das Maximum beider Operanden und zwar genau in der Präzision des jeweils zurückgegebenen Operanden, die von der Current-Precision durchaus verschieden sein kann!

```

MpfrClass Round (const MpfrClass& op);
MpfrClass Floor (const MpfrClass& op);
MpfrClass Ceil (const MpfrClass& op);
MpfrClass Trunc (const MpfrClass& op);
MpfrClass Frac (const MpfrClass& op);

```

Die obigen fünf Rückgabewerte vom Typ `MpfrClass` haben alle die Präzision von `op`, die mit der Current-Precision nicht übereinstimmen muss. Die ersten vier Funktionen liefern den jeweiligen rundungsfehlerfreien ganzzahligen Anteil von `op`, wobei `Round()` weg von der Null rundet, falls `op` genau zwischen benachbarten ganzzahligen Werten liegt. Die letzte Funktion liefert den rundungsfehlerfreien nicht-ganzzahligen Teil von `op`.

```

MpfrClass minfloat (PrecisionType prec);
MpfrClass MaxFloat (PrecisionType prec);

```

Zurückgegeben werden der positive kleinste bzw. größte Zahlenwert im Datenformat der Präzision `prec`. Im Gegensatz zu `minfloat(prec)` fallen die entsprechenden Rückgabewerte von `MaxFloat(prec)` in Abhängigkeit von `prec` unterschiedlich aus! Werden die Funktionen ohne `prec` aufgerufen, so wird das Format mit der Current-Precision benutzt.

```

MpfrClass pred (const MpfrClass& op);
MpfrClass succ (const MpfrClass& op);

```

Beide Funktionen geben den Vorgänger bzw. Nachfolger von `op` zurück und zwar in der **gleichen** Präzision des Operanden `op`. Es macht nämlich keinen mathematischen Sinn, diese Werte anschließend in ein Format mit anderer Präzision zu runden, da sich die Funktionen `pred` und `succ` genau auf das Format von `op` beziehen.

```

void times2pown (MpfrClass& op, long int op1, RoundingMode rnd);

```

Obige Funktion liefert mit dem Eingabewert `op` den Wert  $op \cdot 2^{op1}$  mit gleicher Präzision zurück. Solange kein Über- oder Unterlauf entsteht, wird  $op \cdot 2^{op1}$  **exakt**, d.h. rundungsfehlerfrei berechnet. Tritt jedoch z.B. ein Überlauf ein, so wird gemäß `rnd` gerundet. Wenn `rnd` nicht gesetzt wird, so erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Ist dieser nicht gesetzt, so erfolgt die Rundung weg von der Null.

```

void Number_Scaling (MpfrClass& x, long int& k);

```

Obige Funktion liefert mit dem Eingabewert `x = a` den Rückgabewert `x = A` vom Typ `MpfrClass` in der Präzision von `a`, zusätzlich wird `k` zurückgegeben. Ist `a` unendlich oder ein NaN, so wird NaN und `k = 0` zurückgegeben. Im Fall `a = 0` wird ebenfalls `k = 0` zurückgegeben. Die Berechnung von  $A = a \cdot 2^{-k}$  wird stets **rundungsfehlerfrei** ausgeführt, wobei  $2 \cdot (A \cdot A)$  garantiert ohne Überlauf berechnet werden kann. Es gilt also stets:  $a = A \cdot 2^k$ .

```

void Number_Scaling_S (MpfrClass& x, long int& k);

```

Obige Funktion liefert mit dem Eingabewert `x = a` den Rückgabewert `x = A` vom Typ `MpfrClass` in der Präzision von `a`, zusätzlich wird `k` zurückgegeben. Ist `a` unendlich oder ein NaN, so wird NaN und `k = 0` zurückgegeben. Im Fall `a = 0` wird ebenfalls `k = 0` zurückgegeben. Die Berechnung von  $A = a \cdot 2^{-k}$  wird stets **rundungsfehlerfrei** ausgeführt, wobei  $4 \cdot A^4$  garantiert ohne Überlauf berechnet werden kann. Es gilt also stets:  $a = A \cdot 2^k$ .

```
void set_nan (MpfrClass& x);
```

Setzt  $x$  auf NaN, wobei die Präzision von  $x$  erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_inf (MpfrClass& x, int k);
```

Setzt  $x$  auf  $\pm\text{Inf}$ , wobei die Präzision von  $x$  erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss. Das Vorzeichen wird durch  $k$  festgelegt.

```
void set_zero (MpfrClass& x, int k);
```

Setzt  $x$  auf  $\pm 0$ , wobei die Präzision von  $x$  erhalten bleibt und daher mit der Current-Precision nicht übereinstimmen muss. Das Vorzeichen wird durch  $k$  festgelegt.

```
void swap(MpfrClass& x, MpfrClass& y);
```

Tauscht den Wert und die Präzision von  $x$  und  $y$ .

```
void swap(MpfrClass& x, mpfr_t& y);
```

Tauscht den Wert und die Präzision von  $x$  und  $y$ .

```
void random(MpfrClass& x, gmp_randstate_t state);
```

Geliefert wird ein Zufallszahl  $x \in [0, +1)$ , wobei  $x$  die Präzision des vorher deklarierten Objekts  $x$  erhält. Das folgende Programm zeigt eine mögliche Anwendung der Funktion `void random(MpfrClass& x, gmp_randstate_t state)`.

```
1 // MPFR-05.cpp
2 #include "mpfrclass.hpp"
3
4 using namespace MPFR;
5 using namespace std;
6
7 int main(void)
8 {
9     MpfrClass::SetCurrRndMode (RoundNearest);
10    cout << "\nCurrent-RoundingMode = " << MpfrClass::GetCurrRndMode() << endl;
11    MpfrClass::SetCurrPrecision (60);
12    cout << "Current-Precision = " << MpfrClass::GetCurrPrecision() << endl;
13
14    gmp_randstate_t state; // Declaration of state;
15    gmp_randinit_default (state); // Initialization of state;
16
17    MpfrClass x(0, RoundNearest, 50); // Declaration of class object x
18                                     // with a precision of 50 bits
19    cout.precision(x.GetPrecision()/3.321928095);
20    cout << "x = " << x << endl;
21    cout << "x.GetPrecision() = " << x.GetPrecision() << endl;
22    random (x, state); // Delivers the first random number x
23    cout << "x = " << x << endl;
24    random (x, state); // Delivers the second random number x
25    cout << "x = " << x << endl;
26
27    return 0;
28 }
```

Das Programm liefert die Ausgabe:

```
Current-RoundingMode = 0
Current-Precision = 60
x = 0
x.GetPrecision() = 50
x = 6.14544775142452e-1
x = 9.88050499009929e-1
```

Weitere Informationen bez. der Initialisierungsfunktion in Zeile 15 findet man unter

<http://gmplib.org/manual/Random-State-Initialization.html>

**Konstanten:**

```
static MpfrClass Pi      (RoundingMode rnd  = CurrRndMode,
                          PrecisionType prec = CurrPrecision);
static MpfrClass Ln2    (RoundingMode rnd  = CurrRndMode,
                          PrecisionType prec = CurrPrecision);
static MpfrClass Euler  (RoundingMode rnd  = CurrRndMode,
                          PrecisionType prec = CurrPrecision);
static MpfrClass Catalan(RoundingMode rnd  = CurrRndMode,
                          PrecisionType prec = CurrPrecision);
```

`Pi(rnd, prec)` rundet  $\pi$  mittels `rnd` in ein Format der Präzision `prec`. Wird `prec` nicht angegeben, so wird mittels `rnd` in ein Format mit der Current-Precision gerundet. Wird auch `rnd` weggelassen, so wird mittels des Current-Rundungsmodus in ein Format mit der Current-Precision gerundet. Entsprechendes gilt für die drei anderen Konstanten. `Ln2()` rundet also  $\ln(2) = 0.693147\dots$  mittels des Current-Rundungsmodus in ein Format mit der voreingestellten Current-Precision.

### 3.11.3. Elementarfunktionen

Tabelle 3.1.: Elementarfunktionen mit  $x, y, a, b$  vom Typ `MpfrClass`,  $n$ : `unsigned long int`;

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\ln(1+x)$	<code>lnp1(x)</code>
$x^2$	<code>sqr(x)</code>	$\log_2(x)$	<code>log2(x)</code>
$x^2 + y^2$	<code>x2py2(x,y)</code>	$\log_{10}(x)$	<code>log10(x)</code>
$x^2 - y^2$	<code>x2my2(x,y)</code>	$\ln(\sin(x))$	<code>ln_sin(x)</code>
$1/x$	<code>reci(x)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$x/(x^2 + y^2)$	<code>x_div_x2py2(x,y)</code>	$\ln(\sqrt{x^2 + y^2})$	<code>ln_sqrtx2y2(x,y)</code>
$(x^2 - y^2)/(x^2 + y^2)^2$	<code>Re_rz2(x,y)</code>	$\ln(\sqrt{(1+x)^2 + y^2})$	<code>ln_sqrtxp1_2y2(x,y)</code>
$2xy/(x^2 + y^2)^2$	<code>mIm_rz2(x,y)</code>	$x^k, k \in \mathbb{Z}$	<code>power(x,k)</code>
$\frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Re_r1pz2(x,y)</code>	$x^y$	<code>pow(x,y)</code>
$\frac{2xy}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Im_r1pz2(x,y)</code>	$\sin(x)$	<code>sin(x)</code>
$\sqrt{x}$	<code>sqr(x)</code>	$1/\sin(x)$	<code>csc(x)</code>
$\sqrt[n]{n}, n \in \mathbb{N}_0$	<code>sqrtn(n)</code>	$\cos(x)$	<code>cos(x)</code>
$1/\sqrt{x}$	<code>sqrtr(x)</code>	$1/\cos(x)$	<code>sec(x)</code>
$\sqrt[3]{x}, x \in \mathbb{R}$	<code>cbrt(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqrtn(x,n)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{x+1} - 1$	<code>sqrtp1m1(x)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt{1+x^2}$	<code>sqrtp1x2(x)</code>	$\arccos(x)$	<code>acos(x)</code>
$\sqrt{1-x^2}$	<code>sqrtp1mx2(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{x^2-1}$	<code>sqrtpx2m1(x)</code>	$\arctan(y/x)$	<code>atan2(y,x)</code>
$x/\sqrt{1+x^2}$	<code>xdsqrtp1x2(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$x/\sqrt{1-x^2}$	<code>xdsqrtp1mx2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$x/\sqrt{x^2-1}$	<code>xdsqrtpx2m1(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>

Fortsetzung auf der nächsten Seite

Fortsetzung der vorherigen Seite			
Funktion	Aufruf	Funktion	Aufruf
$\sqrt{x^2 + y^2}$	<code>sqrtox2y2(x, y)</code>	$\cosh(x)$	<code>cosh(x)</code>
$e^x$	<code>exp(x)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$2^x$	<code>exp2(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$10^x$	<code>exp10(x)</code>	$\coth(x)$	<code>coth(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{arcosh}(1 + x)$	<code>acoshp1(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\operatorname{AGM}(x, y)$	<code>agm(x, y)</code>

**Anmerkungen:**

1. Die Funktion `atan2(y, x)` ist wie folgt definiert:

$$\operatorname{atan2}(y, x) = \begin{cases} \arctan\left(\frac{y}{x}\right) & x > 0 \\ \pi + \arctan\left(\frac{y}{x}\right) & y \geq 0, x < 0 \\ -\pi + \arctan\left(\frac{y}{x}\right) & y < 0, x < 0 \\ \frac{\pi}{2} & y > 0, x = 0 \\ \frac{\pi}{2} & y = +\infty, x = 0 \\ -\frac{\pi}{2} & y < 0, x = 0 \\ -\frac{\pi}{2} & y = -\infty, x = 0 \\ 0 & y = 0, x = 0 \text{ Vorsicht!} \\ 0 & |y| < +\infty, x = +\infty \\ \pi & y \geq 0, x = -\infty \\ -\pi & y < 0, x = -\infty \\ \frac{\pi}{4} & y = +\infty, x = +\infty \\ -\frac{\pi}{4} & y = -\infty, x = +\infty \\ \frac{3\pi}{4} & y = +\infty, x = -\infty \\ -\frac{3\pi}{4} & y = -\infty, x = -\infty \\ \text{NaN} & y = \text{NaN} \text{ oder } x = \text{NaN}. \end{cases}$$

2. Mit `sqr(x)`; wird der exakte Funktionswert  $x^2$  mit dem Current-Rundungsmodus in die Current-Precision gerundet. Weitere Informationen zur Implementierung findet man auf Seite 32.

3. Mit `sqr(x, rnd)`; wird wie unter 2. verfahren, jedoch wird der exakte Funktionswert  $x^2$  jetzt mittels `rnd` in die Current-Precision gerundet, vgl. ebenfalls Seite 32.
4. Alle Funktionen in obiger Tabelle sind nach 2. bzw. 3. implementiert. Für `rnd` stehen dabei die folgenden Rundungsmodi zur Verfügung<sup>1</sup>: `RoundNearest`, `RoundUp`, `RoundDown`.
5. Bei allen Funktionen aus obiger Tabelle wird vorausgesetzt, dass  $x, y, a, b$  exakte Maschinenzahlen sind, die also nicht schon durch vorhergehende Rechnungen gerundet worden sind. Beispielsweise kann daher eine Maschinenzahl  $a$  nicht den Wert  $\sqrt{3}$  annehmen, wohl aber eine geeignete Näherung, die dann aber als exakt anzusehen ist.
6. `abs(x)` kann zusätzlich noch mit einem Präzisionsparameter `prec` aufgerufen werden. Mit `abs(x, rnd, x.GetPrecision())`; wird dabei der exakte Absolutbetrag  $|x|$ , und zwar unabhängig von `rnd`, zurückgegeben, vgl. auch Seite 33.
7. `power(x, k)` kann bez. `k` mit den Datentypen `int`, `long int` aufgerufen werden.
8. `pow(x, y)` kann bez. `y` mit folgenden Datentypen aufgerufen werden: `MpfrClass`, `real`, `double`.
9. Die Funktion `agm(x, y)` rundet das exakte Arithmetisch-Geometrische Mittel  $AGM(x, y)$  der beiden `MpfrClass`-Objekte  $x, y \geq 0$  mit dem Current-Rundungsmodus optimal in ein `MpfrClass`-Objekt mit der Current-Precision. Der absolute Fehler ist dabei höchstens 0.5 ulp. Mit `agm(x, y, RoundDown)` wird das AGM optimal in die Current-Precision abgerundet, d.h. der absolute Fehler ist kleiner als 1 ulp. Das AGM ist bez.  $x, y$  streng monoton wachsend und es gilt:

$$0 \leq x \leq AGM(x, y) = AGM(y, x), \text{ agm}(x, y) = \text{agm}(y, x) \leq y.$$

Das Arithmetisch-geometrische Mittel AGM spielt bei der Auswertung elliptischer Integrale eine zentrale Rolle.

---

<sup>1</sup>Für alle Funktionen, die in MPFR direkt implementiert sind, gibt es bez. `rnd` die zusätzlichen Rundungsmodi: `RoundToZero`, `RoundFromZero`.

### 3.11.4. Skalarprodukt aus zwei Teilprodukten

Mit den folgenden Funktionen lassen sich Skalarprodukte der Form  $a \cdot b + c \cdot d$  ohne vorzeitigen Overflow berechnen, wobei nach Bedarf gerundet werden kann.

```
void prod_H1(MpfrClass& r, long int& k, const MpfrClass& a,  
             const MpfrClass& b, RoundingMode rnd);
```

$r$  wird zunächst auf die Current-Precision gesetzt, dann wird  $a \cdot b$  bez.  $rnd$  nach  $r$  gerundet, so dass gilt:  $r \cdot 2^k \approx a \cdot b$ .

Wird  $rnd$  nicht gesetzt, so wird mit dem voreingestellten Current-Rundungsmodus nach  $r$  gerundet. Es gilt:  $|r| \in [0.5, 2)$ .

Es gilt:  $r = 0$ ,  $r = \text{NaN}$ ,  $r = +\text{Inf}$ ,  $r = -\text{Inf}$  --->  $k = 0$ ;

```
void Prod_H1(MpfrClass& r, long int& k, const MpfrClass& a,  
             const MpfrClass& b, RoundingMode rnd);
```

$r$  wird zunächst auf die Current-Precision gesetzt, dann wird  $a \cdot b$  bez.  $rnd$  nach  $r$  gerundet, so dass gilt:  $r \cdot 2^k \approx a \cdot b$ .

Wird  $rnd$  nicht gesetzt, so wird mit dem voreingestellten Current-Rundungsmodus nach  $r$  gerundet. Es gilt:  $\text{MaxFloat}()/32 \leq |r| < \text{MaxFloat}()/8$ .

Es gilt außerdem:  $r = 0$ ,  $r = \text{NaN}$ ,  $r = +\text{Inf}$ ,  $r = -\text{Inf}$  --->  $k = 0$ ;

```
void sum_k_H1(MpfrClass& r, long int& k, const MpfrClass& x, long int& kx,  
             const MpfrClass& y, long int& ky, RoundingMode rnd);
```

Voraussetzung:  $kx \geq ky$ , falls  $(x, y \neq 0 \ \&\& \ \text{isNumber}(x, y) = \text{True})$ .

$r$  wird zunächst auf das Maximum der Präzisionen von  $x$  und  $y$  gesetzt. Ist diese kleiner als die Current-Precision, so erhält  $r$  als Präzision diese Current-Precision. Danach wird durch geeignete Skalierung  $x \cdot 2^{kx} + y \cdot 2^{ky} = 2^{kx} \cdot (x + y \cdot 2^{ky-kx})$  die Klammer  $(x + y \cdot 2^{ky-kx})$  berechnet und mittels  $rnd$  nach  $r$  gerundet. Man erhält dann im Fall  $x, y \neq 0$  das Ergebnis:  $r \cdot 2^k = r \cdot 2^{kx} \approx x \cdot 2^{kx} + y \cdot 2^{ky}$ .

Wird  $rnd$  nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus.

Es gilt außerdem:  $r = 0$ ,  $r = \text{NaN}$ ,  $r = +\text{Inf}$ ,  $r = -\text{Inf}$  --->  $k = 0$ ;

```
void scal_prod_k(MpfrClass& r, long int& k,  
                const MpfrClass& a, const MpfrClass& b,  
                const MpfrClass& c, const MpfrClass& d,  
                RoundingMode rnd);
```

$r$  wird zunächst auf das Maximum der Präzisionen von  $a, b, c, d$  gesetzt. Ist diese kleiner als die Current-Precision, so erhält  $r$  als Präzision diese Current-Precision.

Danach wird mit Hilfe der obigen Funktionen  $\text{prod\_H1}()$  und  $\text{sum\_k\_H1}()$  das Skalarprodukt  $a \cdot b + c \cdot d = r \cdot 2^k$  mittels  $rnd$  nach  $r$  gerundet und entsprechend  $k$  berechnet.

Wird  $rnd$  nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus.

```
void Scal_prod_k(MpfrClass& r, long int& k,  
                const MpfrClass& a, const MpfrClass& b,  
                const MpfrClass& c, const MpfrClass& d,  
                RoundingMode rnd);
```

$r$  wird zunächst auf das Maximum der Präzisionen von  $a, b, c, d$  gesetzt. Ist diese kleiner als die Current-Precision, so erhält  $r$  als Präzision diese Current-Precision.

Um Auslöschung möglichst zu vermeiden, wird dann diese Präzision von  $r$  noch einmal mehr als verdoppelt. Danach wird mit Hilfe der obigen Funktionen  $\text{Prod\_H1}()$  und  $\text{sum\_k\_H1}()$  das Skalarprodukt  $a \cdot b + c \cdot d = r \cdot 2^k$  mittels  $rnd$  nach  $r$  gerundet und entsprechend  $k$  berechnet. Wird  $rnd$  nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus.



```
void scal_prod(MpfrClass& r, const MpfrClass& a, const MpfrClass& b,
               const MpfrClass& c, const MpfrClass& d,
               RoundingMode rnd);
```

$a \cdot b + c \cdot d$  wird intern mit `Skal_prod.k()` in so hoher Präzision berechnet, dass es danach mit `rnd` nach `r` höchstens zweimal in die Current-Präzision gerundet werden muss. Mit `rnd = RoundDown` und `rnd = RoundUp` erhält man damit für  $a \cdot b + c \cdot d$  eine nahezu maximalgenau Einschließung. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus. Ein vorzeitiger interner Über- oder Unterlauf ist ausgeschlossen.

```
MpfrClass ComplRe (const MpfrClass& a, const MpfrClass& b,
                  const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

Rundet den Realteil von  $(a + i \cdot b)/(x + i \cdot y)$  in die Current-Precision und gibt diesen gerundeten Wert zurück. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus. Ein vorzeitiger interner Über- oder Unterlauf wird vermieden.

```
MpfrClass ComplIm (const MpfrClass& a, const MpfrClass& b,
                  const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

Rundet den Imaginärteil von  $(a + i \cdot b)/(x + i \cdot y)$  in die Current-Precision und gibt diesen gerundeten Wert zurück. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus. Ein vorzeitiger interner Über- oder Unterlauf wird vermieden.

```
MpfrClass plus_ab (const MpfrClass& x, const MpfrClass& a,
                  const MpfrClass& b, RoundingMode rnd);
```

$(x + a \cdot b)$  wird bez. `rnd` in den Rückgabewert mit der Current-Precision gerundet. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus. Durch zweimalige Anwendung dieser Funktion:

$$x = \text{plus\_ab}(0, a, b); \quad y = \text{plus\_ab}(x, c, d);$$

wird das Skalarprodukt  $a \cdot b + c \cdot d$  jetzt mit bis zu zwei Rundungen berechnet, und im Gegensatz zur Funktion `scal_prod()` kann ein vorzeitiger Überlauf nicht vermieden werden. Die Funktion `scal_prod()` benötigt darüber hinaus i.a. nur eine Rundung und liefert damit für  $a \cdot b + c \cdot d$  i.a. die **bessere** Näherung.

### 3.11.4.1. Beispiel

Als Beispiel für die Berechnung eines Skalarprodukts der Form  $x = a \cdot b + c \cdot d$  wählen wir:

$$a := \text{MaxFloat}(), \quad b := e^1 \cdot 2^k, \quad k \in \mathbb{Z}, \quad c := -b, \quad d := \text{pred}(a), \quad \text{und damit gilt:}$$

$$(3.3) \quad x := a \cdot b + c \cdot d = e^1 \cdot 2^k (a - \text{pred}(a)) =: y,$$

wobei man zur Kontrolle  $y = e^1 \cdot 2^k (a - \text{pred}(a))$  auch ohne Skalarprodukt berechnen kann.

```
1 // MPFR-03.cpp
2 #include <iostream>
3 #include <iomanip>
4 #include "mpfrclass.hpp"
5
6 using namespace MPFR;
7 using namespace std;
8
9 int main(void)
10 {
11     MpfrClass::SetCurrRndMode (RoundNearest);
12     cout << "\nCurrent-RoundingMode = " << MpfrClass::GetCurrRndMode() << endl;
13     MpfrClass::SetCurrPrecision (5000000);
14     cout << "Current-Precision = " << MpfrClass::GetCurrPrecision() << endl;
15
16     MpfrClass a,b,c,d,x,y;
17
18     long int k = -2; // 1234567;
19     a = MaxFloat(); b = exp(MpfrClass(1)); times2pown(b, k);
20     c = -b; d = pred(a);
21
22     y = b*(a-pred(a));
23     scal_prod (x, a, b, c, d, RoundDown);
24
25     cout.precision ();
26     cout << "x = " << x << endl;
27
28     if (x==y) cout << "Kontrollrechnung: x == y" << endl;
29     else cout << "x != y" << endl;
30
31     y = 0;
32     y = plus_ab(y, a, b, RoundDown);
33     x = plus_ab(y, c, d, RoundDown);
34
35     cout << "Zweimalige Anwendg. von plus_ab(), x = " << x << endl;
36
37     return 0;
38 }
```

Mit RoundDown und RoundUp in Zeile 23 liefert das Programm MPFR-03.cpp die folgende Ausgabe

```
Current-RoundingMode = 0
Current-Precision = 5000000
x = 1.49913e321723346
Kontrollrechnung: x == y
Zweimalige Anwendg. von plus_ab(), x = 7.92265e321723345
```

Man erkennt, dass durch die zweimalige Anwendung von `plus_ab(..., RoundDown)`, verbunden mit mindestens einer Rundung (`RoundDown`), im Vergleich zum exakten Wert `x` eine deutlich kleinerer Näherung `y` für das Skalarprodukt berechnet wird. Die Ausgabe der Kontrollrechnung zeigt, dass das Skalarprodukt mit `scal_prod(x, a, b, c, d, RoundDown);` sogar **exakt** berechnet wurde. Dies wird auch durch die Rechnung mit `RoundUp` in Zeile 23 bestätigt.

Wählt man im Programm in Zeile 18 `k = 1234567;` so werden beide Teilprodukte  $a \cdot b$  und  $c \cdot d$ , jeweils einzeln berechnet, einen Overflow erzeugen. Das Programm liefert jedoch die Ausgabe

```
Current-RoundingMode = 0
Current-Precision = 5000000
x = 2.99610e322094988
Kontrollrechnung: x == y
Zweimalige Anwendg. von plus_ab(), x = -@Inf@
```

Mit `scal_prod(x, a, b, c, d, RoundDown);` erhält man wieder das **exakte** Skalarprodukt  $x = 2.996\dots \cdot 10^{322094988}$ , jetzt jedoch ohne vorzeitigen internen Overflow, der aber schon beim ersten Aufruf der Funktion `plus_ab(..., RoundDown)` natürlich nicht vermieden werden kann!

### Anmerkungen:

- Wählt man im obigen Programm die `CurrentPrecision` mit `prec = 500000` um den Faktor 10 kleiner, so wird das exakte Skalarprodukt zu groß und kann wegen notwendiger Rundungen nicht mehr exakt berechnet werden, d.h. `x` erhält in Abhängigkeit von `RoundUp` und `RoundDown` **verschiedene** Werte. Beachten Sie in diesem Zusammenhang, dass der Faktor  $(a - \text{pred}(a))$  im Ausdruck für  $y$  nur dann hinreichend klein wird, wenn man die `Current-Precision` hinreichend groß gewählt.
- Der Ausdruck  $b \cdot (a - \text{pred}(a))$  kann rundungsfehlerfrei berechnet werden, da  $a - \text{pred}(a)$  die Binärdarstellung  $1.0000e\dots$  besitzt und damit eine reine Zweierpotenz ist, mit der  $b$  rundungsfehlerfrei multipliziert werden kann, solange dabei kein Über- oder Unterlauf entsteht.

### 3.11.5. Funktionen der Mathematischen Physik

Tabelle 3.2.: Funktionen der Mathematischen Physik mit x vom Typ MpfrClass

Funktionsterm	Aufruf	Anmerkung
$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>	monoton wachsend
$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>	monoton fallend
$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt, x > 0$	<code>gamma(x)</code>	Pole: x=0, -1, -2, ...
$\Gamma'(x)$	<code>gamma_D(x)</code>	Pole: x=0, -1, -2, ...
$\frac{1}{\Gamma(x)}$	<code>gamma_reci(x)</code>	überall differenzierbar
$\left(\frac{1}{\Gamma(x)}\right)'$	<code>gamma_reci_D(x)</code>	überall differenzierbar
$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$	<code>digamma(x)</code>	Pole: x=0, -1, -2, ...
$\ln(\Gamma(x))$	<code>lngamma(x)</code>	$2k-1 \leq x \leq -2k \rightsquigarrow \text{NaN}, k=0,1,2,\dots$
$\ln( \Gamma(x) )$	<code>int k; lngamma(x,k);</code>	$k = \begin{cases} +1, & \Gamma(x) > 0 \\ -1, & \Gamma(x) < 0 \end{cases}$
$k!$	<code>factorial(k)</code>	unsigned long int k
$\zeta(x) = \sum_{k=1}^\infty k^{-x}, x > 1;$	<code>zeta(x)</code>	$x \neq +1$
$\zeta(k), k = 0, 2, 3, 4, \dots$	<code>zeta(k)</code>	unsigned long int k
$\operatorname{Ei}(x) = \gamma + \ln(x) + \sum_{k=1}^\infty \frac{x^k}{k \cdot k!}, x > 0;$	<code>Ei(x)</code>	$x=0 \rightsquigarrow -\text{Inf}; x < 0 \rightsquigarrow \text{NaN};$
$\operatorname{Li2}(x) = -\int_0^x \frac{\ln(1-t)}{t} dt, x < 1;$	<code>Li2(x)</code>	$x > 1 \rightsquigarrow \text{Nur Realteil!}$
$J_n(x) = \sum_{k=0}^\infty \frac{(-1)^k}{k! \Gamma(k+n+1)} \left(\frac{x}{2}\right)^{2k+n};$	<code>Jn(n,x)</code>	Bessel-Fkt. 1. Art; $n \in \mathbb{Z};$
$J_0(x);$	<code>J0(x)</code>	Bessel-Fkt. 1. Art;
$J_1(x);$	<code>J1(x)</code>	Bessel-Fkt. 1. Art;
$Y_n(x);$	<code>Yn(n,x)</code>	Bessel-Fkt. 2. Art; $n \in \mathbb{Z};$
$Y_0(x);$	<code>Y0(x)</code>	Bessel-Fkt. 2. Art;
$Y_1(x);$	<code>Y1(x)</code>	Bessel-Fkt. 2. Art;

### Anmerkungen:

1. Alle Funktionen aus obiger Tabelle können mit einem zusätzlichen Rundungsparameter `rnd` aufgerufen werden, womit eine vom Current-Rundungsmodus abweichende Rundung realisiert werden kann. Bezüglich der Implementierung gelten die gleichen Anmerkungen wie auf Seite 32.
2. Die Funktion  $\text{Li2}(x)$  ist für  $x < 1$  zunächst definiert durch

$$\text{Li2}(x) := - \int_0^x \frac{\ln(1-t)}{t} dt = \sum_{k=1}^{\infty} \frac{x^k}{k^2}, \quad x < 1;$$

Durch analytische Fortsetzung lässt sich diese Definition auf weitere  $z \in \mathbb{C}$  in der ab  $x = 1$  längs der positiven reellen Achse aufgeschnittenen komplexen Ebene ausdehnen. Für  $x > 1$  liefert  $\text{Li2}(x)$  dabei nur den Realteil der dann komplexwertigen Funktion. Der Imaginärteil kann mit Hilfe der Beziehung

$$\text{Li2}(x) = -\text{Li2}(1/x) + \frac{\pi^2}{3} - \frac{1}{2} \ln^2(x) - i \cdot \pi \ln(x), \quad x > 1;$$

einfach berechnet werden, da  $\pi \ln(x)$  für  $x > 1$  positiv und streng monoton wachsend ist.

Den Zusammenhang mit der Polylogarithmus-Funktion  $\text{Li}_s(z) = \sum_{k=1}^{\infty} \frac{z^k}{k^s}$  findet man unter

<http://de.wikipedia.org/wiki/Polylogarithmus>

<http://en.wikipedia.org/wiki/Polylogarithm>

3. Die Besselfunktionen 2. Art  $Y_n(x)$  (Neumann-Funktionen) definiert man zunächst für nichtganzzahlige  $\nu$  als Linearkombination von  $J_\nu$  und  $J_{-\nu}$

$$Y_\nu(x) := \frac{J_\nu(x) \cdot \cos \pi \nu - J_{-\nu}(x)}{\sin \pi \nu}, \quad \nu \notin \mathbb{Z},$$

und bildet für  $n \in \mathbb{Z}$  die Funktionen  $Y_n(x)$  durch Grenzübergang

$$Y_n(x) := \lim_{\nu \rightarrow n} Y_\nu(x), \quad n \in \mathbb{Z};$$

Dies ist lediglich eine Anmerkung zur Definition, nicht aber zur numerischen Auswertung von  $Y_n(x)$ . Die Funktionen  $Y_n(x)$  und  $J_n(x)$  sind für festes  $n \in \mathbb{Z}$  linear unabhängig und bilden damit eine Basis für die Lösungen der Besselschen Differentialgleichung.



## 4. MpfiClass-Interface zur Anbindung der MPFI-Bibliothek an C-XSC

### 4.1. MPFI-Bibliothek

#### 4.1.1. Entwickler

Die MPFI-Bibliothek wurde 2002 von der INRIA-Gruppe an der Universität von Lyon entwickelt. Die Gruppe besteht aus Nathalie Revol, Fabrice Rouillier, Sylvain Chevillard, Hong Diep NGUYEN und Christoph Lauter. Unterstützt wurden sie von den Entwicklern der MPFR-Bibliothek, zu der auch Nathalie Revol gehört.

#### 4.1.2. Allgemein

Die MPFI-Bibliothek wurde für langzahlige Intervallberechnungen entwickelt und ist nur in C implementiert. Sie basiert auf der GMP- und MPFR-Bibliothek. Die Bibliothek entstand im Jahre 2002 und wird seitdem immer weiter entwickelt. Die aktuelle Version ist "1.5". Um die Bibliothek nutzen zu können, müssen die GMP-Bibliothek (Version 4.1 oder höher) und die MPFR-Bibliothek (Version 3.0.0 oder höher) vorhanden sein. Da die Bibliothek auf der MPFR-Bibliothek basiert, profitiert MPFI von den korrekten Rundungen der MPFR-Bibliothek. Die Bibliothek befolgt den IEEE-754-Standard für Gleitkommaarithmetik.

Intern wird die Bibliothek mit Hilfe von zwei MPFR-Variablen realisiert. Die beiden Variablen repräsentieren die beiden Endpunkte des Intervalls. Ein Intervall kann endliche oder unendliche Endpunkte haben. Ebenso kann ein Endpunkt - oder beide - @NaN@ sein. Dies zeigt an, dass eine unzulässige Berechnung durchgeführt wurde. Ein leeres Intervall ist dadurch definiert, dass der rechte Endpunkt kleiner als der linke Endpunkt ist.

#### 4.1.3. Installation

Die Installation der MPFI-Bibliothek ist nur unter einem Linux/Unix-System möglich und erfolgt völlig analog zur Beschreibung auf Seite 13. Die aktuelle Version der MPFI-Bibliothek kann unter:

`http://gforge.inria.fr/projects/mpfi/`

bezogen werden.

## 4.2. Grundlegendes

Das Mpficlass-Interface ist eine in `mpficlass.hpp` und `mpficlass.cpp` implementierte C++-Wrapper-Klasse `Mpficlass` für die C-Bibliothek MPFI, deren C-Funktionen über die implementierten Operatoren und Funktionen aufgerufen werden. Alle Funktionen mit der möglichen Übergabe eines RoundingModes oder eines PrecisionTypes haben als Default-Wert den Wert von CurrRndMode bzw. CurrPrecision, die beide beliebig gesetzt werden können. Dies gilt in vielen Fällen auch für die Konstruktoren.

### 4.2.1. Allgemein

Um in C-XSC das Interface verwenden zu können, muss der Header `mpficlass.hpp` eingebunden werden. In ihm sind die benötigten Header-Dateien der MPFI-Bibliothek enthalten. Die Mpficlass-Klasse liegt im Namensraum "MPFI".

### 4.2.2. Aufbau

Die Klasse besteht intern aus einer `mpfi_t`-Variablen. Diese dient zum Speichern eines Intervalls. Zusätzlich gibt es static Elemente, um den Standard-Rundungsmodus, die Standard-Präzision und die aktuelle Basis zu speichern.

### 4.2.3. Präzision

Die Präzision gibt die Anzahl der binären Mantissenstellen einer Mpficlass-Variablen an. Der Wert muss mindestens 2 betragen. Die Current-Precision kann global gesetzt werden; wenn dies nicht geschieht, so wird die Default-Precision mit 53 Bits benutzt. Unabhängig davon kann die Präzision für jede Mpficlass-Variable auch einzeln festgelegt werden.

### 4.2.4. Variablentyp PrecisionType

Mit Hilfe des Variablentyps `PrecisionType`, dessen Variablen i.a. mit `prec` bezeichnet werden, kann der Präzisionswert einer Mpficlass-Variablen eingestellt werden. Der Variablentyp ist ein typedef für eine `mpfr_prec_t`-Variable.



## 4.3. Konstruktoren / Destruktoren

### 4.3.1. Konstruktoren

```
MpfiClass ();
```

Der Default-Konstruktor legt ein neues Intervall-Objekt mit Current-Precision an.

Aufruf: `MpfiClass y;`

Es wird kein Wert initialisiert und man erhält: `y = [NaN, NaN];`

Alle folgenden Konstruktoren erzeugen ein `MpfiClass`-Intervall `y` mit der Präzision `prec`, das die jeweiligen Werte `op` bzw. das Intervall `[op1,op2]` optimal einschließen. Wird beim Konstruktoraufruf `prec` nicht angegeben, so erhält `y` die voreingestellte Current-Precision. Falls diese mit `SetCurrPrecision` nicht gesetzt wurde, erhält `y` die Präzision `prec = 53`.

```
MpfiClass (const MpfiClass& op, PrecisionType prec);  
MpfiClass (const mpfi_t& op, PrecisionType prec);  
MpfiClass (const MFR::MpfrClass& op, PrecisionType prec);  
MpfiClass (const mpfr_t& op, PrecisionType prec);  
MpfiClass (const cxsc::interval& op, PrecisionType prec);  
MpfiClass (const cxsc::real& op, PrecisionType prec);  
MpfiClass (const double& op, PrecisionType prec);  
MpfiClass (int op, PrecisionType prec);  
MpfiClass (const MFR::MpfrClass& op1,  
const MFR::MpfrClass& op2, PrecisionType prec);
```

```
MpfiClass (const std::string& op, PrecisionType prec);
```

Auch hier gelten sinngemäß die gleichen Anmerkungen wie für die oberen neun Konstruktoren, wobei unter dem Wert von `op` der Dezimalwert des Strings zu verstehen ist. Mit den beiden Strings `op = "[1.0e-1, 1.0e-1]"` bzw. `op = "1.0e-1"` liefert der Konstruktoraufruf `MpfiClass y(op);` mit der voreingestellten Current-Precision Einschließungen des Dezimalwertes 0.1, wobei zu beachten ist, dass 0.1 im binären System nicht darstellbar ist, so dass alle Einschließungen von 0.1 auch mit noch so großen `prec`-Werten keine Punktintervalle sein können. Der Konstruktoraufruf `MpfiClass y(op, 140);` liefert eine Einschließung von 0.1 mit einer Präzision von 140 Bits, wobei auch jetzt 0.1 natürlich nicht durch eine Punktintervall eingeschlossen werden kann.

### 4.3.2. Destruktor

```
~MpfiClass ();
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

## 4.4. Leeres Intervall

In seltenen Situationen benötigt man ein Intervall, dessen untere Grenze größer ist als sein obere Grenze. Dies wird realisiert durch die Funktion

```
MpfiClass EmptyIntVal ();
```

welche das Intervall `[999999999.0, -999999999.0]` erzeugt, vgl. auf Seite 58 die Member-Funktion `bool isEmpty()`.

## 4.5. Zuweisungs-Operatoren

Unabhängig von der Präzision des rechten Operanden erhält bei allen folgenden Zuweisungsoperatoren der linke Operand  $y$  stets die Current-Precision und schließt die jeweiligen `op`-Werte optimal ein.

```
MpfiClass& operator = (const MpfiClass& op);  
MpfiClass& operator = (const MPFR::MpfrClass& op);  
MpfiClass& operator = (const mpfi_t& op);  
MpfiClass& operator = (const mpfr_t& op);  
MpfiClass& operator = (const int& op);  
MpfiClass& operator = (const cxsc::real& op);  
MpfiClass& operator = (const double& op);  
MpfiClass& operator = (const cxsc::interval& op);  
MpfiClass& operator = (const std::string& op);
```

### Anmerkungen:

- Ist z.B. `op` ein `real`-Wert und wurde die Current-Precision mit `SetCurrPrecision` zu klein gewählt, so ist der linke Operand  $y$  i.a. kein Punktintervall. Ist die Current-Precision jedoch größer oder gleich 53, so ist das einschließende Intervall  $y$  stets ein Punktintervall.
- Ist `op` ein String, z.B. `op = "[0.1,0.1]"` oder `op = "0.1"`, so kann  $y$  bei noch so großer Current-Precision kein Punktintervall sein, da 0.1 im vorliegenden Binärsystem nicht exakt darstellbar ist.

## 4.6. Arithmetische Operatoren

Für alle arithmetischen Intervalloperationen gilt:

Das exakte Ergebnis einer arithmetischen Intervalloperation wird unabhängig von der Präzision der Operanden mit der voreingestellten Current-Präzision optimal eingeschlossen.

Die Operatoren  $\odot =$ , mit  $\odot \in \{+, -, \cdot, /\}$ , bedeuten  $u \odot = v \iff u = u \odot v$ . Dabei wird  $u \odot v$  optimal durch  $u$  eingeschlossen, wobei  $u$  als Präzision die Current-Präzision erhält.

### 4.6.1. Addition

```
MpfiClass operator + (const MpfiClass& op1, const MpfiClass& op2);
MpfiClass operator + (const MpfiClass& op1, const double& op2);
MpfiClass operator + (const MpfiClass& op1, const cxsc::real& op2);
MpfiClass operator + (const MpfiClass& op1, const cxsc::interval& op2);
MpfiClass operator + (const MpfiClass& op1, const int op2);
MpfiClass operator + (const MpfiClass& op1, const mpfr_t& op2);
MpfiClass operator + (const MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass operator + (const MpfiClass& op1, const mpfi_t& op2);
MpfiClass operator + (const double& op1, const MpfiClass& op2);
MpfiClass operator + (const cxsc::real& op1, const MpfiClass& op2);
MpfiClass operator + (const cxsc::interval& op1, const MpfiClass& op2);
MpfiClass operator + (const int op1, const MpfiClass& op2);
MpfiClass operator + (const mpfr_t& op1, const MpfiClass& op2);
MpfiClass operator + (const MPFR::MpfrClass& op1, const MpfiClass& op2);
MpfiClass operator + (const mpfi_t& op1, const MpfiClass& op2);
```

```
MpfiClass& operator += (MpfiClass& op1, const MpfiClass& op2);
MpfiClass& operator += (MpfiClass& op1, const double& op2);
MpfiClass& operator += (MpfiClass& op1, const cxsc::real& op2);
MpfiClass& operator += (MpfiClass& op1, const cxsc::interval& op2);
MpfiClass& operator += (MpfiClass& op1, const int op2);
MpfiClass& operator += (MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass& operator += (MpfiClass& op1, const mpfr_t& op2);
MpfiClass& operator += (MpfiClass& op1, const mpfi_t& op2);
```

## 4.6.2. Subtraktion

Beachten Sie den Hinweis auf Seite 51.

```
MpfiClass operator - (const MpfiClass& op1, const MpfiClass& op2);
MpfiClass operator - (const MpfiClass& op1, const double& op2);
MpfiClass operator - (const MpfiClass& op1, const cxsc::real& op2);
MpfiClass operator - (const MpfiClass& op1, const cxsc::interval& op2);
MpfiClass operator - (const MpfiClass& op1, const int op2);
MpfiClass operator - (const MpfiClass& op1, const mpfr_t& op2);
MpfiClass operator - (const MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass operator - (const MpfiClass& op1, const mpfi_t& op2);
MpfiClass operator - (const double& op1, const MpfiClass& op2);
MpfiClass operator - (const cxsc::real& op1, const MpfiClass& op2);
MpfiClass operator - (const cxsc::interval& op1, const MpfiClass& op2);
MpfiClass operator - (const int op1, const MpfiClass& op2);
MpfiClass operator - (const mpfr_t& op1, const MpfiClass& op2);
MpfiClass operator - (const MPFR::MpfrClass& op1, const MpfiClass& op2);
MpfiClass operator - (const mpfi_t& op1, const MpfiClass& op2);
```

```
MpfiClass& operator -= (MpfiClass& op1, const MpfiClass& op2);
MpfiClass& operator -= (MpfiClass& op1, const double& op2);
MpfiClass& operator -= (MpfiClass& op1, const cxsc::real& op2);
MpfiClass& operator -= (MpfiClass& op1, const cxsc::interval& op2);
MpfiClass& operator -= (MpfiClass& op1, const int op2);
MpfiClass& operator -= (MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass& operator -= (MpfiClass& op1, const mpfr_t& op2);
MpfiClass& operator -= (MpfiClass& op1, const mpfi_t& op2);
```

## 4.6.3. Multiplikation

Beachten Sie den Hinweis auf Seite 51.

```
MpfiClass operator * (const MpfiClass& op1, const MpfiClass& op2);
MpfiClass operator * (const MpfiClass& op1, const double& op2);
MpfiClass operator * (const MpfiClass& op1, const cxsc::real& op2);
MpfiClass operator * (const MpfiClass& op1, const cxsc::interval& op2);
MpfiClass operator * (const MpfiClass& op1, const int op2);
MpfiClass operator * (const MpfiClass& op1, const mpfr_t& op2);
MpfiClass operator * (const MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass operator * (const MpfiClass& op1, const mpfi_t& op2);
MpfiClass operator * (const double& op1, const MpfiClass& op2);
MpfiClass operator * (const cxsc::real& op1, const MpfiClass& op2);
MpfiClass operator * (const cxsc::interval& op1, const MpfiClass& op2);
MpfiClass operator * (const int op1, const MpfiClass& op2);
MpfiClass operator * (const mpfr_t& op1, const MpfiClass& op2);
MpfiClass operator * (const MPFR::MpfrClass& op1, const MpfiClass& op2);
MpfiClass operator * (const mpfi_t& op1, const MpfiClass& op2);
```

```
MpfiClass& operator *= (MpfiClass& op1, const MpfiClass& op2);
MpfiClass& operator *= (MpfiClass& op1, const double& op2);
MpfiClass& operator *= (MpfiClass& op1, const cxsc::real& op2);
MpfiClass& operator *= (MpfiClass& op1, const cxsc::interval& op2);
MpfiClass& operator *= (MpfiClass& op1, const int op2);
MpfiClass& operator *= (MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass& operator *= (MpfiClass& op1, const mpfr_t& op2);
MpfiClass& operator *= (MpfiClass& op1, const mpfi_t& op2);
```

#### 4.6.4. Division

Beachten Sie den Hinweis auf Seite 51.

```
MpfiClass operator / (const MpfiClass& op1, const MpfiClass& op2);  
MpfiClass operator / (const MpfiClass& op1, const double& op2);  
MpfiClass operator / (const MpfiClass& op1, const cxsc::real& op2);  
MpfiClass operator / (const MpfiClass& op1, const cxsc::interval& op2);  
MpfiClass operator / (const MpfiClass& op1, const int op2);  
MpfiClass operator / (const MpfiClass& op1, const mpfr_t& op2);  
MpfiClass operator / (const MpfiClass& op1, const MPFR::MpfrClass& op2);  
MpfiClass operator / (const MpfiClass& op1, const mpfi_t& op2);  
MpfiClass operator / (const double& op1, const MpfiClass& op2);  
MpfiClass operator / (const cxsc::real& op1, const MpfiClass& op2);  
MpfiClass operator / (const cxsc::interval& op1, const MpfiClass& op2);  
MpfiClass operator / (const int op1, const MpfiClass& op2);  
MpfiClass operator / (const mpfr_t& op1, const MpfiClass& op2);  
MpfiClass operator / (const MPFR::MpfrClass& op1, const MpfiClass& op2);  
MpfiClass operator / (const mpfi_t& op1, const MpfiClass& op2);
```

```
MpfiClass& operator /= (MpfiClass& op1, const MpfiClass& op2);  
MpfiClass& operator /= (MpfiClass& op1, const double& op2);  
MpfiClass& operator /= (MpfiClass& op1, const cxsc::real& op2);  
MpfiClass& operator /= (MpfiClass& op1, const cxsc::interval& op2);  
MpfiClass& operator /= (MpfiClass& op1, const int op2);  
MpfiClass& operator /= (MpfiClass& op1, const MPFR::MpfrClass& op2);  
MpfiClass& operator /= (MpfiClass& op1, const mpfr_t& op2);  
MpfiClass& operator /= (MpfiClass& op1, const mpfi_t& op2);
```

## 4.7. Eingabe / Ausgabe

```
std::ostream& operator << (std::ostream& os, const MpfiClass& x);
```

Ermöglicht die Ausgabe einer MpfiClass-Variablen `x` über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(x.GetPrecision()/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Der Wert der Variablen `x` wird dabei in der mit `SetBase(k)` voreingestellten Basis ausgegeben. Für die meist dezimale Ausgabe ist natürlich `k = 10` zu wählen.

```
std::ostream& operator << (std::ostream& os, mpfi_t& x);
```

Ermöglicht die Ausgabe einer Variablen `x` vom Typ `mpfi_t` über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(mpfi_get_prec(x)/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Der Intervallwert `x` wird dabei in der mit `SetBase(k)` voreingestellten Basis ausgegeben. Für die meist dezimale Ausgabe ist `k = 10` zu wählen.

```
std::istream& operator >> (std::istream& is, MpfiClass& x);
```

Ermöglicht das Einlesen einer MpfiClass-Variablen über den Standard-Eingabestrom "cin". Das eingegebene Intervall ist auf keine Stellenanzahl begrenzt. Folgende Eingabeformate sind zulässig:

- `[-1.23e-4401,2.3E+2000]`      Vorsicht: Leerzeichen sind nicht erlaubt!
- `1.1`

Mit der letzten Eingabe entsteht ein nicht-punktförmiges Intervall, das die nicht-darstellbare Dezimalzahl 1.1 bezüglich der Current-Precision optimal einschließt. Die Zeichenkette muss mit der voreingestellten Basis übereinstimmen, sonst entsteht eine Fehlermeldung.

## 4.8. Base und Precision Handling

```
PrecisionType GetPrecision () const;
```

Diese Memberfunktion gibt die Präzision des aktuellen Objekts in Bits zurück. Als Beispiel entsprechen dabei 302 Bits  $302/\log_2(10) \approx 91$  Dezimalstellen.

```
void SetPrecision (PrecisionType prec);
```

Diese Memberfunktion setzt die Präzision des aktuellen Objektes auf `prec`. Sein Wert bleibt nicht erhalten.

```
void RoundPrecision (PrecisionType prec);
```

Diese Memberfunktion schließt das aktuelle Objekt mit der neuen Präzision `prec` ein. Sollte die Präzision des ursprünglichen Objektes größer sein als `prec`, so erhält man eine gröbere Einschließung. Ist die Präzision jedoch kleiner als `prec`, so werden die restlichen binären Stellen mit Nullen aufgefüllt, so dass die Werte der Intervallrandpunkte erhalten bleiben.

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision in `MpfrClass` auf `prec`. Diese Current-Precision wird dann automatisch auch in der Klasse `MpfiClass` benutzt. Wird die Current-Precision nicht gesetzt, so wird in beiden Klassen mit der Default-Precision von 53 Bits gerechnet. Das Setzen der Current-Precision hat auf die Präzision der bis dahin benutzten Variablen **keinerlei** Einfluss.

```
const int MpfrClass::GetBase ();
```

Gibt die aktuelle Basis zurück, diese hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die Funktion ist definiert in `mpfrclass.cpp`.

```
void MpfrClass::SetBase (int b);
```

Setzt die aktuelle Basis auf `b`. Dies muss ein Wert zwischen 2 und 36 sein. Die Basis hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die internen Rechnungen erfolgen stets im Binärsystem! Die Funktion ist definiert in `mpfrclass.cpp`.

## 4.9. Typ-Umwandlungen

Um ein flexibles Arbeiten zwischen C-XSC und der MPFI-Bibliothek zu ermöglichen, wurden möglichst viele Typ-Umwandlungs-Funktionen bereitgestellt:

### 4.9.1. real, double, ... → MPFI

```
MpfiClass real2Mpfi      (const cxsc::real& op);  
MpfiClass double2Mpfi   (const double& op);  
MpfiClass int2Mpfi      (const int& op);  
MpfiClass MpfrClass2Mpfi (const MPFR::MpfrClass& op);  
MpfiClass mpfr_t2Mpfi   (const mpfr_t& op);  
MpfiClass mpfi_t2Mpfi   (const mpfi_t& op);  
MpfiClass interval2Mpfi (const cxsc::interval& op);
```

Obige Funktionen liefern mit dem jeweiligen Eingabewert `op` einen Rückgabewert vom Typ `MpfiClass` in einer Präzision, die gewährleistet, dass der Rückgabewert **genau** dem Wert von `op` entspricht. Die Präzision des Rückgabewertes wird also i.a. nicht mit der Current-Precision übereinstimmen! Die obigen sieben Funktionen kommen u.a. bei den Vergleichsoperatoren zur Anwendung.

### 4.9.2. MPFI → interval

Die folgende Funktion liefert mit einem Objekt `op` vom Typ `MpfiClass` eine i.a. **gerundete** Einschließung von `op` vom C-XSC Typ `interval`;

```
cxsc::interval to_interval(const MpfiClass& op);
```

Eine optimale Einschließung von `op` wird nur erreicht, wenn `op` im IEEE-System darstellbar ist.

### 4.9.3. MPFI → mpfi\_t

```
const mpfi_t& getvalue(const MpfiClass& v)
```

Obige Funktion liefert von einem als `const` deklarierten Objekt `v` eine Referenz auf den Wert seines Attributs `mpfi_rep` vom Typ `mpfi_t`.

**Anwendung:** `const`-Parameter in einer Konstruktor-Parameterliste, vgl. z.B. den Konstruktor `MpfiClass::MpfiClass (const MpfiClass& x, PrecisionType prec)` in der Datei `mpficlass.cpp`.

### 4.9.4. MPFI ↔ mpfi\_t

```
mpfi_t& MpfiClass::GetValue();
```

Mithilfe der Memberfunktion `GetValue()`, die eine Referenz auf den Wert `mpfi_rep` vom Typ `mpfi_t` des aktuellen Objekts liefert, können sowohl `MpfiClass`-Objekte an eine Funktion übergeben als auch referenzierte Rückgabewerte vom Typ `mpfi_t` von einer Funktion übernommen werden. Mithilfe von `GetValue()` kann man daher Funktionen mit referenzierten `mpfi_t`-Parametern einfach aufrufen, vergleiche dazu das Programm `MPFR-02` auf Seite 30, in dem die analoge Übergabe an referenzierte Parameter vom Typ `mpfr_t` gezeigt wird.



#### 4.9.5. mpfi\_t → MPFI

```
void SetValue(const mpfi_t& t);
```

Mithilfe der Memberfunktion `SetValue()` wird der `mpfi_rep`-Wert des aktuellen Objekts auf `t` gesetzt, wobei `mpfi_rep` die Präzision von `t` übernimmt, d.h. der Wert von `t` wird ohne Rundung exakt übernommen.

#### 4.9.6. string → MPFI

```
MpfiClass string2Mpfi(const std::string& op, PrecisionType prec);
```

Der Aufruf `string2Mpfi(op);` rundet den String `op` mittels der voreingestellten Current-Precision in ein Klassenobjekt vom Typ `MpfiClass`. Der zweite mögliche Aufruf `string2Mpfi(op,140);` rundet `op` ebenfalls in ein Klassenobjekt vom Typ `MpfiClass` der Präzision 140 Bits. Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in ein binäres Format konvertiert werden kann, so dass i.a. Rundungen nicht zu vermeiden sind, es werden jedoch stets Einschließungen von `op` zurückgegeben. Eine weitere Möglichkeit, einen String in ein `MpfiClass`-Objekt zu verwandeln, besteht in einem entsprechenden Konstruktor-Aufruf, vgl. Seite 49. Dort findet man auch Hinweise auf mögliche Intervall-String-Formate.

#### 4.9.7. MPFI → string

```
std::string to_string (const MpfiClass& x, PrecisionType prec);
```

`x` wird in einen String `s` mit `prec` Dezimalstellen gerundet, wenn Base gleich 10 ist, dabei ist der String stets eine optimale Einschließung von `x`. Wählt man `prec` hinreichend groß, so stellt der String den Wert von `x` sogar **exakt** dar, weil eine binäre Zahl stets exakt in eine Dezimalzahl umgewandelt werden kann.

Wird `prec` nicht angegeben, so wird `x` mit einem String eingeschlossen, der im Fall Base=10 so viele Dezimalstellen besitzt, wie es der Präzision von `x` entspricht. Ist die Präzision von `x` z.B. 302, so wird ein String von  $302/\log_2(10) \approx 302/3.321928095 \approx 91$  Dezimalstellen generiert.

#### 4.9.8. MPFI → MPFI

Beachten Sie, dass bei einer Wertzuweisung an eine `MpfiClass`-Variable mit dem Operator `=` der linke Operand stets auf die Current-Precision gesetzt wird und dass der rechte `MpfiClass`-Operand dabei vom linken Operanden **stets** optimal eingeschlossen wird, vgl. dazu auch Seite 50. Will man jedoch, dass der linke einschließende Operand auch eine andere Präzision `prec` erhält, so kann dies mit folgender Funktion realisiert werden:

```
void set_Mpfi (MpfiClass& op, const MpfiClass& op1, PrecisionType prec);
```

Folgende Funktionsaufrufe sind möglich:

1. `set_Mpfi (op, op1, prec);`  
`op` erhält die Präzision `prec` und den i.a. gerundeten Wert von `op1`, wobei `op1` von `op` stets optimal eingeschlossen wird. Setzt man `prec` gleich der Präzision von `op1`, so erhält `op` den **exakten** Wert von `op1`.
2. `set_Mpfi (op, op1);`  
`op` wird auf die Current-Precision gesetzt und erhält den i.a. gerundeten Wert von `op1`, wobei auch hier `op1` von `op` stets optimal eingeschlossen wird.

## 4.10. Abfragen

Bei allen folgenden Abfragefunktionen braucht die Präzision von `x` nicht mit der Current-Precision übereinzustimmen.

```
bool isNaN (const MpfiClass& x);  
bool isInf (const MpfiClass& x);  
bool isPoint (const MpfiClass& x);
```

`isNaN` und `isInf` überprüfen, ob ein Randpunkt von `x` gleich NaN bzw.  $\pm\text{Inf}$  ist.  
`isPoint` überprüft, ob `x` ein Punktintervall ist.

```
bool isBounded (const MpfiClass& x);
```

`isBounded` überprüft, ob `x` ein normales Intervall ist, d.h. kein Randpunkt ist NaN oder  $\pm\text{Inf}$ .

```
bool isZero (const MpfiClass& x);
```

Überprüft, ob `x` das Null-Intervall ist.

```
bool hasZero (const MpfiClass& x);
```

Überprüft, ob `x` die Null enthält, nicht notwendig im Innern von `x`.

```
bool isPos (const MpfiClass& x)
```

Überprüft, ob die Elemente von `x` größer oder gleich Null sind, d.h. der linke Randpunkt kann Null sein.

```
bool isStrictlyPos(const MpfiClass& x)
```

Überprüft, ob alle Elemente von `x` positiv sind.

```
bool isNonNeg(const MpfiClass& x);
```

Überprüft, ob alle Elemente von `x` nicht-negativ sind.

```
bool isNeg(const MpfiClass& x)
```

Überprüft, ob alle Elemente von `x` negativ sind, wobei der rechte Randpunkt Null sein kann.

```
bool isStrictlyNeg(const MpfiClass& x)
```

Überprüft, ob alle Elemente von `x` negativ sind.

```
bool isNonPos(const MpfiClass& x)
```

Überprüft, ob alle Elemente von `x` kleiner oder gleich Null sind, d.h. der rechte Randpunkt kann Null sein.

```
bool isEmpty ();
```

Die obige Member-Funktion liefert True, falls die Randpunkte des aktuellen Objekts in falscher Ordnung vorliegen und False sonst.

```
bool Disjoint(const MpfiClass& x, const MpfiClass& y);
```

Die Funktion liefert True, falls `x` und `y` keine gemeinsamen Elemente besitzen und False sonst.

## 4.11. Vergleiche

Mit den folgenden Vergleichsfunktionen werden die üblichen Vergleichsoperatoren implementiert.

### 4.11.1. Vergleichsfunktionen

```
bool compare_equal (const MpfiClass& x, const MpfiClass& y);
bool compare_less (const MpfiClass& x, const MpfiClass& y);
bool compare_lessequal (const MpfiClass& x, const MpfiClass& y);
```

Die obigen Funktionen überprüfen, ob jeweils gilt:

$x = y$ ,  $x < y$ ,  $x \leq y$ , Ist  $x$  oder  $y$  NaN oder Inf, so wird False zurückgegeben.

### 4.11.2. Vergleichsoperatoren =, ≠, >, ≥, <, ≤

```
bool operator == (const MpfiClass& op1, const MpfiClass& op2);
bool operator == (const MpfiClass& op1, const double& op2);
bool operator == (const MpfiClass& op1, const int op2);
bool operator == (const MpfiClass& op1, const cxsc::real& op2);
bool operator == (const MpfiClass& op1, const cxsc::interval& op2);
bool operator == (const MpfiClass& op1, const MPFR::MpfrClass& op2);
bool operator == (const MpfiClass& op1, const mpfr_t& op2);
bool operator == (const MpfiClass& op1, const mpfi_t& op2);
bool operator == (const double& op1, const MpfiClass& op2);
bool operator == (const int op1, const MpfiClass& op2);
bool operator == (const cxsc::real& op1, const MpfiClass& op2);
bool operator == (const mpfr_t& op1, const MpfiClass& op2);
bool operator == (const mpfi_t& op1, const MpfiClass& op2);
bool operator == (const MPFR::MpfrClass& op1, const MpfiClass& op2);
bool operator == (const cxsc::interval& op1, const MpfiClass& op2);
```

```
bool operator != (const MpfiClass& op1, const MpfiClass& op2);
bool operator != (const MpfiClass& op1, const double& op2);
bool operator != (const MpfiClass& op1, const int op2);
bool operator != (const MpfiClass& op1, const cxsc::real& op2);
bool operator != (const MpfiClass& op1, const cxsc::interval& op2);
bool operator != (const MpfiClass& op1, const MPFR::MpfrClass& op2);
bool operator != (const MpfiClass& op1, const mpfr_t& op2);
bool operator != (const MpfiClass& op1, const mpfi_t& op2);
bool operator != (const double& op1, const MpfiClass& op2);
bool operator != (const int op1, const MpfiClass& op2);
bool operator != (const cxsc::real& op1, const MpfiClass& op2);
bool operator != (const mpfr_t& op1, const MpfiClass& op2);
bool operator != (const mpfi_t& op1, const MpfiClass& op2);
bool operator != (const MPFR::MpfrClass& op1, const MpfiClass& op2);
bool operator != (const cxsc::interval& op1, const MpfiClass& op2);
```

```

bool operator < (const MpfiClass& op1, const MpfiClass& op2);
bool operator < (const MpfiClass& op1, const mpfi_t& op2);
bool operator < (const MpfiClass& op1, const cxsc::interval& op2);
bool operator < (const mpfi_t& op1, const MpfiClass& op2);
bool operator < (const cxsc::interval& op1, const MpfiClass& op2);

```

'interval' < 'MpfiClass': 'interval' liegt ganz im Innern von 'MpfiClass'

```

bool operator < (const MPFR::MpfrClass& op1, const MpfiClass& op2);
bool operator < (const mpfr_t& op1, const MpfiClass& op2);
bool operator < (const double& op1, const MpfiClass& op2);
bool operator < (const cxsc::real& op1, const MpfiClass& op2);
bool operator < (const int op1, const MpfiClass& op2);

```

'real' < 'MpfiClass': 'real' liegt ganz im Innern von 'MpfiClass'

```

bool operator <= (const MpfiClass& op1, const MpfiClass& op2);
bool operator <= (const MpfiClass& op1, const mpfi_t& op2);
bool operator <= (const MpfiClass& op1, const cxsc::interval& op2);
bool operator <= (const mpfi_t& op1, const MpfiClass& op2);
bool operator <= (const cxsc::interval& op1, const MpfiClass& op2);

```

'interval' <= 'MpfiClass': 'interval'  $\subseteq$  'MpfiClass'

```

bool operator <= (const MpfiClass& op1, const MPFR::MpfrClass& op2);
bool operator <= (const MpfiClass& op1, const mpfr_t& op2);
bool operator <= (const MpfiClass& op1, const double& op2);
bool operator <= (const MpfiClass& op1, const cxsc::real& op2);
bool operator <= (const MpfiClass& op1, const int op2);

```

'MpfiClass' <= 'real': Nur wahr, wenn Punktintervall 'MpfiClass' = 'real'

```

bool operator <= (const MPFR::MpfrClass& op1, const MpfiClass& op2);
bool operator <= (const mpfr_t& op1, const MpfiClass& op2);
bool operator <= (const double& op1, const MpfiClass& op2);
bool operator <= (const cxsc::real& op1, const MpfiClass& op2);
bool operator <= (const int op1, const MpfiClass& op2);

```

'real' <= 'MpfiClass': 'real'  $\in$  'MpfiClass'

```

bool operator > (const MpfiClass& op1, const MpfiClass& op2);
bool operator > (const MpfiClass& op1, const mpfi_t& op2);
bool operator > (const MpfiClass& op1, const cxsc::interval& op2);
bool operator > (const mpfi_t& op1, const MpfiClass& op2);
bool operator > (const cxsc::interval& op1, const MpfiClass& op2);

```

'interval' > 'MpfiClass': 'MpfiClass' liegt ganz im Innern von 'interval'

```

bool operator > (const MpfiClass&, const double&);
bool operator > (const MpfiClass&, const int);
bool operator > (const MpfiClass&, const cxsc::real&);
bool operator > (const MpfiClass&, const mpfr_t&);
bool operator > (const MpfiClass&, const MPFR::MpfrClass&);

```

'MpfiClass' > 'real': 'real' liegt ganz im Innern von 'MpfiClass'

```

bool operator >= (const MpfiClass& op1, const MpfiClass& op2);
bool operator >= (const MpfiClass& op1, const mpfi_t& op2);
bool operator >= (const MpfiClass& op1, const cxsc::interval& op2);
bool operator >= (const mpfi_t& op1, const MpfiClass& op2);
bool operator >= (const cxsc::interval& op1, const MpfiClass& op2);

```

'interval' >= 'MpfiClass': 'interval'  $\supseteq$  'MpfiClass'

```

bool operator >= (const MpfiClass& op1, const MPFR::MpfrClass& op2);
bool operator >= (const MpfiClass& op1, const mpfr_t& op2);
bool operator >= (const MpfiClass& op1, const double& op2);
bool operator >= (const MpfiClass& op1, const cxsc::real& op2);
bool operator >= (const MpfiClass& op1, const int op2);

```

'MpfiClass' >= 'real': 'MpfiClass'  $\ni$  'real'

```

bool operator >= (const MPFR::MpfrClass& op1, const MpfiClass& op2);
bool operator >= (const mpfr_t& op1, const MpfiClass& op2);
bool operator >= (const double& op1, const MpfiClass& op2);
bool operator >= (const cxsc::real& op1, const MpfiClass& op2);
bool operator >= (const int op1, const MpfiClass& op2);

```

'real' >= 'MpfiClass': Nur wahr, wenn Punktintervall 'MpfiClass' = 'real'

```

int in (const MPFR::MpfrClass::MpfrClass& x, const MpfiClass& y);
int in (const mpfr_t& x, const MpfiClass& y);
int in (const double& x, const MpfiClass& y);
int in (const cxsc::real& x, const MpfiClass& y);
int in (const int& x, const MpfiClass& y);
int in (const MpfiClass& x, const MpfiClass& y);
int in (const mpfi_t& x, const MpfiClass& y);
int in (const cxsc::interval& x, const MpfiClass& y);

```

Zurückgegeben wird die Eins, wenn  $x$  ganz im Innern von  $y$  enthalten ist, sonst wird die Null zurückgegeben. Ist einer der Operanden ein NaN oder sind beide Operanden unbegrenzt, so wird ebenfalls die Null zurückgegeben. Die Präzisionen beider Operanden können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen.

#### Hinweise:

- Alle Vergleichsoperatoren sind wie in C-XSC definiert.
- Die Präzisionen der Operanden können verschieden sein und müssen mit der Current-Präzision nicht übereinstimmen.

## 4.12. Durchschnitt

Berechnet wird der Durchschnitt zweier Intervalle. Ist einer der nachfolgenden Operanden eine Zahl, so ist diese als Punktintervall zu verstehen. Die Präzisionen der Intervalle oder Zahlen können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen. Die Präzision des zurückgegebenen Durchschnitt-Intervalls ist das Maximum der Präzisionen beider Operanden und muss daher mit der Current-Präzision **nicht** übereinstimmen. Dadurch wird der Durchschnitt **rundungsfehlerfrei** berechnet. Ist der Durchschnitt leer, so wird NaN zurückgegeben.

```
MpfiClass operator & (const MpfiClass& op1, const MpfiClass& op2);
MpfiClass operator & (const MpfiClass& op1, const mpfi_t& op2);
MpfiClass operator & (const MpfiClass& op1, const cxsc::interval& op2);
MpfiClass operator & (const mpfi_t& op1, const MpfiClass& op2);
MpfiClass operator & (const cxsc::interval& op1, const MpfiClass& op2);
```

```
MpfiClass operator & (const MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass operator & (const MpfiClass& op1, const mpfr_t& op2);
MpfiClass operator & (const MpfiClass& op1, const double& op2);
MpfiClass operator & (const MpfiClass& op1, const cxsc::real& op2);
MpfiClass operator & (const MpfiClass& op1, int op2);
```

```
MpfiClass operator & (const MPFR::MpfrClass& op1, const MpfiClass& op2);
MpfiClass operator & (const mpfr_t& op1, const MpfiClass& op2);
MpfiClass operator & (const double& op1, const MpfiClass& op2);
MpfiClass operator & (const cxsc::real& op1, const MpfiClass& op2);
MpfiClass operator & (int op1, const MpfiClass& op2);
```

```
MpfiClass& operator &= (MpfiClass& op1, const MpfiClass& op2);
MpfiClass& operator &= (MpfiClass& op1, const mpfi_t& op2);
MpfiClass& operator &= (MpfiClass& op1, const cxsc::interval& op2);
MpfiClass& operator &= (MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass& operator &= (MpfiClass& op1, const mpfr_t& op2);
MpfiClass& operator &= (MpfiClass& op1, const double& op2);
MpfiClass& operator &= (MpfiClass& op1, const cxsc::real& op2);
MpfiClass& operator &= (MpfiClass& op1, int op2);
```

Die Anweisung `op1 &= op2;` liefert an `op1` den Durchschnitt (`op1 & op2`), wobei die neue Präzision von `op1` gleich dem Maximum der Präzisionen der ursprünglichen Operanden `op1` und `op2` ist. Auch hier wird wieder erreicht, dass der Durchschnitt **rundungsfehlerfrei** an `op1` zurückgegeben wird.

### Hinweis:

Wenn beide Operanden als Präzision die Current-Präzision besitzen, so wird der Durchschnitt stets **rundungsfehlerfrei** in der gleichen Präzision zurückgegeben.

## 4.13. Konvexe Hülle

Berechnet wird die konvexe Hülle zweier Intervalle. Ist einer der nachfolgenden Operanden eine Zahl, so ist diese als Punktintervall zu verstehen. Die Präzisionen der Intervalle oder Zahlen können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen. Die Präzision des zurückgegebenen Hüllen-Intervalls ist das Maximum der Präzisionen beider Operanden und muss mit der Current-Präzision **nicht** übereinstimmen. Dadurch wird die konvexe Hülle **rundungsfehlerfrei** berechnet.

```
MpfiClass operator | (const MpfiClass& op1, const MpfiClass& op2);
MpfiClass operator | (const double& op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, const double& op2);
MpfiClass operator | (const cxsc::real& op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, const cxsc::real& op2);
MpfiClass operator | (const cxsc::interval& op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, const cxsc::interval& op2);
MpfiClass operator | (int op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, int op2);
MpfiClass operator | (const MpfiClass& op1, const mpfr_t& op2);
MpfiClass operator | (const mpfr_t& op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, const mpfi_t& op2);
MpfiClass operator | (const mpfi_t& op1, const MpfiClass& op2);
MpfiClass operator | (const MpfiClass& op1, const MPFR::MpfrClass& op2);
MpfiClass operator | (const MPFR::MpfrClass& op1, const MpfiClass& op2);
```

```
MpfiClass& operator |= (MpfiClass& op1, const MpfiClass& op2);
MpfiClass& operator |= (MpfiClass& op1, const double& op2);
MpfiClass& operator |= (MpfiClass& op1, const cxsc::real& op2);
MpfiClass& operator |= (MpfiClass& op1, const cxsc::interval& op2);
MpfiClass& operator |= (MpfiClass& op1, int op2);
MpfiClass& operator |= (MpfiClass& op1, const mpfr_t& op2);
MpfiClass& operator |= (MpfiClass& op1, const mpfi_t& op2);
MpfiClass& operator |= (MpfiClass& op1, const MPFR::MpfrClass& op2);
```

Die Anweisung `op1 |= op2;` liefert an `op1` die konvexe Hülle (`op1 | op2`), wobei die neue Präzision von `op1` gleich dem Maximum der Präzisionen der ursprünglichen Operanden `op1` und `op2` ist. Auch hier wird wieder erreicht, dass die konvexe Hülle **rundungsfehlerfrei** zurückgegeben wird.

### Hinweis:

Wenn beide Operanden als Präzision die Current-Präzision besitzen, so wird die konvexe Hülle stets **rundungsfehlerfrei** in der gleichen Präzision zurückgegeben.



## 4.14. Mathematische Intervall-Funktionen

Zu einer gegebenen Funktion  $f: \mathbb{R} \rightarrow \mathbb{R}$  liefert eine Intervall-Funktion  $f[x]$  zu einem vorgegebenen Intervall  $[x] = [a, b]$  die Menge aller Funktionswerte  $f(x)$ , mit  $x \in [x]$ :

$$f([x]) = \{y \in \mathbb{R} \mid y = f(x) \wedge x \in [x] = [a, b]\}.$$

Gilt dann für ein Maschinenintervall  $\mathbf{x} \supseteq [x]$ , so liefert eine implementierte Intervallfunktion, z.B.  $\mathbf{exp}(\mathbf{x})$ , eine garantierte Einschließung aller Funktionswerte  $y = f(x) = e^x$ , mit  $x \in [x] = [a, b] \subseteq \mathbf{x}$ .

$$\{y \in \mathbb{R} \mid y = e^x \wedge x \in \mathbf{x}\} \subseteq \mathbf{exp}(\mathbf{x}).$$

### 4.14.1. Standard-Implementierung

Die MPFI-Bibliothek stellt eine Vielzahl von Elementarfunktionen als Intervallfunktionen zur Verfügung, wobei zu einem Argument  $\mathbf{x}$  mit beliebiger Präzision zunächst das exakte Funktionswert-Intervall  $y_0$  berechnet wird. Danach wird dann  $y_0$  durch ein Ergebnisintervall  $y \supseteq y_0$  mit einer möglichen anderen Präzision eingeschlossen. Die Deklaration z.B. der Exponentialfunktion ist mit diesen Bezeichnungen gegeben durch:

```
int mpfi_exp (mpfi_t y, mpfr_t x);
```

Berechnet eine Funktion der MPFI-Bibliothek die Randpunkte einer Einschließung ohne Über- oder Unterlauf, so ist die berechnete Einschließung maximal-genau, d.h. die exakten Randpunkte unterscheiden sich von den berechneten auf- bzw. abgerundeten Randpunkten um weniger als ein ulp.

Mit obiger MPFI-Funktion  $\mathbf{mpfi\_exp}(\dots)$  wird die Exponentialfunktion für das C-XSC Interface wie folgt implementiert:

```
MpfiClass exp(const MpfiClass& x)
{
    MpfiClass y(0);
    mpfi_exp(res.mpfi_rep, v.mpfi_rep);
    return res;
}
```

Mit Hilfe des Konstruktoraufrufs wird also zunächst die Präzision der Ergebnisvariablen  $y$  auf die Current-Precision gesetzt, vgl. dazu auch Seite 49. Danach wird mit dem Argument  $\mathbf{x}$  und seiner Präzision das exakte Intervall  $y_0 = e^x$  berechnet und dann durch  $y$  optimal eingeschlossen.

Bei allen im MpfiClass-Interface implementierten mathematischen Intervallfunktionen wird das berechnete Einschließungsintervall unabhängig von der Präzision des Arguments in der Current-Präzision zurückgegeben.

Es gibt jedoch einige Intervallfunktionen, bei denen es sinnvoll ist, das Ergebnis mit der gleichen Präzision des Arguments zurückzugeben. So wird man beispielsweise verlangen, dass die Randpunkte eines Intervalls  $\mathbf{x}$  mit der gleichen Präzision von  $\mathbf{x}$  berechnet werden. Dies erreicht man z.B. mit den Funktionen  $\mathbf{GetLeft}(\dots)$  bzw.  $\mathbf{GetRight}(\dots)$ . Mit den Funktionen  $\mathbf{Inf}(\dots)$  bzw.  $\mathbf{Sup}(\dots)$  werden die Randpunkte jedoch in die Current-Präzision gerundet.

Im folgenden Abschnitt werden diejenigen Intervallfunktionen beschrieben, deren Rückgabewerte Präzisionen besitzen, die nicht mit der Current-Präzision übereinstimmen müssen.



#### 4.14.2. Davon abweichende Funktionen und Konstanten

Bei den folgenden Funktionen kann die Präzision des Rückgabewertes von der Current-Präzision abweichen.

```
MPFR::MpfrClass diam(const MpfiClass& op);
```

Liefert den aufgerundeten absoluten Durchmesser von `op`, wobei der Rückgabewert die Präzision von `op` erhält.

```
MPFR::MpfrClass RelDiam(const MpfiClass& op);
```

Liefert den aufgerundeten absoluten Durchmesser von `op`, falls Null in `op` enthalten ist. Sonst wird der aufgerundete relative Durchmesser von `op` berechnet. In beiden Fällen erhält der Rückgabewert die Präzision von `op`.

```
MPFR::MpfrClass AbsMax(const MpfiClass& op);
```

```
MPFR::MpfrClass AbsMin(const MpfiClass& op);
```

Zurückgegeben wird das Maximum bzw. das Minimum aller Absolutbeträge von `op`, wobei der Rückgabewert die Präzision von `op` erhält.

```
MPFR::MpfrClass mid(const MpfiClass& op);
```

Zurückgegeben wird der Mittelpunkt von `op`, wobei der Rückgabewert die Präzision von `op` erhält.

Die folgenden sechs Funktionen haben die jeweils gleiche Bedeutung wie oben, wobei aber jetzt `op2` der Rückgabewert mit der Präzision von `op1` ist.

```
void diam (const MpfiClass& op1, mpfr_t& op2);
```

```
void RelDiam (const MpfiClass& op1, mpfr_t& op2);
```

```
void AbsMax (const MpfiClass& op1, mpfr_t& op2);
```

```
void AbsMin (const MpfiClass& op1, mpfr_t& op2);
```

```
void mid (const MpfiClass& op1, mpfr_t& op2);
```

```
MPFR::MpfrClass Inf(const MpfiClass& op, PrecisionType prec);
```

```
MPFR::MpfrClass Sup(const MpfiClass& op, PrecisionType prec);
```

Zurückgegeben wird der linke bzw. rechte Randpunkt von `op`, jeweils gerundet in ein Format mit der Präzision `prec`. Ohne Angabe von `prec` wird dabei nach  $-\infty$  bzw. nach  $+\infty$  in die Current-Präzision gerundet. Mit `Inf(op, op.GetPrecision())` wird der linke Randpunkt **rundungsfehlerfrei** mit der Präzision von `op` zurückgegeben. Mit `Sup(op)` erhält man den rechten Randpunkt in der Current-Präzision, wobei im Bedarfsfall nach  $+\infty$  gerundet wird.

```
void GetLeft (mpfr_t& op);
```

```
void GetRight(mpfr_t& op);
```

Die obigen Memberfunktionen liefern vom aktuellen Objekt den linken bzw. rechten Randpunkt mit der gleichen Präzision des aktuellen Objekts.

```
MpfiClass Blow(const MpfiClass& op1, const MPFR::MpfrClass& op2);
```

Der Rückgabewert ist ein mit `op2` aufgeblähtes Intervall `op1` mit gleicher Präzision. `Blow(...)` ist analog zur gleichnamigen Funktion in C-XSC definiert.

```
int swap_endpoints ();
```

Die obige Memberfunktion tauscht am aktuellen Objekt beide Randpunkt, falls der linke Randpunkt größer ist als der rechte. Der Rückgabewert vom Typ `int` ist positiv, wenn ein Tausch notwendig war, sonst ist der Wert gleich Null.

```
void swap(MpfiClass& x, MpfiClass& y);
```

Tauscht den Wert und die Präzision von `x` und `y`.

```
void swap(MpfiClass& x, mpfi_t& y);
```

Tauscht den Wert und die Präzision von `x` und `y`.

```
int common_decimals(const MpfiClass& op);
```

Liefert die Anzahl der übereinstimmenden Dezimalziffern, in denen die Randpunkte von `op` übereinstimmen.

```
MPFR::MpfiClass random(const MpfiClass& op);
```

Zurückgegeben wird eine Zufallszahl aus `op`, wobei der Rückgabewert die Präzision von `op` erhält.

```
void random(MpfiClass& x, gmp_randstate_t state);
```

Zurückgegeben wird ein Zufallsintervall  $x \subseteq [0, 1]$ , wobei `x` die Präzision des vorher deklarierten Objekts `x` erhält. Das nachfolgende Programm zeigt eine mögliche Anwendung der Funktion `void random(MpfiClass& x, gmp_randstate_t state)`.

```
1 // MPFR-04.cpp
2 #include "mpficlass.hpp"
3
4 using namespace MPFI;
5 using namespace std;
6
7 int main(void)
8 {
9     MpfiClass::SetCurrRndMode (RoundNearest);
10    cout << "\nCurrent-RoundingMode = " << MpfiClass::GetCurrRndMode() << endl;
11    MpfiClass::SetCurrPrecision (60);
12    cout << "Current-Precision = " << MpfiClass::GetCurrPrecision() << endl;
13
14    gmp_randstate_t state; // Declaration of state;
15    gmp_randinit_default (state); // Initialization of state;
16
17    MpfiClass x(0,50); // Declaration of interval class object x
18                    // with a precision of 50 bits
19    cout.precision(x.GetPrecision()/3.321928095);
20    cout << "x = " << x << endl;
21    cout << "x.GetPrecision() = " << x.GetPrecision() << endl;
22    random (x, state); // Delivers the first random interval x
23    cout << "x = " << x << endl;
24    random (x, state); // Delivers the second random interval x
25    cout << "x = " << x << endl;
26
27    return 0;
28 }
```

Das Programm liefert die Ausgabe:

```
Current-RoundingMode = 0
Current-Precision = 60
x = [0, -0]
x.GetPrecision() = 50
x = [6.14544775142451e-1, 9.88050499009929e-1]
x = [2.04581622647135e-1, 4.38359857288966e-1]
```

Weitere Informationen bez. der Initialisierungsfunktion in Zeile 15 findet man unter

<http://gmplib.org/manual/Random-State-Initialization.html>

```
void times2pown (MpfClass& op, long int op1);
```

Obige Funktion liefert mit dem Eingabewert `op` eine optimale Einschließung von  $op \cdot 2^{op1}$  mit gleicher Präzision zurück. Solange kein Über- oder Unterlauf entsteht, wird  $op \cdot 2^{op1}$  **exakt**, d.h. rundungsfehlerfrei eingeschlossen.

```
void set_nan (MpfClass& x);
```

Setzt `x` auf NaN, wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_inf (MpfClass& x);
```

Setzt `x` auf  $[-\text{Inf}, +\text{Inf}]$ , wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_zero (MpfClass& x);
```

Setzt `x` auf  $[0, 0]$ , wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void MpfClass::SetInterval(const MPFR::MpfrClass& a,  
                           const MPFR::MpfrClass& b);  
void MpfClass::SetInterval(const mpfr_t&      a, const mpfr_t&      b);  
void MpfClass::SetInterval(const double&      a, const double&      b);  
void MpfClass::SetInterval(const cxsc::real&   a, const cxsc::real&   b);  
void MpfClass::SetInterval(const int&         a, const int&         b);
```

Die obigen Memberfunktionen liefern im Fall  $a \leq b$  für das jeweils aktuelle Objekt stets eine Einschließung von  $[a, b]$ , wobei die Präzision des aktuellen Objekts erhalten bleibt. Die Reihenfolge der Randpunkte `a` und `b` kann aber beliebig gewählt werden.

```
static MpfClass Pi      (PrecisionType prec = CurrPrecision);  
static MpfClass Ln2    (PrecisionType prec = CurrPrecision);  
static MpfClass Euler (PrecisionType prec = CurrPrecision);  
static MpfClass Catalan (PrecisionType prec = CurrPrecision);
```

`Pi(prec)` schließt  $\pi$  mit der Präzision `prec` ein. Wird `prec` nicht angegeben, so wird  $\pi$  mittels der Current-Precision eingeschlossen. Entsprechendes gilt für die drei anderen Konstanten. `Ln2()` schließt also  $\ln(2) = 0.693147\dots$  mittels der voreingestellten Current-Precision ein. Vergessen Sie nicht, beim Aufruf der Konstanten `Ln2()` die Klammern zu setzen.

### 4.14.3. Elementarfunktionen

Tabelle 4.1.: Elementarfunktionen mit  $x, y$  vom Typ `MpfiClass`,  $a$  vom Typ `MpfrClass`;

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\ln(x)$	<code>ln(x)</code>
$x^2$	<code>sqr(x)</code>	$\ln(1+x)$	<code>lnp1(x)</code>
$x^2 + y^2$	<code>x2py2(x,y)</code>	$\log_2(x)$	<code>log2(x)</code>
$x^2 - y^2$	<code>x2my2(x,y)</code>	$\log_{10}(x)$	<code>log10(x)</code>
$1/x$	<code>reci(x)</code>	$\ln(\sin(x))$	<code>ln_sin(x)</code>
$x/(x^2 + y^2)$	<code>x_div_x2py2(x,y)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$(x^2 - y^2)/(x^2 + y^2)^2$	<code>Re_rz2(x,y)</code>	$\ln(\sqrt{x^2 + y^2})$	<code>ln_sqrtx2y2(x,y)</code>
$2xy/(x^2 + y^2)^2$	<code>mIm_rz2(x,y)</code>	$\ln(\sqrt{(1+x)^2 + y^2})$	<code>ln_sqrtxp1_2y2(x,y)</code>
$\frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Re_r1pz2(x,y)</code>	$x^k, k \in \mathbb{Z}$	<code>power(x,k)</code>
$\frac{2xy}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Im_r1pz2(x,y)</code>	$x^y$	<code>pow(x,y)</code>
$\sqrt{x}$	<code>sqr(x)</code>	$\sin(x)$	<code>sin(x)</code>
$\sqrt{n}, n \in \mathbb{N}_0$	<code>sqrtn(n)</code>	$1/\sin(x)$	<code>csc(x)</code>
$a/\sqrt{n}, n \in \mathbb{N}$	<code>xdivsqrtn(a,n)</code>	$\cos(x)$	<code>cos(x)</code>
$x/\sqrt{n}, n \in \mathbb{N}$	<code>xdivsqrtn(x,n)</code>	$1/\cos(x)$	<code>sec(x)</code>
$1/\sqrt{x}$	<code>sqrtr(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt[3]{x}$	<code>cbrt(x)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqrtn(x,n)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt{x+1} - 1$	<code>sqrtp1m1(x)</code>	$\arccos(x)$	<code>acos(x)</code>
$\sqrt{1+x^2}$	<code>sqrtp1px2(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{1-x^2}$	<code>sqrtp1mx2(x)</code>	$\arg(x + i \cdot y)$	<code>atan2(y,x)</code>
$\sqrt{x^2-1}$	<code>sqrtpx2m1(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$x/\sqrt{1+x^2}$	<code>xdsqrtp1px2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>

Fortsetzung auf der nächsten Seite

Fortsetzung der vorherigen Seite			
Funktion	Aufruf	Funktion	Aufruf
$x/\sqrt{1-x^2}$	<code>xdsqrt1mx2(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$x/\sqrt{x^2-1}$	<code>xdsqrtx2m1(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\sqrt{x^2+y^2}$	<code>sqrtox2y2(x,y)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$e^x$	<code>exp(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$2^x$	<code>exp2(x)</code>	$\coth(x)$	<code>coth(x)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{AGM}(x,y)$	<code>agm(x,y)</code>

### Anmerkungen:

1. Mit `sqr(x)`; erhält man eine Einschließung aller  $x^2$ -Werte für  $x \in \mathbf{x}$ , wobei in die Current-Precision gerundet wird. Weitere Informationen zur Implementierung findet man auf Seite 64. Beachten Sie, dass z.B. für  $x = [-4, 1]$  die normale Multiplikation  $x \diamond x = [-4, +16]$  und `sqr(x) = [0, +16]` **verschiedene** Ergebnisse liefern!
2. Alle Funktionen der Tabelle auf Seite 68 sind nach Punkt 2. implementiert.
3. `abs(x)` kann zusätzlich noch mit einem Präzisionsparameter `prec` aufgerufen werden. Mit `abs(x, x.GetPrecision())` erhält man dabei eine **rundungsfehlerfreie** Einschließung aller  $|x|$ , mit  $x \in \mathbf{x}$ .
4. `power(x, k)` kann bez. `k` mit folgenden Datentypen aufgerufen werden: `int`, `long int`.
5. `pow(x, y)` kann bez. `y` mit folgenden Datentypen aufgerufen werden: `MpfiClass`, `interval`, `MpfrClass`, `real`, `double`.
6. Mit dem Aufruf `agm(X,Y)` erhält man in Current-Präzision eine optimale Einschließung aller  $\operatorname{AGM}(x,y)$ -Werte für  $x \in X$  und  $y \in Y$ . Das Arithmetisch-geometrische Mittel AGM spielt bei der Auswertung elliptischer Integrale eine zentrale Rolle.
7. Mit `u = xdivsqrtn(a,n)` wird  $a/\sqrt{n}$  durch `u` in der Current-Präzision **optimal** eingeschlossen. Dabei ist `a` vom Typ `MpfrClass` und `n` vom Typ `unsigned long int`.
8. Mit `u = xdivsqrtn(x,n)` wird  $x/\sqrt{n}$  durch `u` in der Current-Präzision **optimal** eingeschlossen. Dabei ist `x` vom Typ `MpfiClass` und `n` vom Typ `unsigned long int`.
9. Das Argument `n` der Funktion `sqrt_N(n)` ist vom Typ `unsigned long int`.

10. Die Funktion `MpfiClass atan2(const MpfiClass& Y, const MmpfiClass& X)` berechnet eine Einschließung der Argumente  $\arg(x+i*y)$  für alle  $x \in X$  und für alle  $y \in Y$ . Falls das Intervall  $(X, Y)$  die negative reelle Achse schneidet, wird das Intervall  $[-\pi, +\pi]$  zurückgegeben.

Die Funktion `MpfiClass ATAN2(const MpfiClass& Y, const MmpfiClass& X)` berechnet ebenfalls eine Einschließung der Argumente  $\arg(x+i*y)$  für alle  $x \in X$  und  $y \in Y$ . Falls das Intervall  $(X, Y)$  jedoch die negative reelle Achse von unten berührt oder schneidet, erfolgt eine entsprechende Fehlermeldung mit Programmabbruch.

#### 4.14.4. Funktionen der Mathematischen Physik

Tabelle 4.2.: Funktionen der Mathematischen Physik mit  $x$  vom Typ `MpfiClass`

Funktionsterm	Aufruf	Anmerkung
$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>	monoton wachsend
$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>	monoton fallend
$\Gamma'(x)$	* <code>gamma_D(x)</code>	Pole: $x=0, -1, -2, \dots$
$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$	<code>digamma(x)</code>	Pole: $x=0, -1, -2, \dots$
$k!$	<code>ifactorial(k)</code>	unsigned long int k
$\zeta(k), k = 0, 2, 3, 4, \dots$	<code>izeta(k)</code>	unsigned long int k
$\operatorname{Ei}(x) = \gamma + \ln(x) + \sum_{k=1}^{\infty} \frac{x^k}{k \cdot k!}, x > 0;$	<code>Ei(x)</code>	<code>Inf(x)=0</code> $\rightsquigarrow$ <code>-Inf</code> ; <code>Inf(x)&lt;0</code> $\rightsquigarrow$ <code>NaN</code> ;

#### Anmerkungen:

1. Mit den Funktionen aus obiger Tabelle werden zu einem Maschinenintervall  $x$  maximalgenaue Einschließungen aller Funktionswerte  $f(x)$ , mit  $x \in x$  berechnet. Lediglich die mit \* gekennzeichneten Funktionen liefern eine etwas gröbere Einschließung, da z.B. eine Einschließung der  $\Gamma'(x)$ -Funktionswerte mittels  $\Gamma'(x) = \psi(x) \cdot \Gamma(x)$  berechnet wird, wobei zwei Funktionen und eine Multiplikation auszuwerten sind, was mit maximal drei Rundungen verbunden ist.

#### 4.14.5. Skalarprodukt aus zwei Intervallprodukten

Um die komplexe Intervallmultiplikation realisieren zu können, benötigt man für die reellen Intervalle  $a, b, c, d$  vom Typ `MpfiClass` eine möglichst optimale Einschließung des Skalarproduktes

$$(4.1) \quad r := a \cdot b + c \cdot d,$$

wobei die reellen Variablen  $a, b, c, d$  als voneinander **unabhängig** anzusehen sind, so dass eine optimale Einschließung von  $r$  schon durch eine einfache Intervallauswertung der rechten Seite von (4.1) erreicht werden kann. Dabei sind aber noch zwei Probleme zu beachten:

1. Insbesondere bei schmalen Intervallen kann es bei der Intervallauswertung der rechten Seite von (4.1) wegen Auslöschungseffekten zu starken Überschätzungen kommen, die oft vermieden werden können, wenn man die Intervallauswertung z.B. mit doppelter Präzision durchführt.
2. Wenn die Randpunkte von z.B.  $a, b$  zu groß werden, kann es bei der Auswertung von z.B.  $a \cdot b$  zu einem vorzeitigen Überlauf kommen, obwohl  $r$  im Zahlenraster noch darstellbar ist. Das Problem kann wie folgt gelöst werden:

- Jedes Intervall, also auch  $a$ , wird so skaliert, dass gilt  $a \subseteq A \cdot 2^{k_A}$ . Dabei wird  $k_A \in \mathbb{Z}$  so gewählt, dass wenigstens einer der Randpunkte von  $2 \cdot (A \cdot A)$  möglichst groß ist aber gerade noch keinen Überlauf erzeugt. Im Fall  $a \neq [0, 0]$  wird  $k_A$  wie folgt berechnet:  $k_A = 536870910 - \max\{\text{expo}(\text{Inf}(a)), \text{expo}(\text{Sup}(a))\}$ , und im Fall  $a = [0, 0]$  wird  $k_A$  auf Null gesetzt. Mit diesen Überlegungen kann das Skalarprodukt rechts in (4.1) wie folgt berechnet werden:  $a \cdot b + c \cdot d \subseteq (A \cdot B) \cdot 2^{k_A+k_B} + (C \cdot D) \cdot 2^{k_C+k_D} = P_1 \cdot 2^{k_1} + P_2 \cdot 2^{k_2}$ , wobei die Teilprodukte  $P_1 := (A \cdot B)$  und  $P_2 := (C \cdot D)$  beide ohne Überlauf berechnet werden.
- Wir setzen jetzt voraus, dass  $k_1$  das Maximum der beiden Zweier-Exponenten ist. Dann kann die Summe  $S := P_1 \cdot 2^{k_1} + P_2 \cdot 2^{k_2}$  erst nach der folgenden Skalierung  $S := P_1 \cdot 2^{k_1} + P_2 \cdot 2^{k_2} = P_1 \cdot 2^{k_1} + (P_1 \cdot 2^{k_2-k_1}) \cdot 2^{k_1} = (P_1 + P_3) \cdot 2^{k_1}$  ausgewertet werden, wobei die Intervallsumme  $P_4 := (P_1 + P_3)$  **ohne** Überlauf berechnet werden kann. Ein Überlauf kann somit erst bei der Multiplikation  $P_4 \cdot 2^{k_1}$  eintreten, der dann aber unvermeidbar ist, sonst gilt:  $a \cdot b + c \cdot d \subseteq P_4 \cdot 2^{k_1}$ .

Abschließend noch ein Hinweis auf einen möglichen Nachteil dieses Verfahrens. Wenn einer der Randpunkte von  $a \subseteq (a \cdot 2^{-k_A}) \cdot 2^{k_A} = A \cdot 2^{k_A}$  schon selbst dicht vor dem Überlauf liegt, so wird  $k_A$  positiv sein, so dass es bei der Skalierung  $a \cdot 2^{-k_A}$  zu einer Überschätzung kommen kann, wenn der andere Randpunkt von  $a$  zu dicht am Unterlaufbereich liegt. Da aber diese Überschätzung nur eintreten wird, wenn das Intervall  $a$  praktisch ganz in einer der beiden komplexen Halbebenen liegt, kann aus praktischen Gründen eine solche, vergleichsweise geringe Überschätzung vermutlich immer in Kauf genommen werden.

Die oben beschriebene Skalierung  $a \subseteq A \cdot 2^{k_A}$  wird mit folgender Funktion realisiert.

```
void Interval_Scaling(MpfiClass& z, long int& k)
```

Obige Funktion liefert mit dem Eingabewert  $z = a$  den Rückgabewert  $z = A$  vom Typ `MpfiClass` in der Präzision von  $a$ , zusätzlich wird  $k_A$  zurückgegeben. Ist einer der beiden Randpunkte von  $a$  unendlich oder ein NaN, so wird  $[NaN, NaN]$  und  $k_A = 0$  zurückgegeben, sonst gilt  $a \subseteq A \cdot 2^{k_A}$ , wobei  $2 \cdot (A \cdot A)$  garantiert ohne Überlauf berechnet werden kann.  $a \subseteq A \cdot 2^{k_A}$  gilt nur in den ganz seltenen Fällen, wenn bei der Berechnung von  $a \cdot 2^{-k_A}$  ein Unterlauf eintritt, d.h. wenn gilt  $a \cdot 2^{-k_A} \subset A$ . Weitere Einzelheiten siehe oben unter 2.



```
MpfiClass scal_prod(const MpfiClass& a, const MpfiClass& b,
                    const MpfiClass& c, const MpfiClass& d);
```

Der Rückgabewert vom Typ `MpfiClass` ist eine meist optimale Einschließung des Skalarproduktes  $a \cdot b + c \cdot d$  in der Current-Präzision, wobei ein vorzeitiger Über- oder Unterlauf vermieden wird. Die Präzisionen der 4 Eingabeintervalle können verschieden sein. Weitere Einzelheiten siehe Seite 72 unter Punkt 2. Zur Anwendung kommt diese Funktion bei der Multiplikation komplexer Intervalle.

Das folgende Programm zeigt die Einschließung spezieller Skalarprodukte mithilfe der Funktion `scal_prod()`.

```
1 // MPFR-06.cpp
2 #include "mpficlass.hpp"
3
4 using namespace MPFR;
5 using namespace MPFI;
6 using namespace std;
7
8 int main(void)
9 {
10  MpfiClass::SetCurrPrecision (300);
11  cout << "Current-Precision = " << MPFR::MpfrClass::GetCurrPrecision() << endl;
12
13  MpfrClass x(372130600.0);
14  x = exp(x);
15  MpfiClass a(x,x), b(succ(x),succ(x)), c(-x,-x), d(x,x), r;
16  r = scal_prod(a,b,c,d);
17  cout.precision( r.GetPrecision()/3.321928095);
18  cout << "r = " << r << endl;
19  cout << "r = " << a*b + c*d << endl;
20
21  return 0;
22 }
```

In Zeile 14 wird mit  $x = 372130600.0$  und  $\exp(x)$  eine hinreichend große Maschinenzahl erzeugt. In Zeile 15 werden  $a \approx b, c = -a, d = a$  als Punktintervalle definiert, so dass bei der Auswertung ihres Skalarproduktes  $r := a \cdot b + c \cdot d$  mit starker Auslöschung und mit einem vorzeitigen Überlauf zu rechnen ist. Dies wird durch die folgende Programmausgabe bestätigt

```
r = [1.554622378...836989271e323228442,1.554622378...836989271e323228442]
r = [-@Inf@, @Inf@]
```

wobei die Einschließung  $r \subseteq \mathbf{r}$  realisiert wird.

Mit dem kleineren Argument  $x = 372130500.0$  erhält man die Ausgabe

```
r = [1.296663701...349242307e323228355,1.296663701...349242308e323228355]
r = [0,2.753514027...992982088e323228355]
```

Die erste Ausgabe liefert eine optimale Einschließung von  $r$ , die intern mit doppelter Präzision berechnet wird. Die zweite Ausgabe wird mit einfacher Präzision `prec = 300` berechnet und zeigt bez. des linken Randpunktes Null eine erhebliche Überschätzung der exakten Einschließung.

## 4.15. Optimale Intervall-Einschließungen

Mit den Intervallfunktionen  $f(\mathbf{x})$  aus den Tabellen auf Seite 68 und 71 erhält man zu einem vorgegebenen, reellen Eingangsintervall  $\mathbf{x}$  i.a. maximalgenaue Einschließungen des entsprechenden Wertebereichs

$$W_{\mathbf{x}} := \{y \in \mathbb{R} \mid y = f(x) \wedge x \in \mathbf{x}\}.$$

Mit dem Funktionsaufruf  $\mathbf{y} = f(\mathbf{x})$ , also z.B.  $\mathbf{y} = \mathbf{exp}(\mathbf{x})$ , erhält man die i.a. maximalgenaue Einschließung

$$W_{\mathbf{x}} \subseteq \mathbf{y}.$$

In der Praxis müssen aber oft auch Intervallausdrücke ausgewertet werden, die aus den vier Grundoperationen zusammen mit den oben genannten Intervallfunktionen  $f(\mathbf{x})$  aufgebaut sind. Als Beispiel betrachten wir nach [6, S. 28] die beiden vereinfachten Funktionen

$$f_1(x) = \frac{x}{1-x}, \quad x \neq 0, \quad x \neq 1; \quad f_2(x) = \frac{1}{\frac{1}{x} - 1}, \quad x \neq 0, \quad x \neq 1.$$

Man bestätigt leicht, dass mit  $\mathbf{x} = [2, 3]$  beide Funktionen den gleichen Wertebereich  $W_{\mathbf{x}} = [-2, -1.5]$  besitzen. Wertet man dann beide Funktionen intervallmäßig aus, so erhält man mit den Einschließungen  $\mathbf{y}_1, \mathbf{y}_2$  das Ergebnis

$$\mathbf{y}_2 := \frac{1}{1/[2, 3] - 1} = \frac{1}{[-2/3, -1/2]} = [-2, -3/2] = W_{\mathbf{x}} \subset \mathbf{y}_1 := \frac{[2, 3]}{1 - [2, 3]} = \frac{[2, 3]}{[-2, -1]} = [-3, -1].$$

$f_2(x)$  liefert damit die **optimale** Einschließung  $\mathbf{y}_2$ , da die Variable  $x$  im Funktionsterm von  $f_2(x)$  explizit **nur einmal** vorkommt, [6, S. 32]. Wir fassen zusammen:

Wird ein aus stetigen Funktionen zusammengesetzter Funktionsterm über einem Intervall  $\mathbf{x}$  intervallmäßig ausgewertet, so ist die berechnete Einschließung  $\mathbf{y}$  i.a. **nicht optimal**, d.h. es gilt  $W_{\mathbf{x}} \subset \mathbf{y}$ .  
Wenn jedoch die Variable  $\mathbf{x}$  im Funktionsterm explizit **nur einmal** vorkommt, so gilt  $W_{\mathbf{x}} = \mathbf{y}$ , d.h. die Einschließung ist dann **optimal**.

### Anmerkungen:

1. Damit die Variable  $\mathbf{x}$  im Funktionsterm nur einmal vorkommt, muss der ursprüngliche Funktionsterm entsprechend äquivalent umgeformt werden, vgl. obiges Beispiel und das Beispiel in [6, S. 36]. Eine solche Umformung ist jedoch nicht immer möglich!
2. In [6, S. 32] findet man für Polynome eine weitere Bedingung für eine optimale Einschließung.
3. In [6] findet man ab Seite 33 noch qualitative Aussagen über die Güte der Einschließung  $\mathbf{y}$  im Vergleich zum Wertebereich  $W_{\mathbf{x}}$ .
4. Für die komplexe Intervallmultiplikation gilt:

$$\mathbf{z} \cdot \mathbf{w} = (X_{\mathbf{z}} + i \cdot Y_{\mathbf{z}}) \cdot (X_{\mathbf{w}} + i \cdot Y_{\mathbf{w}}) = X_{\mathbf{z}} \cdot X_{\mathbf{w}} - Y_{\mathbf{z}} \cdot Y_{\mathbf{w}} + i \cdot (X_{\mathbf{z}} \cdot Y_{\mathbf{w}} + Y_{\mathbf{z}} \cdot X_{\mathbf{w}}).$$

Betrachtet man z.B. den Realteil  $X_{\mathbf{z}} \cdot X_{\mathbf{w}} - Y_{\mathbf{z}} \cdot Y_{\mathbf{w}}$ , so kommt dort jedes der vier reellen Intervalle jeweils nur einmal vor, so dass die Intervallauswertung automatisch eine optimale Einschließung des Realteils von  $\mathbf{z} \cdot \mathbf{w}$  liefert, entsprechendes gilt für den Imaginärteil. Aber beachten Sie, dass für die komplexe Division  $\mathbf{z}/\mathbf{w}$  diese Überlegungen nicht mehr gelten, so dass für die optimale komplexe Intervalldivision ein neuer Algorithmus gebraucht wird.

# 5. MpfcClass-Interface für komplexe Langzahlrechnungen in C-XSC

## 5.1. Grundlegendes

Das MpfcClass-Interface ist eine in `mpfcclass.hpp` und `mpfcclass.cpp` implementierte C++-Wrapper-Klasse `MpfcClass` für die C-Bibliothek MPFR, deren C-Funktionen über die implementierten Operatoren und Funktionen aufgerufen werden. Alle Funktionen mit der möglichen Übergabe eines `RoundingModes` oder eines `PrecisionTypes` besitzen als Standard die Werte von `CurrRndMode` bzw. `CurrPrecision`, die beide beliebig gesetzt werden können. Dies gilt auch für alle Konstruktoren.

### 5.1.1. Allgemein

Um in C-XSC das Interface verwenden zu können, muss der Header `mpfcclass.hpp` eingebunden werden. In ihm sind die benötigten Header-Dateien der MPFR-Bibliothek enthalten. Die `MpfcClass`-Klasse liegt im Namensraum "MPFR".

### 5.1.2. Aufbau

Die Klasse besteht intern aus zwei "mpfr\_t"-Variablen. Diese dienen zum Speichern von Real- und Imaginärteil. Zusätzlich gibt es static Elemente, um den Standard-Rundungsmodus, die Standard-Precision und die aktuelle Basis zu speichern.

### 5.1.3. Präzision

Die Präzision gibt die Anzahl der binären Mantissenstellen einer `MpfcClass`-Variablen an, Real- und Imaginärteil erhalten stets die gleiche Präzision. Der Wert `prec` muss mindestens 2 betragen. Die Current-Precision kann global gesetzt werden; wenn dies nicht geschieht, so wird mit der Default-Precision von 53 Bits gerechnet. Unabhängig davon kann die Präzision für jede `MpfcClass`-Variable auch einzeln festgelegt werden.

### 5.1.4. Variablentyp `PrecisionType`

Mithilfe des Variablentyps `PrecisionType` (Name der Variablen meist `prec`) kann der Präzisionswert einer `MpfcClass`-Variablen eingestellt werden. Der Variablentyp ist ein typedef für eine "mp\_prec\_t"-Variable.

### 5.1.5. Variablentyp `RoundingMode`

Mit dem Variablentyp "RoundingMode" (Name der Variablen meist `rnd`) wird der gewünschte Rundungsmodus eingestellt. Der Variablentyp ist ein typedef für eine `mpfr_rnd_t`-Variable.

- `RoundNearest` – Der Wert wird zur nächsten Rasterzahl gerundet (0)
- `RoundUp` – Der Wert wird aufgerundet (2)
- `RoundDown` – Der Wert wird abgerundet (3)

- RoundToZero – Der Wert wird in Richtung Null gerundet (1)
- RoundFromZero – Rundung weg von der Null (4)

## 5.2. Rundungsmodi und Precision Handling

```
PrecisionType GetPrecision () const;
```

Diese Memberfunktion gibt die maximale Präzision von Real- und Imaginärteil des aktuellen Objekts in Bits zurück. Dabei entsprechen z.B. 302 Bits  $302/\log_2(10) \approx 91$  Dezimalstellen.

```
void SetPrecision (PrecisionType prec);
```

Diese Memberfunktion setzt die Präzision des aktuellen Objektes auf `prec`. Sein Wert bleibt nicht erhalten.

```
void RoundPrecision (PrecisionType prec, RoundingMode rnd);
```

Diese Memberfunktion rundet das aktuelle Objekt auf die neue Precision `prec`, sein Wert bleibt dabei i.a. nicht erhalten. Sollte die Präzision des Objektes größer sein als `prec`, wird das Objekt mit Hilfe des eingestellten Rundungsmodus so gerundet, dass es in das Format der Präzision `prec` passt. Ist die Präzision kleiner als `prec`, werden die restlichen binären Stellen mit Nullen aufgefüllt.

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision in `MpfrClass`, `MpfiClass` und `MpfcClass` auf `prec`. Wird die Current-Precision nicht gesetzt, so wird mit der Default-Precision von 53 Bits gerechnet. Durch das Setzen der Current-Precision werden die Präzisionen der bis dahin benutzten Variablen **nicht** geändert.

```
const int MpfrClass::GetBase ();
```

Gibt die aktuelle Basis zurück, diese hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die Funktion ist definiert in `mpfrclass.cpp`.

```
void MpfrClass::SetBase (int b);
```

Setzt die aktuelle Basis auf `b`. Dies muss ein Wert zwischen 2 und 36 sein. Die Basis hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die internen Rechnungen erfolgen stets im Binärsystem! Die Funktion ist definiert in `mpfrclass.cpp`.

```
static const RoundingMode GetCurrRndMode ();
```

Gibt den aktuellen Rundungsmodus zurück mit den Werten: (0, 1, 2, 3, 4).

```
static void SetCurrRndMode (RoundingMode rnd);
```

Setzt den Current-Rundungsmodus in der Klasse `MpfrClass` und damit automatisch auch in der Klasse `MpfcClass` auf `rnd`. Für `rnd` sind fünf verschiedene Modi möglich:

- RoundNearest: Rundung zur nächsten Rasterzahl (0)
- RoundToZero: Rundung in Richtung Null (1)
- RoundFromZero: Rundung weg von der Null (4)
- RoundUp: Aufrunden (2)
- RoundDown: Abrunden (3)

Wird der Rundungsmodus mit `SetCurrRndMode` nicht gesetzt, so wird als Default-Rundungsmodus `RoundNearest` benutzt.

## 5.3. Konstruktoren / Destruktor

### 5.3.1. Konstruktoren

```
MpfcClass ();
```

Der Default-Konstruktor legt ein neues Element mit der Current-Precision an.  
Der Aufruf `MpfcClass y;` initialisiert den Wert: `y = (NaN, NaN);`

```
MpfcClass::MpfcClass(const MpfcClass& z,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const cxsc::complex& z,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const MpfrClass& x, const MpfrClass& y,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const mpfr_t& x, const mpfr_t& y,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const cxsc::real& x, const cxsc::real& y,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const double& x, const double& y,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(int x, int y,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const MpfrClass& x,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const mpfr_t& x,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const cxsc::real& x,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(const double& x,
                    RoundingMode rnd, PrecisionType prec);
MpfcClass::MpfcClass(int x,
                    RoundingMode rnd, PrecisionType prec);
```

Mit den Konstruktoraufrufen `MpfcClass w(z);` oder `MpfcClass w(x, y);` werden die Objekte `z` bzw. `x, y` nach `w` in die Current-Präzision gerundet. Mit `rnd` und `prec` sind Rundungen auch in eine andere Präzision möglich. Real- und Imaginärteil von `w` erhalten stets die **gleiche** Präzision.

Mögliche Konstruktor-Aufrufe sind:

1. `MpfcClass w(z); MpfcClass w(x, y);`
2. `MpfcClass w(z, RoundNearest); MpfcClass w(x, y, RoundNearest);`
3. `MpfcClass w(z, RoundDown, 3); MpfcClass w(x, y, RoundDown, 3);`

**Zu 1.** Die Objekte `z, x, y` werden mit dem Current-Rundungsmodus in das Objekt `w` mit der Current-Precision gerundet.

**Zu 2.** Die Objekte `z, x, y` werden mit `RoundNearest` in das Objekt `w` mit der Current-Precision gerundet.

**Zu 3.** Die Objekte `z, x, y` werden in das Objekt `w` auf die (sehr kleine) Präzision `prec = 3` abgerundet.

Der Aufruf `MpfcClass w(x, y, 3);` führt zu einer Fehlermeldung, da in der Parameterliste der Rundungsmodus fehlt. Die oberen 18 Konstruktoren erlauben also eine sehr flexible Initialisierung von `MpfcClass`-Objekten.

```
MpfcClass::MpfcClass (const std::string& s, RoundingMode rnd,  
                    PrecisionType prec);
```

Der Aufruf `MpfcClass z(s);` rundet den String `s` mittels der voreingestellten Current-Precision und mit dem Current-Rundungsmodus in ein Klassenobjekt `z` vom Typ `MpfcClass`. Mit dem Aufruf `MpfcClass z(s, RoundUp);` wird der String `s` mit der Current-Precision in ein Klassenobjekt `z` vom Typ `MpfcClass` aufgerundet. Der Aufruf `MpfcClass z(s, RoundNearest, 140);` rundet `s` in ein Klassenobjekt `z` vom Typ `MpfcClass` mit einer Präzision von 140 Bits. Gerundet wird dabei zur nächsten Rasterzahl dieses Formats. Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in eine binäre Zahl konvertiert werden kann, so dass daher Rundungen i.a. nicht zu vermeiden sind.

Der String `s` muss das Format `(Number,Number)` ohne Leerzeichen besitzen.

### 5.3.2. Destruktor

```
~MpfcClass ();
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

## 5.4. Zuweisungs-Operatoren

Bei den folgenden sechs Zuweisungs-Operationen erhält der linke Operand als Präzision stets die aktuelle Current-Precision, und der rechte Operand `op` wird nach dem aktuellen Current-Rundungsmodus in den linken Operanden gerundet.

```
MpfcClass& operator = (const MpfcClass& op);  
MpfcClass& operator = (const MpfrClass& op);  
MpfcClass& operator = (const double& op);  
MpfcClass& operator = (const cxsc::real& op);  
MpfcClass& operator = (const cxsc::complex& op);  
MpfcClass& operator = (const int opt);
```

Ist z.B. `op` vom Typ `real`, `double` oder `complex` und ist die Current-Precision kleiner als 53, so wird `op` i.a. in den linken Operanden gerundet. Nur wenn die Current-Precision größer oder gleich 53 ist, erhält der linke Operand genau den Wert des rechten Operanden. Für andere Typen des rechten Operanden gelten ganz entsprechende Aussagen.

Ist z.B. `op` vom Typ `MpfrClass` und ist seine Präzision größer als die Current-Precision, so wird `op` in den Realteil des linken Operanden bez. des Current-Rundungsmodus gerundet, d.h. der Realteil des linken Operanden wird dann i.a. vom Wert des rechten Operanden **verschieden** sein!

### Zusammenfassung:

Bei einer Wertzuweisung erhält der linke Operand stets die Current-Präzision. Ist diese kleiner als die Präzision des rechten Operanden, so wird dieser mittels des Current-Rundungsmodus in den linken Operanden gerundet. Ist die Current-Präzision größer oder **gleich** der Präzision des rechten Operanden, so wird dieser **rundungsfehlerfrei** nach links übertragen. Real- und Imaginärteil erhalten stets die gleiche Präzision.

```
MpfcClass& operator = (const std::string& s);
```

Der String `s` muss die Form `(Number,Number)` haben und darf keine Leerzeichen enthalten. Real- und Imaginärteil von `s` werden mit dem Current-Rundungsmodus in die Current-Präzision gerundet. Beachten Sie, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in das interne Binärformat gerundet wird, so dass dann Rundungen unvermeidbar sind.



## 5.5. Eingabe / Ausgabe

```
std::ostream& operator << (std::ostream& os, const MpfcClass& z);
```

Ermöglicht die Ausgabe einer MpfcClass-Variablen  $z$  über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(z.GetPrecision()/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Real- und Imaginärteil von  $z$  werden in der mit `SetBase(k)` voreingestellten Basis im Format

```
(Realteil, Imaginärteil)
```

ausgegeben, wobei Real- und Imaginärteil beide mit dem voreingestellten Current-Rundungsmodus gerundet werden. Für die meist dezimale Ausgabe ist natürlich  $k = 10$  zu wählen.

```
std::istream& operator >> (std::istream& is, MpfiClass& z);
```

Ermöglicht das Einlesen einer MpfcClass-Variablen über den Standard-Eingabestrom "cin". Die eingegebene komplexe Zahl ist auf keine Stellenanzahl begrenzt. Nur das folgende Eingabeformat ist zulässig, wobei mit dem Current-Rundungsmodus in das interne Binärsystem gerundet wird:

- $(-1.23e-4401, +2.3E+2000)$

**Vorsicht:** Leerzeichen sind nur nach dem Komma erlaubt, und die runden Klammern müssen beide gesetzt werden, sonst erfolgt eine Fehlermeldung.

## 5.6. Typ-Umwandlungen

Um ein flexibles Arbeiten zwischen C-XSC und der Klasse `MpfcClass` zu ermöglichen, wurden möglichst viele Typ-Umwandlungs-Funktionen bereitgestellt:

### 5.6.1. `real, double, complex, ...` → **MPFC**

```
MpfcClass MpfrClass2Mpfc (const MpfrClass& x);
MpfcClass mpfr_t2Mpfc (const mpfr_t& x);
MpfcClass real2Mpfc (const cxsc::real& x);
MpfcClass double2Mpfc (const double& x);
MpfcClass complex2Mpfc (const cxsc::complex& x);
MpfcClass int2Mpfc (const int& x);
```

Obige Funktionen liefern mit dem jeweiligen Eingabewert `x` einen Rückgabewert vom Typ `MpfcClass` in einer Präzision, die gewährleistet, dass der Rückgabewert **genau** dem Wert von `x` entspricht. Die Präzision des Rückgabewertes wird also i.a. nicht mit der Current-Precision übereinstimmen!

### 5.6.2. **MPFC** → `complex`

```
cxsc::complex to_complex (const MpfcClass& z, RoundingMode rnd);
```

Obige Funktion rundet mit `rnd` den Eingabewert `z` in einen C-XSC Rückgabewert vom Typ `complex`.

### 5.6.3. **MPFC** → `mpfr_t`

```
mpfr_t& MpfcClass::GetValueRe();
```

Die Memberfunktion `GetValueRe()` liefert für das jeweils aktuelle Objekt vom Typ `MpfcClass` eine Referenz auf seinen Realteilwert `mpfr_re` vom Typ `mpfr_t`. Damit kann dieser Realteil z.B. an einen nicht-konstanten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
mpfr_t& MpfcClass::GetValueIm();
```

Die Memberfunktion `GetValueIm()` liefert für das jeweils aktuelle Objekt vom Typ `MpfcClass` eine Referenz auf seinen Imaginärteil `mpfr_im` vom Typ `mpfr_t`. Damit kann dieser Imaginärteil z.B. an einen nicht-konstanten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
const mpfr_t& getvalueRe(const MpfcClass& z)
```

Die Funktion `getvalueRe(...)` liefert für das als `const` definierte Objekt `z` vom Typ `MpfcClass` eine Referenz auf seinen Realteil `mpfr_re` vom Typ `mpfr_t`. Damit kann dieser Realteil z.B. an einen als `const` deklarierten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
const mpfr_t& getvalueIm(const MpfcClass& z)
```

`getvalueIm(...)` liefert für das als `const` definierte Objekt `z` vom Typ `MpfcClass` eine Referenz auf seinen Imaginärteil `mpfr_im` vom Typ `mpfr_t`. Damit kann dieser Imaginärteil z.B. an einen als `const` deklarierten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

#### 5.6.4. mpfr\_t → MPFC

```
void MpfcClass::SetValueRe(const mpfr_t& v)
```

Mit dieser Memberfunktion wird der Realteil des aktuellen Objekts exakt auf  $v$  gesetzt. Der Imaginärteil des aktuellen Objekts bleibt dabei erhalten. Die Präzision des aktuellen Objekts wird gesetzt auf das Maximum der Präzisionen von  $v$  und `mpfr_im`, d.h. die Präzisionen von Real- und Imaginärteil sind nach dem Funktionsaufruf wieder gleich.

```
void MpfcClass::SetValueIm(const mpfr_t& v)
```

Mit dieser Memberfunktion wird der Imaginärteil des aktuellen Objekts exakt auf  $v$  gesetzt. Der Realteil des aktuellen Objekts bleibt dabei erhalten. Die Präzision des aktuellen Objekts wird gesetzt auf das Maximum der Präzisionen von  $v$  und `mpfr_re`, d.h. die Präzisionen von Real- und Imaginärteil sind nach dem Funktionsaufruf wieder gleich.

```
void MpfcClass::SetValue(const mpfr_t& re, const mpfr_t& im)
```

Mit dieser Memberfunktion werden Real- und Imaginärteil des aktuellen Objekts exakt auf  $re$  bzw.  $im$  gesetzt. Das aktuelle Objekt erhält als Präzision das Maximum der Präzisionen von  $re$  und  $im$ , d.h. die Präzisionen von Real- und Imaginärteil sind nach dem Funktionsaufruf wieder gleich.

#### 5.6.5. MPFC → string

```
std::string to_string(const MpfcClass& z, RoundingMode rnd,  
PrecisionType prec);
```

$z = x + i \cdot y$  wird mittels `rnd` in einen String  $s$  mit `prec` Dezimalstellen gerundet, wenn Base gleich 10 ist. Der String besitzt das Format

(Number, Number)

wobei intern keine Leerzeichen auftreten. Wählt man `prec` hinreichend groß, so stellt der String den Wert von  $z$  **exakt** dar, weil eine Binärzahl stets exakt in eine Dezimalzahl umgewandelt werden kann.

Wird `prec` nicht angegeben, so werden  $x, y$  mittels `rnd` in einen String gerundet, der bei Base=10 so viele Dezimalstellen besitzt, wie es der Präzision von  $z$  entspricht. Ist die Präzision von  $z$  z.B. 302, so werden im String Real- und Imaginärteil mit  $302/\log_2(10) \approx 91$  Dezimalstellen generiert. Wird neben `prec` auch `rnd` nicht angegeben, so wird mit dem Current-Rundungsmodus in den String mit gleicher Dezimalstellenzahl gerundet.

### 5.6.6. string → MPFC

```
MpfcClass string2Mpfc(const std::string& s, Roundingmode rnd,  
                    PrecisionType prec);
```

Der Aufruf `string2Mpfc(s);` rundet den String `s` mittels der voreingestellten Current-Precision und mit dem Current-Rundungsmodus in ein Klassenobjekt vom Typ `MpfcClass`. Mit dem Aufruf `string2Mpfc(s, RoundUp);` wird der String `s` mit der Current-Precision in ein Klassenobjekt vom Typ `MpfcClass` aufgerundet. Der Aufruf `string2Mpfc(s, RoundNearest, 140);` rundet `s` in ein Klassenobjekt vom Typ `MpfcClass` mit der Präzision 140 Bits. Gerundet wird dabei zur nächsten Rasterzahl dieses Formats. Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in eine binäre Zahl konvertiert werden kann, so dass daher Rundungen i.a. nicht zu vermeiden sind. Eine weitere Möglichkeit, einen String in ein `MpfcClass`-Objekt zu verwandeln, besteht in einem entsprechenden Konstruktor-Aufruf, vgl. Seite 79.

Der String `s` muss das Format `(Number,Number)` ohne Leerzeichen besitzen.

### 5.6.7. MPFC → MPFC

Beachten Sie bitte, dass bei einer Wertzuweisung an eine `MpfcClass`-Variable mit Hilfe des Operators `=` der linke Operand mit Real- und Imaginärteil stets auf die Current-Precision gesetzt wird und dass der rechte `MpfcClass`-Operand dabei **stets** bez. des Current-Rundungsmodus in den linken Operanden gerundet wird, vgl. dazu auch Seite 80. Will man jedoch abweichend von dieser Rundung einen anderen Rundungsmodus benutzen, so kann dies mit folgender Funktion ohne Rückgabewert realisiert werden:

```
void set_Mpfc (MpfcClass& op, const MpfcClass& op1, RoundingMode rnd,  
             PrecisionType prec);
```

Folgende Funktionsaufrufe sind möglich:

1. `set_Mpfc (op, op1, RoundUp, prec);`  
`op` erhält die Präzision `prec` und den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls aufgerundet wird. Setzt man `prec` gleich der Präzision von `op1`, so erhält `op` den **exakten** Wert von `op1`, und zwar unabhängig vom gewählten Rundungsmodus `rnd`.
2. `set_Mpfc (op, op1, RoundDown);`  
`op` wird auf die Current-Precision gesetzt und erhält den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls abgerundet wird.
3. `set_Mpfc (op, op1);`  
`op` wird auf die Current-Precision gesetzt und erhält den i.a. gerundeten Wert von `op1`, wobei hier gegebenenfalls bez. des Current-Rundungsmodus gerundet wird.

## 5.7. Abfragen

Bei allen folgenden Abfragefunktionen braucht die Präzision von `x` nicht mit der Current-Precision übereinzustimmen.

```
bool isNaN (const MpfcClass& x);  
bool isInf (const MpfcClass& x);  
bool isNumber(const MpfcClass& x);  
bool isZero (const MpfcClass& x);
```

`isNaN` und `isInf` prüfen, ob Real- oder Imaginärteil von `x` gleich NaN bzw.  $\pm\text{Inf}$  sind. `isNumber` überprüft, ob Real- und Imaginärteil von `x` normale `MpfcClass` Zahlen ungleich NaN und ungleich  $\pm\text{Inf}$  sind, und `isZero` überprüft, ob Real- und Imaginärteil von `x` gleich Null sind.

## 5.8. Vergleiche

Es werden die üblichen Vergleichsoperatoren implementiert, wobei wenigstens ein Operand vom Typ `MpfcClass` sein muss. Die Operanden können unterschiedliche Präzisionen besitzen.

### 5.8.1. Vergleichsoperatoren `=`, `!=`

```
bool operator == (const MpfcClass& y, const MpfcClass& x);
bool operator == (const MpfcClass& y, const cxsc::complex& x);
bool operator == (const MpfcClass& y, const MpfrClass& x);
bool operator == (const MpfcClass& y, const mpfr_t& x);
bool operator == (const MpfcClass& y, const double& x);
bool operator == (const MpfcClass& y, const cxsc::real& x);
bool operator == (const MpfcClass& y, const int x);
bool operator == (const MpfrClass& y, const MpfcClass& x);
bool operator == (const cxsc::complex& y, const MpfcClass& x);
bool operator == (const mpfr_t& y, const MpfcClass& x);
bool operator == (const double& y, const MpfcClass& x);
bool operator == (const cxsc::real& y, const MpfcClass& x);
bool operator == (const int y, const MpfcClass& x);
```

```
bool operator != (const MpfcClass& y, const MpfcClass& x);
bool operator != (const MpfcClass& y, const cxsc::complex& x);
bool operator != (const MpfcClass& y, const MpfrClass& x);
bool operator != (const MpfcClass& y, const mpfr_t& x);
bool operator != (const MpfcClass& y, const double& x);
bool operator != (const MpfcClass& y, const cxsc::real& x);
bool operator != (const MpfcClass& y, const int x);
bool operator != (const MpfrClass& y, const MpfcClass& x);
bool operator != (const cxsc::complex& y, const MpfcClass& x);
bool operator != (const mpfr_t& y, const MpfcClass& x);
bool operator != (const double& y, const MpfcClass& x);
bool operator != (const cxsc::real& y, const MpfcClass& x);
bool operator != (const int y, const MpfcClass& x);
```

## 5.9. Arithmetische Operatoren

Für alle arithmetischen Operatoren gilt:

Das exakte Ergebnis einer arithmetischen Operation wird unabhängig von der Präzision der Operanden mit dem voreingestellten Current-Rundungsmodus optimal gerundet. Die Ergebnis-Präzision ist dabei stets gleich der voreingestellten Current-Precision.

Die Operatoren  $\odot =$ , mit  $\odot \in \{+, -, \cdot, /\}$ , bedeuten  $u \odot = v \iff u = u \odot v$ . Dabei wird  $u \odot v$  mit dem Current-Rundungsmodus in die Current-Precision gerundet und in  $u$  gespeichert, wobei  $u$  als Präzision die Current-Precision erhält.

### 5.9.1. Addition

```
MpfcClass operator + (const MpfcClass& x, const MpfcClass& y);
MpfcClass operator + (const MpfcClass& x, const cxsc::complex& y);
MpfcClass operator + (const MpfcClass& x, const MpfrClass& y);
MpfcClass operator + (const MpfcClass& x, const mpfr_t& y);
MpfcClass operator + (const MpfcClass& x, const double& y);
MpfcClass operator + (const MpfcClass& x, const cxsc::real& y);
MpfcClass operator + (const MpfcClass& x, const int y);
MpfcClass operator + (const cxsc::complex& x, const MpfcClass& y);
MpfcClass operator + (const MpfrClass& x, const MpfcClass& y);
MpfcClass operator + (const mpfr_t& x, const MpfcClass& y);
MpfcClass operator + (const double& x, const MpfcClass& y);
MpfcClass operator + (const cxsc::real& x, const MpfcClass& y);
MpfcClass operator + (const int x, const MpfcClass& y);
```

```
MpfcClass& operator += (MpfcClass&, const MpfcClass&);
MpfcClass& operator += (MpfcClass&, const cxsc::complex&);
MpfcClass& operator += (MpfcClass&, const MpfrClass&);
MpfcClass& operator += (MpfcClass&, const mpfr_t&);
MpfcClass& operator += (MpfcClass&, const double&);
MpfcClass& operator += (MpfcClass&, const cxsc::real&);
MpfcClass& operator += (MpfcClass&, const int);
```

## 5.9.2. Subtraktion

Beachten Sie den Hinweis auf Seite 87.

```
MpfcClass operator - (const MpfcClass& x, const MpfcClass& y);
MpfcClass operator - (const MpfcClass& x, const cxsc::complex& y);
MpfcClass operator - (const MpfcClass& x, const MpfrClass& y);
MpfcClass operator - (const MpfcClass& x, const mpfr_t& y);
MpfcClass operator - (const MpfcClass& x, const double& y);
MpfcClass operator - (const MpfcClass& x, const cxsc::real& y);
MpfcClass operator - (const MpfcClass& x, const int y);
MpfcClass operator - (const cxsc::complex& x, const MpfcClass& y);
MpfcClass operator - (const MpfrClass& x, const MpfcClass& y);
MpfcClass operator - (const mpfr_t& x, const MpfcClass& y);
MpfcClass operator - (const double& x, const MpfcClass& y);
MpfcClass operator - (const cxsc::real& x, const MpfcClass& y);
MpfcClass operator - (const int x, const MpfcClass& y);
```

```
MpfcClass& operator -= (MpfcClass&, const MpfcClass&);
MpfcClass& operator -= (MpfcClass&, const cxsc::complex&);
MpfcClass& operator -= (MpfcClass&, const MpfrClass&);
MpfcClass& operator -= (MpfcClass&, const mpfr_t&);
MpfcClass& operator -= (MpfcClass&, const double&);
MpfcClass& operator -= (MpfcClass&, const cxsc::real&);
MpfcClass& operator -= (MpfcClass&, const int);
```



### 5.9.3. Multiplikation

Beachten Sie den Hinweis auf Seite 87.

```
MpfcClass operator * (const MpfcClass& x, const MpfcClass& y);
MpfcClass operator * (const MpfcClass& x, const cxsc::complex& y);
MpfcClass operator * (const MpfcClass& x, const MpfrClass& y);
MpfcClass operator * (const MpfcClass& x, const mpfr_t& y);
MpfcClass operator * (const MpfcClass& x, const double& y);
MpfcClass operator * (const MpfcClass& x, const cxsc::real& y);
MpfcClass operator * (const MpfcClass& x, const int y);
MpfcClass operator * (const cxsc::complex& x, const MpfcClass& y);
MpfcClass operator * (const MpfrClass& x, const MpfcClass& y);
MpfcClass operator * (const mpfr_t& x, const MpfcClass& y);
MpfcClass operator * (const double& x, const MpfcClass& y);
MpfcClass operator * (const cxsc::real& x, const MpfcClass& y);
MpfcClass operator * (const int x, const MpfcClass& y);
```

```
MpfcClass& operator *= (MpfcClass&, const MpfcClass&);
MpfcClass& operator *= (MpfcClass&, const cxsc::complex&);
MpfcClass& operator *= (MpfcClass&, const MpfrClass&);
MpfcClass& operator *= (MpfcClass&, const mpfr_t&);
MpfcClass& operator *= (MpfcClass&, const double&);
MpfcClass& operator *= (MpfcClass&, const cxsc::real&);
MpfcClass& operator *= (MpfcClass&, const int);
```

### 5.9.4. Division

Beachten Sie den Hinweis auf Seite 87.

```
MpfcClass operator / (const MpfcClass& x, const MpfcClass& y);
MpfcClass operator / (const MpfcClass& x, const cxsc::complex& y);
MpfcClass operator / (const MpfcClass& x, const MpfrClass& y);
MpfcClass operator / (const MpfcClass& x, const mpfr_t& y);
MpfcClass operator / (const MpfcClass& x, const double& y);
MpfcClass operator / (const MpfcClass& x, const cxsc::real& y);
MpfcClass operator / (const MpfcClass& x, const int y);
MpfcClass operator / (const cxsc::complex& x, const MpfcClass& y);
MpfcClass operator / (const MpfrClass& x, const MpfcClass& y);
MpfcClass operator / (const mpfr_t& x, const MpfcClass& y);
MpfcClass operator / (const double& x, const MpfcClass& y);
MpfcClass operator / (const cxsc::real& x, const MpfcClass& y);
MpfcClass operator / (const int x, const MpfcClass& y);
```

```
MpfcClass& operator /= (MpfcClass&, const MpfcClass&);
MpfcClass& operator /= (MpfcClass&, const cxsc::complex&);
MpfcClass& operator /= (MpfcClass&, const MpfrClass&);
MpfcClass& operator /= (MpfcClass&, const mpfr_t&);
MpfcClass& operator /= (MpfcClass&, const double&);
MpfcClass& operator /= (MpfcClass&, const cxsc::real&);
MpfcClass& operator /= (MpfcClass&, const int);
```

## 5.10. Mathematische Funktionen

### 5.10.1. Standard-Implementierung

Die Implementierung komplexwertiger Funktionen mit komplexen Punktargumenten der Klasse `MpfcClass` erfolgt ganz analog zu den Funktionen der Klasse `MpfrClass`, vgl. Seite 32.

Unabhängig von der Präzision des Eingangsarguments werden Real- und Imaginärteil des Funktionswertes mit dem Current-Rundungsmodus oder mit dem Rundungsparameter `rnd` in die Current-Precision gerundet.

Es gibt jedoch Funktionen, wie z.B. die komplexe Konjugation, bei denen man von dieser Standard-Implementierung abweichen sollte. Im folgenden Abschnitt werden diese Funktionen kurz beschrieben.

### 5.10.2. Davon abweichende Funktionen

```
MpfcClass conj (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  ist der Rückgabewert von `conj(z)` gegeben durch  $x - i \cdot y$ , wobei die Präzision von  $x$  und  $y$  nicht geändert wird und daher mit der Current-Precision auch nicht übereinstimmen muss.

```
MpfrClass abs (const MpfcClass& z, RoundingMode rnd, PrecisionType prec);
```

Mit  $z = x + i \cdot y$  wird  $|z| = \sqrt{x^2 + y^2}$  i.a. gerundet zurückgegeben. Für die Funktion gibt es drei verschiedene Aufrufmöglichkeiten:

1. `abs (z);`  
 $|z|$  wird bez. des Current-Rundungsmodus in die Current-Precision gerundet und zurückgegeben.
2. `abs (z, RoundUp);`  
 $|z|$  wird in die Current-Precision aufgerundet zurückgegeben.
3. `abs (z, RoundDown, prec);`  
 $|z|$  wird in ein Format mit der Präzision `prec` abgerundet zurückgegeben. Die Präzision des zurückgegebenen Wertes wird also i.a. von der voreingestellten Current-Precision verschieden sein!  
Will man jedoch  $|z|$  **rundungsfehlerfrei** mit der gleichen Präzision von  $z$  zurückgeben, so gelingt dies in allen Fällen, unabhängig von der voreingestellten Current-Precision, mit dem Funktionsaufruf:

```
abs (z, RoundNearest, z.GetPrecision());
```

wobei der Rundungsmodus (hier `RoundNearest`) natürlich beliebig gesetzt werden kann. Stimmt die Präzision von  $z$  mit der Current-Precision überein, so liefert der Aufruf `abs (z);` ebenfalls den **rundungsfehlerfreien** Wert von  $|z|$ . In der Praxis wird die rundungsfehlerfreie Rückgabe von  $|z|$  vermutlich immer im Vordergrund stehen.

Die Funktion `abs` kann also sehr flexibel eingesetzt werden und funktioniert nach 1. und 2. wie bei der Standard-Implementierung von Seite 32. Lediglich der Punkt 3. weicht von dieser Standard-Implementierung ab, um den exakten Wert von  $|z|$  in jedem Fall garantieren zu können.

```
MpfrClass arg (const MpfcClass& z, RoundingMode rnd, PrecisionType prec);
```

Mit  $z = x + i \cdot y$  wird  $\arg(z) = \text{atan2}(y,x)$  i.a. gerundet zurückgegeben. Zur Definition von  $\text{atan2}(y,x)$  vgl. auch Seite 38. Für die Funktion gibt es ganz analog zur `abs(z)`-Funktion drei verschiedene Aufrufmöglichkeiten, siehe dazu Seite 90.

```
MpfrClass Re (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  wird  $\text{Re}(z) = x$  ohne Rundung in der Präzision von `z` zurückgegeben.

```
MpfrClass Im (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  wird  $\text{Im}(z) = y$  ohne Rundung in der Präzision von `z` zurückgegeben.

```
void times2pown (MpfcClass& op, long int k, RoundingMode rnd);
```

Die Funktion liefert mit dem Eingabewert `op` den Wert  $op \cdot 2^k$  mit der ursprünglichen Präzision zurück, d.h. Real- und Imaginärteil werden mit  $2^k$  multipliziert. Solange kein Über- oder Unterlauf entsteht, wird  $op \cdot 2^k$  **exakt**, d.h. rundungsfehlerfrei berechnet. Tritt jedoch z.B. ein Überlauf ein, so wird gemäß `rnd` gerundet. Wenn `rnd` nicht gesetzt wird, so erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Ist dieser nicht gesetzt, so erfolgt die Rundung weg von der Null.

```
void set_nan (MpfcClass& x);
```

Setzt Real- und Imaginärteil von `x` auf NaN, wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_inf (MpfcClass& x, const int k);
```

Setzt `x` auf  $(\pm\text{Inf}, \pm\text{Inf})$ , wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss. Das Vorzeichen beim Real- und Imaginärteil wird durch `k` festgelegt.

```
void set_zero (MpfcClass& x);
```

Setzt `x` auf  $(0,0)$ , wobei die Präzision von `x` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

### 5.10.3. Elementarfunktionen

Tabelle 5.1.: Elementarfunktionen mit  $z = x + i \cdot y$ ,  $a, b$  vom Typ `MpfcClass`

Funktion	Aufruf	Funktion	Aufruf
$\text{conj}(z) = x - i \cdot y$	<code>conj(z)</code>	$e^z - 1$	<code>expm1(z)</code>
$\text{Re}(z) = x$	<code>Re(z)</code>	$e^{z^2}$	<code>expx2(z)</code>
$\text{Im}(z) = y$	<code>Im(z)</code>	$e^{z^2} - 1$	<code>expx2m1(z)</code>
$ z $	<code>abs(z)</code>	$e^{-z^2}$	<code>expmx2(z)</code>
$\arg(z)$	<code>arg(z)</code>	$e^{-z^2} - 1$	<code>expmx2m1(z)</code>
$\arg(1 + z)$	<code>argp1(z)</code>	$2^z$	<code>exp2(z)</code>
$z^2$	<code>sqr(z)</code>	$10^z$	<code>exp10(z)</code>
$z^2 + a \cdot z + b$	<code>poly2(z,a,b)</code>	$z^n, n \in \mathbb{Z}$	<code>power(z,n)</code>
$1/z$	<code>reci(z)</code>	$z^r, r \in \mathbb{R}$	<code>pow(z,r)</code>
$1/z^2$	<code>reci_z2(z)</code>	$z^w, w \in \mathbb{C}$	<code>pow(z,w)</code>
$1/(1 + z^2)$	<code>reci_1pz2(z)</code>	$\sin(z)$	<code>sin(z)</code>
$1/(1 - z^2)$	<code>reci_1mz2(z)</code>	$\cos(z)$	<code>cos(z)</code>
$\sqrt{z}$	<code>sqrt(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z}$	<code>sqrt_all(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{z+1} - 1$	<code>sqrtp1m1(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$\sqrt{1 + z^2}$	<code>sqrt1px2(z)</code>	$\arccos(z)$	<code>acos(z)</code>
$\sqrt{1 - z^2}$	<code>sqrt1mx2(z)</code>	$\arctan(z)$	<code>atan(z)</code>
$\sqrt{z^2 - 1}$	<code>sqrtx2m1(z)</code>	$\text{arccot}(z)$	<code>acot(z)</code>
$\sqrt{-z^2 - 1}$	<code>sqrtnx2m1(z)</code>	$\sinh(z)$	<code>sinh(z)</code>
$\sqrt[n]{z}, n \in \mathbb{Z}$	<code>sqrtn(z,n)</code>	$\cosh(z)$	<code>cosh(z)</code>
$\sqrt[n]{z}, n \in \mathbb{Z}$	<code>sqrtn_all(z,n)</code>	$\tanh(z)$	<code>tanh(z)</code>
$\log(z)$	<code>ln(z)</code>	$\text{coth}(z)$	<code>coth(z)</code>

Fortsetzung auf der nächsten Seite

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\log(1 + z)$	<code>lnp1(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$\log_2(z)$	<code>log2(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$\log_{10}(z)$	<code>log10(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$e^z$	<code>exp(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>

### Anmerkungen:

1. Bei den Funktionen `Re`, `Im` und `conj` werden die Rückgabewerte rundungsfehlerfrei in der Präzision des jeweiligen Eingabewertes zurückgegeben. Bei allen anderen Funktionen aus der Tabelle kann mit dem Rundungsparameter `rnd` entsprechend gerundet werden. Wird dieser Parameter nicht gesetzt, so wird nach dem Current-Rundungsmodus gerundet. Ist dieser nicht gesetzt, so erfolgt die Rundung mit `RoundNearest`.
2. Bei den Funktionen `abs`, `arg` und `argp1` kann mit dem zusätzlichen Parameter `prec` die Präzision festgelegt werden, in den der Rückgabewert zu runden ist. Ohne `prec` wird in die Current-Präzision gerundet.
3. Bei den mehrdeutigen Funktionen, wie z.B. `ln(z)`, `asin(z)` oder `atan(z)`, werden nur die Funktionswerte des Hauptzweiges berechnet.

Das nachfolgende Programm MPFR-10.cpp berechnet in einer Liste alle abgerundeten dritten Wurzeln aus  $z = -1 + i$  und gibt den Inhalt dieser Liste auf dem Bildschirm aus.

```

1 // MPFR-10.cpp
2 #include "mpfclass.hpp"
3
4 using namespace MPFR;
5 using namespace cxsc;
6 using namespace std;
7
8 int main(void)
9 {
10  MPFR::MpfrClass::SetCurrPrecision(70);
11  cout << "GetCurrPrecision() = " << MPFR::MpfrClass::GetCurrPrecision() << endl;
12  int n = 3; // Alle n-ten Wurzeln gerundet berechnen
13  real re(-1), im(1);
14  MpfcClass z(re, im, RoundNearest, 53);
15  cout.precision(70/3.321928095); // Ausgabe mit 21 Dez.-Stellen
16  cout << "z = " << z << endl;
17  cout << "z.GetPrecision() = " << z.GetPrecision() << endl;
18  cout << "Berechnung aller " << n << "-ten Wurzeln aus z" << endl;
19
20  list<MpfcClass> res;
21  res = sqrt_all(z, n, RoundDown);
22
23  list<MpfcClass>::iterator pos;
24  // Ausgabe der n n-ten Wurzeln:
25  for (pos = res.begin(); pos != res.end(); ++pos )
26  {
27    cout << *pos << endl; // Jede Einschliessung in neue Zeile
28    cout << "Praezision = " << (*pos).GetPrecision() << endl;
29  }
30
31  return 0;
32 }

```

Das Programm liefert die Ausgabe

```

GetCurrPrecision() = 70

z = (-1.00000000000000000000, 1.00000000000000000000)

z.GetPrecision() = 53
Berechnung aller 3-ten Wurzeln aus z
(7.93700525984099737374e-1, 7.93700525984099737374e-1)
Praezision = 70
(-1.08421508149135118188, 2.90514555507251444495e-1)
Praezision = 70
(2.90514555507251444499e-1, -1.08421508149135118188)
Praezision = 70

```

mit den drei abgerundeten (RoundDown) dritten Wurzeln aus  $z = -1 + i$ .

## 6. MpfcClass-Interface für komplexe Langzahl-Intervallrechnungen in C-XSC

### 6.1. Grundlegendes

Das MpfcClass-Interface ist eine in `mpfciclass.hpp`, `mpfciclass.cpp` implementierte C++-Wrapper-Klasse `MpfcClass` für die C-Bibliotheken MPFR und MPFI, deren C-Funktionen über die implementierten Operatoren und Funktionen aufgerufen werden. Alle Funktionen mit der möglichen Übergabe eines `PrecisionTypes` besitzen als Standard den Wert von `CurrPrecision`, der beliebig gesetzt werden kann. Dies gilt auch für alle Konstruktoren.

#### 6.1.1. Allgemein

Um in C-XSC das Interface verwenden zu können, muss der Header `mpfciclass.hpp` eingebunden werden. In ihm sind die benötigten Header-Dateien der MPFI-Bibliothek enthalten. Die MpfcClass-Klasse liegt im Namensraum "MPFI".

#### 6.1.2. Aufbau

Die Klasse besteht intern aus zwei "mpfi\_t"-Variablen. Diese dienen zum Speichern der Real- und Imaginärteil-Intervalle. Zusätzlich gibt es ein static Element, um die aktuelle Basis für die Ein- und Ausgabe zu speichern.

#### 6.1.3. Präzision

Die Präzision gibt die Anzahl der binären Mantissenstellen einer MpfcClass-Variablen an. Die Real- und Imaginärteil-Intervalle erhalten stets die gleiche Präzision. Der Wert `prec` muss mindestens 2 betragen. Die Current-Precision kann global gesetzt werden; wenn dies nicht geschieht, so wird mit der Default-Precision von 53 Bits gerechnet. Unabhängig davon kann die Präzision für jede MpfcClass-Variable auch einzeln festgelegt werden.

#### 6.1.4. Variablentyp `PrecisionType`

Mithilfe des Variablentyps `PrecisionType` (Name der Variablen meist `prec`) kann der Präzisionswert einer MpfcClass-Variablen eingestellt werden. Der Variablentyp ist ein typedef für eine "mp\_prec\_t"-Variable der MPFR-Bibliothek.

## 6.2. Rundungs und Precision Handling

```
PrecisionType GetPrecision () const;
```

Diese Memberfunktion gibt für das aktuelle Objekt die maximale Präzision seines Real- und Imaginärteil-Intervalls in Bits zurück, wobei hier für beide Intervalle sogar verschiedene Präzisionen erlaubt sind. 302 Bits entsprechen dabei  $302/\log_2(10) \approx 91$  Dezimalstellen.

```
void SetPrecision (PrecisionType prec);
```

Diese Memberfunktion setzt die Präzision des aktuellen Objektes auf `prec`. Die Intervalle für Real- und Imaginärteil bleiben dabei nicht erhalten.

```
void RoundPrecision (PrecisionType prec);
```

Diese Memberfunktion schließt die ursprünglichen Real- und Imaginärteil-Intervalle des aktuellen Objekts durch entsprechende Intervalle mit der neuen Präzision `prec` ein.

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision in `MpfrClass`, `MpfiClass`, `MpfcClass`, `Mpfciclass` auf `prec`. Wird die Current-Precision nicht gesetzt, so kommt die Default-Precision von 53 Bits zur Anwendung. Durch das Setzen der Current-Precision werden die Präzisionen der bis dahin benutzten Variablen jedoch **nicht** geändert.

```
const int MpfrClass::GetBase ();
```

Gibt die aktuelle Basis zurück, diese hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die Funktion ist definiert in `mpfrclass.cpp`.

```
void MpfrClass::SetBase (int b);
```

Setzt die aktuelle Basis auf `b`. Dies muss ein Wert zwischen 2 und 36 sein. Die Basis hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die internen Rechnungen erfolgen stets im Binärsystem! Die Funktion ist definiert in `mpfrclass.cpp`.



## 6.3. Konstruktoren / Destruktor

### 6.3.1. Konstruktoren

```
Mpfciclass ();
```

Der Default-Konstruktor legt ein neues Element mit der Current-Precision an.

Mit `Mpfciclass y;` initialisiert man den Wert:  $y = ([NaN, NaN], [NaN, NaN])$ .

```
Mpfciclass (const Mpfciclass& x, PrecisionType prec);
Mpfciclass (const cxsc::cinterval& x, PrecisionType prec);
Mpfciclass (const mpfi_t& x, const mpfi_t& y, PrecisionType prec);
Mpfciclass (const mpfi_t& x, PrecisionType prec);
Mpfciclass (const Mpfciclass& x, const Mpfciclass& y, PrecisionType prec);
Mpfciclass (const Mpfciclass& x, PrecisionType prec);

Mpfciclass (const cxsc::interval& x, const cxsc::interval& y, PrecisionType prec);

Mpfciclass (const cxsc::interval& x, PrecisionType prec);
Mpfciclass (const MPFR::MpfcClass& x, PrecisionType prec);
Mpfciclass (const cxsc::complex& x, PrecisionType prec);
Mpfciclass (const mpfr_t& x, const mpfr_t& y, PrecisionType prec);
Mpfciclass (const mpfr_t& x, PrecisionType prec);
Mpfciclass (const MpfrClass& x, const MpfrClass& y, PrecisionType prec);
Mpfciclass (const MpfrClass& x, PrecisionType prec);
Mpfciclass (const double& x, const double& y, PrecisionType prec);
Mpfciclass (const double& x, PrecisionType prec);

Mpfciclass (const cxsc::real& x, const cxsc::real& y, PrecisionType prec);

Mpfciclass (const cxsc::real& x, PrecisionType prec);
Mpfciclass (const int& x, const int& y, PrecisionType prec);
Mpfciclass (const int& x, PrecisionType prec);
```

Mit den Konstruktoraufrufen `Mpfciclass w(z);` oder `Mpfciclass w(x, y);` werden die Objekte  $z$  bzw.  $x, y$  durch  $w$  mit der Current-Präzision eingeschlossen. Mit dem zusätzlichen Parameter `prec` lassen sich die Objekte auch mit einer anderen Präzision einschließen. Die Real- und Imaginärteil-Intervalle von  $w$  erhalten stets die gleiche Präzision.

Mögliche Konstruktor-Aufrufe sind:

1. `Mpfciclass w(z); Mpfciclass w(x, y);`
2. `Mpfciclass w(z, 3); Mpfciclass w(x, y, 3);`

**Zu 1.** Die Objekte  $z, x, y$  werden mit der Current-Precision durch die Real- und Imaginärteil-Intervalle des Objekts  $w$  eingeschlossen.

**Zu 2.** Die Objekte  $z, x, y$  werden durch die Real- und Imaginärteil-Intervalle des Objekts  $w$  mit der (sehr kleinen) Präzision `prec = 3` entsprechend grob eingeschlossen.

Die oberen 19 Konstruktoren erlauben also eine sehr flexible Initialisierung von `Mpfciclass`-Objekten.

```
Mpfciclass::Mpfciclass (const std::string& s, PrecisionType prec);
```

Der Aufruf `Mpfciclass w(s);` liefert ein Objekt `w` vom Typ `Mpfciclass` mit Real- und Imaginärteil-Intervallen, welche die entsprechenden Intervalle des String `s` in der Current-Precision garantiert einschließen.

Der Aufruf `Mpfciclass w(s, 140);` liefert ein Objekt `w` vom Typ `Mpfciclass` mit Real- und Imaginärteil-Intervallen der gleichen Präzision `prec = 140`, welche die entsprechenden Intervalle des String `s` garantiert einschließen. Der Präzision von 140 Bits entsprechen dabei  $140/\log_2(10) = 140/3.32192809\dots \approx 42$  Dezimalstellen. Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in eine binäre Zahl konvertiert werden kann, so dass daher bei beiden Konstruktor-Aufrufen Rundungen nach außen i.a. nicht zu vermeiden sind.

Der String `s` muss das Format `([Number,Number],[Number,Number])` besitzen.

### 6.3.2. Destruktor

```
~Mpfciclass ();
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

## 6.4. Zuweisungs-Operatoren

Unabhängig von der Präzision des rechten Operanden erhalten bei allen folgenden Zuweisungsoperatoren die Real- und Imaginärteil-Intervalle des linken Operanden `y` stets die Current-Precision und schließen die Zahlen oder Intervalle des jeweiligen rechten Operanden `op` optimal ein.

```
MpfcIClass& operator = (const MpfcIClass& op);
MpfcIClass& operator = (const cxsc::cinterval& op);
MpfcIClass& operator = (const MPFR::MpfcClass& op);
MpfcIClass& operator = (const cxsc::complex& op);
MpfcIClass& operator = (const mpfi_t& op);
MpfcIClass& operator = (const MPFI::MpfiClass& op);
MpfcIClass& operator = (const cxsc::interval& op);
MpfcIClass& operator = (const mpfr_t& op);
MpfcIClass& operator = (const MPFR::MpfrClass& op);
MpfcIClass& operator = (const double& op);
MpfcIClass& operator = (const cxsc::real& op);
MpfcIClass& operator = (const int& op);
MpfcIClass& operator = (const std::string& op);
```

### Anmerkungen:

- Ist z.B. `op` ein `real`-Wert und wurde die Current-Precision mit `SetCurrPrecision` zu klein gewählt, so ist das Realteil-Intervall des linken Operanden `y` i.a. kein Punktintervall. Ist die Current-Precision jedoch größer oder gleich 53, so ist das einschließende Realteil-Intervall `y` stets ein Punktintervall, das `op` optimal einschließt.
- Ist `op` ein String, so muss dieser die Form `([Number,Number], [Number,Number])` haben. Ist z.B. `op = ([0.1,0.1], [-1.1,-1.1])` so können die Real- und Imaginärteil-Intervalle von `y` bei noch so großer Current-Precision keine Punktintervalle sein, da 0.1 und 1.1 im vorliegenden Binärsystem nicht exakt darstellbar sind.

## 6.5. Eingabe / Ausgabe

```
std::ostream& operator << (std::ostream& os, const MpfcClass& z);
```

Ermöglicht die Ausgabe einer MpfcClass-Variablen  $z$  über den Standard-Ausgabestrom "cout". Die Anzahl der Nachkommastellen ist identisch mit dem Wert, der in `cout.precision()` eingestellt ist. Für die dezimale Ausgabe liefert

```
cout.precision(z.GetPrecision()/3.32192809...);
```

die entsprechende Dezimalstellenzahl. Die Real- und Imaginärteil-Intervalle von  $z$  werden in der mit `SetBase(k)` voreingestellten Basis im Format

```
( [Number,Number] , [Number,Number] )
```

ausgegeben, wobei die Real- und Imaginärteil-Intervalle von  $z$  beide durch die Ausgabeintervalle mit der durch `cout.precision(...)` voreingestellten Ausgabe-Präzision eingeschlossen werden.

Für die meist dezimale Ausgabe ist natürlich mit `SetBase(10)` die richtige Ausgabebasis  $k = 10$  zu wählen.

```
std::istream& operator >> (std::istream& is, MpfcClass& z);
```

Ermöglicht das Einlesen eines MpfcClass-Objekts über den Standard-Eingabestrom "cin" in das Objekt  $z$ . Die eingegebenen Real- und Imaginärteil-Intervalle sind auf keine Stellenzahl begrenzt. Nur das folgende Eingabeformat ist zulässig:

```
( [Number,Number] , [Number,Number] )
```

Die obigen Real- und Imaginärteil-Intervalle werden durch die binären Real- und Imaginärteil-Intervalle des Objekts  $z$  in der Current-Präzision optimal eingeschlossen.

**Vorsicht:** Leerzeichen sind nicht erlaubt, und die runden Klammern müssen beide gesetzt werden, sonst erfolgt eine entsprechende Fehlermeldung.

```
const int MpfcClass::GetBase ();
```

Gibt die aktuelle Basis zurück, diese hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die Funktion ist definiert in `mpfrclass.cpp`.

```
void MpfcClass::SetBase (int b);
```

Setzt die aktuelle Basis auf  $b$ . Dies muss ein Wert zwischen 2 und 36 sein. Die Basis hat nur Einfluss auf die Ein- und Ausgabeoperatoren und auf String-Manipulationen. Die internen Rechnungen erfolgen stets im Binärsystem! Die Funktion ist definiert in `mpfrclass.cpp`.

## 6.6. Typ-Umwandlungen

Um ein flexibles Arbeiten zwischen C-XSC und der Klasse `Mpfciclass` zu ermöglichen, wurden möglichst viele Typ-Umwandlungs-Funktionen bereitgestellt:

### 6.6.1. `real, interval, ...` → `MPFCI`

```
Mpfciclass cinterval2Mpfciclass (const cxsc::cinterval& op);
Mpfciclass mpfi_t2Mpfciclass (const mpfi_t& op);
Mpfciclass Mpfciclass2Mpfciclass (const Mpfciclass& op);
Mpfciclass interval2Mpfciclass (const cxsc::interval& op);
Mpfciclass mpfr_t2Mpfciclass (const mpfr_t& op);
Mpfciclass Mpfciclass2Mpfciclass (const MPFR::Mpfciclass::Mpfciclass& op)
Mpfciclass real2Mpfciclass (const cxsc::real& op);
Mpfciclass double2Mpfciclass (const double& op);
Mpfciclass int2Mpfciclass (const int& op);
```

Obige Funktionen liefern mit dem jeweiligen Eingabewert `op` einen Rückgabewert vom Typ `Mpfciclass` in einer Präzision, die gewährleistet, dass der Rückgabewert **genau** dem Wert von `op` entspricht. Die Präzision des Rückgabewertes wird also i.a. nicht mit der Current-Precision übereinstimmen! Die obigen zwölf Funktionen kommen u.a. bei den Vergleichsoperatoren zur Anwendung.

### 6.6.2. `MPFCI` → `cinterval`

Die folgende Funktion liefert mit einem Objekt `op` vom Typ `Mpfciclass` eine i.a. **gerundete** Einschließung von `op` vom C-XSC Typ `cinterval`;

```
cxsc::cinterval to_cinterval(const Mpfciclass& op);
```

Eine optimale Einschließung von `op` wird nur erreicht, wenn `op` im IEEE-double-Format darstellbar ist.

### 6.6.3. `MPFCI` → `mpfi_t`

```
mpfi_t& Mpfciclass::GetValueRe();
```

Die Memberfunktion `GetValueRe()` liefert für das jeweils aktuelle Objekt vom Typ `Mpfciclass` eine Referenz auf sein Realteilintervall `mpfi_re` vom Typ `mpfi_t`. Damit kann dieses Realteilintervall z.B. an einen nicht-konstanten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
mpfr_t& Mpfciclass::GetValueIm();
```

Die Memberfunktion `GetValueIm()` liefert für das jeweils aktuelle Objekt vom Typ `Mpfciclass` eine Referenz auf sein Imaginärteilintervall `mpfi_im` vom Typ `mpfi_t`. Damit kann dieses Imaginärteilintervall z.B. an einen nicht-konstanten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
const mpfi_t& getvalueRe(const Mpfciclass& z)
```

Die Funktion `getvalueRe(...)` liefert für das als `const` definierte Objekt `z` vom Typ `Mpfciclass` eine Referenz auf sein Realteilintervall `mpfi_re` vom Typ `mpfi_t`. Damit kann dieses Realteilintervall z.B. an einen als `const` deklarierten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

```
const mpfi_t& getvalueIm(const Mpfciclass& z)
```

`getvalueIm(...)` liefert für das als `const` definierte Objekt `z` vom Typ `Mpfciclass` eine Referenz auf sein Imaginärteilintervall `mpfi_im` vom Typ `mpfi_t`. Damit kann dieses Imaginärteilintervall z.B. an einen als `const` deklarierten Referenz-Parameter gleichen Typs in einer Parameterliste an eine Funktion übergeben werden.

#### 6.6.4. mpfi\_t → MPFCI

```
void Mpfciclass::SetValueRe(const mpfi_t& v)
```

Mit dieser Memberfunktion wird das Realteilintervall des aktuellen Objekts exakt auf  $v$  gesetzt. Das Imaginärteilintervall des aktuellen Objekts bleibt dabei erhalten. Die Präzision des aktuellen Objekts wird gesetzt auf das Maximum der Präzisionen von  $v$  und `mpfi_im`, d.h. die Präzisionen von Real- und Imaginärteilintervall sind nach dem Funktionsaufruf wieder gleich.

```
void Mpfciclass::SetValueIm(const mpfi_t& v)
```

Mit dieser Memberfunktion wird das Imaginärteilintervall des aktuellen Objekts exakt auf  $v$  gesetzt. Das Realteilintervall des aktuellen Objekts bleibt dabei erhalten. Die Präzision des aktuellen Objekts wird gesetzt auf das Maximum der Präzisionen von  $v$  und `mpfi_re`, d.h. die Präzisionen von Real- und Imaginärteilintervall sind nach dem Funktionsaufruf wieder gleich.

```
void Mpfciclass::SetValue(const mpfi_t& re, const mpfi_t& im)
```

Mit dieser Memberfunktion werden Real- und Imaginärteilintervall des aktuellen Objekts exakt auf  $re$  bzw.  $im$  gesetzt. Das aktuelle Objekt erhält als Präzision das Maximum der Präzisionen von  $re$  und  $im$ , d.h. die Präzisionen von Real- und Imaginärteilintervall sind nach dem Funktionsaufruf wieder gleich.

#### 6.6.5. string → MPFCI

```
Mpfciclass string2Mpfciclass(const std::string& op, PrecisionType prec);
```

Der Aufruf `string2Mpfciclass(op);` liefert ein Objekt vom Typ `Mpfciclass` dessen Real- und Imaginärteil-Intervalle mit der voreingestellten Current-Präzision die Real- und Imaginärteil-Intervalle von `op` einschließt.

Der zweite mögliche Aufruf `string2Mpfciclass(op,140);` liefert ganz analog eine Einschließung der String-Intervalle mit einer Präzision von jetzt 140 Bits.

Zu beachten ist, dass ein dezimaler String i.a. nicht rundungsfehlerfrei in ein Binärformat konvertiert werden kann, so dass i.a. Überschätzungen nicht zu vermeiden sind, es werden jedoch stets garantierte Einschließungen von `op` zurückgegeben. Eine weitere Möglichkeit, einen String in ein `Mpfciclass`-Objekt zu verwandeln, besteht in einem entsprechenden Konstruktor-Aufruf, vgl. Seite 98. Dort findet man auch Hinweise auf die möglichen Intervall-String-Formate.

#### 6.6.6. MPFCI → string

```
std::string to_string(const Mpfciclass& z, PrecisionType prec);
```

Die Real- und Imaginärteilintervalle des Objekts  $z = x + i \cdot y$  werden durch den zurückgegebenen String `s` mit `prec` Dezimalstellen eingeschlossen, wenn Base gleich 10 ist. Der String besitzt das Format

$$([\text{Number}, \text{Number}], [\text{Number}, \text{Number}])$$

wobei intern keine Leerzeichen auftreten. Wählt man `prec` hinreichend groß, so stellt der String die Real- und Imaginärteil-Intervalle von  $z$  **exakt** dar, weil eine Binärzahl stets exakt in eine Dezimalzahl umgewandelt werden kann.

Wird `prec` nicht angegeben, so werden  $x, y$  durch String-Intervalle eingeschlossen, die bei Base=10 so viele Dezimalstellen besitzen, wie es der Präzision von  $z$  entspricht. Ist die Präzision von  $z$  z.B. 302, so werden im String Real- und Imaginärteil-Intervalle mit  $302/\log_2(10) = 302/3.32192809\dots \approx 91$  Dezimalstellen generiert.

### 6.6.7. MPFCI → MPFCI

Beachten Sie bitte, dass bei einer Wertzuweisung an eine `Mpfciclass`-Variable mit Hilfe des Operators `=` der linke Operand mit seinem Real- und Imaginärteilintervall die entsprechenden Real- und Imaginärteilwerte des rechten Operanden mit der Current-Precision stets einschließt, vgl. dazu auch Seite 99. Will man jedoch abweichend von dieser Current-Precision mit einer anderen Präzision einschließen, so kann dies mit folgender Funktion ohne Rückgabewert realisiert werden:

```
void set_Mpfciclass (Mpfciclass& op, const Mpfciclass& op1, PrecisionType prec);
```

Folgende Funktionsaufrufe sind möglich:

1. `set_Mpfciclass (op, op1, prec);`  
`op` schließt die Real- und Imaginärteilintervalle von `op1` mit der Präzision `prec` ein. Setzt man `prec` gleich der Präzision von `op1`, so erhält `op` den **exakten** Wert von `op1`.
2. `set_Mpfciclass (op, op1);`  
`op` schließt die Real- und Imaginärteilintervalle von `op1` mit der Current-Precision ein.

## 6.7. Abfragen

Bei allen folgenden Abfragefunktionen muss die Präzision von `x` nicht mit der Current-Precision übereinstimmen.

```
bool isNaN (const Mpfciclass& x);  
bool isInf (const Mpfciclass& x);  
bool isBounded (const Mpfciclass& x);  
bool isZero (const Mpfciclass& x);  
bool isPoint (const Mpfciclass& x);
```

`isNaN` und `isInf` überprüfen, ob Real- oder Imaginärteil von `x` gleich NaN bzw.  $\pm\text{Inf}$  sind. `isBounded` überprüft, ob Real- und Imaginärteil von `x` ein normales `Mpfciclass` Intervall ungleich NaN und ungleich  $\pm\text{Inf}$  ist. `isZero` überprüft, ob das Real- und Imaginärteil-Intervall von `x` gleich Null ist, und `isPoint` überprüft, ob das Real- und Imaginärteil-Intervall von `x` jeweils ein Punktintervall ist.

## 6.8. Vergleiche

Es werden die üblichen Vergleichsoperatoren implementiert, wobei wenigstens ein Operand vom Typ `MpfcClass` sein muss. Die Operanden können unterschiedliche Präzisionen besitzen.

### 6.8.1. Vergleichsoperatoren `=`, `!=`, `<`, `<=`,

```
bool operator == (const MpfcClass& x, const MpfcClass& y);
bool operator == (const MpfcClass& x, const cxsc::cinterval& y);
bool operator == (const MpfcClass& x, const MPFR::MpfcClass& y);
bool operator == (const MpfcClass& x, const cxsc::complex& y);
bool operator == (const MpfcClass& x, const mpfi_t& y);
bool operator == (const MpfcClass& x, const MpfiClass& y);
bool operator == (const MpfcClass& x, const cxsc::interval& y);
bool operator == (const MpfcClass& x, const mpfr_t& y);
bool operator == (const MpfcClass& x, const MPFR::MpfrClass& y);
bool operator == (const MpfcClass& x, const double& y);
bool operator == (const MpfcClass& x, const cxsc::real& y);
bool operator == (const MpfcClass& x, const int y);
bool operator == (const cxsc::cinterval& x, const MpfcClass& y);
bool operator == (const MPFR::MpfcClass& x, const MpfcClass& y);
bool operator == (const cxsc::complex& x, const MpfcClass& y);
bool operator == (const mpfi_t& x, const MpfcClass& y);
bool operator == (const MpfiClass& x, const MpfcClass& y);
bool operator == (const cxsc::interval& x, const MpfcClass& y);
bool operator == (const mpfr_t& x, const MpfcClass& y);
bool operator == (const MpfrClass& x, const MpfcClass& y);
bool operator == (const double& x, const MpfcClass& y);
bool operator == (const cxsc::real& x, const MpfcClass& y);
bool operator == (const int x, const MpfcClass& y);
```

```
bool operator != (const MpfcClass& x, const MpfcClass& y);
bool operator != (const MpfcClass& x, const cxsc::cinterval& y);
bool operator != (const MpfcClass& x, const MPFR::MpfcClass& y);
bool operator != (const MpfcClass& x, const cxsc::complex& y);
bool operator != (const MpfcClass& x, const mpfi_t& y);
bool operator != (const MpfcClass& x, const MpfiClass& y);
bool operator != (const MpfcClass& x, const cxsc::interval& y);
bool operator != (const MpfcClass& x, const mpfr_t& y);
bool operator != (const MpfcClass& x, const MPFR::MpfrClass& y);
bool operator != (const MpfcClass& x, const double& y);
bool operator != (const MpfcClass& x, const cxsc::real& y);
bool operator != (const MpfcClass& x, const int y);
bool operator != (const cxsc::cinterval& x, const MpfcClass& y);
bool operator != (const MPFR::MpfcClass& x, const MpfcClass& y);
bool operator != (const cxsc::complex& x, const MpfcClass& y);
bool operator != (const mpfi_t& x, const MpfcClass& y);
bool operator != (const MpfiClass& x, const MpfcClass& y);
bool operator != (const cxsc::interval& x, const MpfcClass& y);
bool operator != (const mpfr_t& x, const MpfcClass& y);
bool operator != (const MpfrClass& x, const MpfcClass& y);
bool operator != (const double& x, const MpfcClass& y);
bool operator != (const cxsc::real& x, const MpfcClass& y);
bool operator != (const int x, const MpfcClass& y);
```



```

bool operator < (const MpfcClass& x, const MpfcClass& y);
bool operator < (const MpfcClass& x, const cxsc::cinterval& y);
bool operator < (const MpfcClass& x, const mpfi_t& y);
bool operator < (const MpfcClass& x, const MpfiClass& y);
bool operator < (const MpfcClass& x, const cxsc::interval& y);
bool operator < (const cxsc::cinterval& x, const MpfcClass& y);
bool operator < (const mpfi_t& x, const MpfcClass& y);
bool operator < (const MpfiClass& x, const MpfcClass& y);
bool operator < (const cxsc::interval& x, const MpfcClass& y);

```

'interval' < 'MpfcClass': 'interval' liegt ganz im Innern von 'MpfcClass', wobei  $[0,0] \subset \text{Im}(y)$  erfüllt sein muss.

```

bool operator < (const MPFR::MpfcClass& x, const MpfcClass& y);
bool operator < (const cxsc::complex& x, const MpfcClass& y);
bool operator < (const mpfr_t& x, const MpfcClass& y);
bool operator < (const MPFR::MpfrClass& x, const MpfcClass& y);
bool operator < (const double& x, const MpfcClass& y);
bool operator < (const cxsc::real& x, const MpfcClass& y);
bool operator < (const int x, const MpfcClass& y);

```

'real' < 'MpfcClass': 'real' liegt ganz im Innern von  $\text{Re}(y)$ , wobei 0 im Innern von  $\text{Im}(y)$  liegen muss.

```

bool operator <= (const MpfcClass& x, const MpfcClass& y);
bool operator <= (const MpfcClass& x, const cxsc::cinterval& y);
bool operator <= (const MpfcClass& x, const mpfi_t& y);
bool operator <= (const MpfcClass& x, const MpfiClass& y);
bool operator <= (const MpfcClass& x, const cxsc::interval& y);
bool operator <= (const cxsc::cinterval& x, const MpfcClass& y);
bool operator <= (const mpfi_t& x, const MpfcClass& y);
bool operator <= (const MpfiClass& x, const MpfcClass& y);
bool operator <= (const cxsc::interval& x, const MpfcClass& y);

```

'interval' <= 'MpfcClass': Es gilt  $x \subseteq \text{Re}(y)$  und  $[0,0] \subseteq \text{Im}(y)$ .

```

bool operator <= (const MPFR::MpfcClass& x, const MpfcClass& y);
bool operator <= (const cxsc::complex& x, const MpfcClass& y);
bool operator <= (const mpfr_t& x, const MpfcClass& y);
bool operator <= (const MPFR::MpfrClass& x, const MpfcClass& y);
bool operator <= (const double& x, const MpfcClass& y);
bool operator <= (const cxsc::real& x, const MpfcClass& y);
bool operator <= (const int x, const MpfcClass& y);

```

'real' <= 'MpfcClass': Es gilt  $x \in \text{Re}(y)$  und  $0 \in \text{Im}(y)$ .

```

bool operator > (const MpfcClass& x, const MpfcClass& y);
bool operator > (const MpfcClass& x, const cxsc::cinterval& y);
bool operator > (const MpfcClass& x, const mpfi_t& y);
bool operator > (const MpfcClass& x, const MpfiClass& y);
bool operator > (const MpfcClass& x, const cxsc::interval& y);
bool operator > (const cxsc::cinterval& x, const MpfcClass& y);
bool operator > (const mpfi_t& x, const MpfcClass& y);
bool operator > (const MpfiClass& x, const MpfcClass& y);
bool operator > (const cxsc::interval& x, const MpfcClass& y);

```

'interval' > 'MpfcClass': 'MpfcClass' liegt ganz im Innern von 'interval', wobei  $\text{Im}(y) = [0,0]$  erfüllt sein muss.

```

bool operator > (const MpfcClass& x, const MPFR::MpfcClass& y);
bool operator > (const MpfcClass& x, const cxsc::complex& y);
bool operator > (const MpfcClass& x, const mpfr_t& y);
bool operator > (const MpfcClass& x, const MPFR::MpfrClass& y);
bool operator > (const MpfcClass& x, const double& y);
bool operator > (const MpfcClass& x, const cxsc::real& y);
bool operator > (const MpfcClass& x, const int y);

```

'MpfcClass' > 'real': Es gilt  $y$  liegt ganz im Innern von  $\text{Re}(x)$  und 0 liegt ganz im Innern von  $\text{Im}(x)$ .

```

bool operator >= (const MpfcClass& x, const MpfcClass& y);
bool operator >= (const MpfcClass& x, const cxsc::cinterval& y);
bool operator >= (const MpfcClass& x, const mpfi_t& y);
bool operator >= (const MpfcClass& x, const MpfiClass& y);
bool operator >= (const MpfcClass& x, const cxsc::interval& y);
bool operator >= (const cxsc::cinterval& x, const MpfcClass& y);
bool operator >= (const mpfi_t& x, const MpfcClass& y);
bool operator >= (const MpfiClass& x, const MpfcClass& y);
bool operator >= (const cxsc::interval& x, const MpfcClass& y);

```

'interval' >= 'MpfcClass': Es gilt  $\text{Re}(y) \subset x$  und  $\text{Im}(y) = [0, 0]$ .

```

bool operator >= (const MpfcClass& x, const MPFR::MpfcClass& y);
bool operator >= (const MpfcClass& x, const cxsc::complex& y);
bool operator >= (const MpfcClass& x, const mpfr_t& y);
bool operator >= (const MpfcClass& x, const MPFR::MpfrClass& y);
bool operator >= (const MpfcClass& x, const double& y);
bool operator >= (const MpfcClass& x, const cxsc::real& y);
bool operator >= (const MpfcClass& x, const int y);

```

'MpfcClass' >= 'real': Es gilt  $y \in \text{Re}(y)$  und  $0 \in \text{Im}(y)$ .

```

int in (const MpfcClass& x, const MpfcClass& y);
int in (const cxsc::cinterval& x, const MpfcClass& y);
int in (const MPFR::MpfcClass::MpfcClass& x, const MpfcClass& y);
int in (const cxsc::complex& x, const MpfcClass& y);
int in (const MpfiClass& x, const MpfcClass& y);
int in (const mpfi_t& x, const MpfcClass& y);
int in (const cxsc::interval& x, const MpfcClass& y);
int in (const MPFR::MpfrClass::MpfrClass& x, const MpfcClass& y);
int in (const mpfr_t& x, const MpfcClass& y);
int in (const double& x, const MpfcClass& y);
int in (const cxsc::real& x, const MpfcClass& y);
int in (const int& x, const MpfcClass& y);

```

Zurückgegeben wird die Eins, wenn  $x$  ganz im Innern von  $y$  enthalten ist, sonst wird die Null zurückgegeben. Ist einer der Operanden ein NaN oder sind beide Operanden unbegrenzt, so wird ebenfalls die Null zurückgegeben. Die Präzisionen beider Operanden können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen.

## 6.9. Durchschnitt

Berechnet wird der Durchschnitt zweier komplexer Intervalle. Ist einer der nachfolgenden Operanden eine reelle oder komplexe Zahl, so ist diese als Punktintervall zu interpretieren. Die Präzisionen der Intervalle oder Zahlen können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen. Die Präzision des zurückgegebenen komplexen Durchschnittsintervalls ist das Maximum der Präzisionen beider Operanden und muss daher mit der Current-Präzision **nicht** übereinstimmen. Dadurch wird der Durchschnitt stets **rundungsfehlerfrei** berechnet. Ist der Durchschnitt leer, so wird  $([\text{NaN}, \text{NaN}], [\text{NaN}, \text{NaN}])$  zurückgegeben.

```
MpfcIClass operator & (const MpfcIClass& z, const MpfcIClass& x);
MpfcIClass operator & (const MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass operator & (const MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass operator & (const MpfcIClass& z, const cxsc::complex& x);
MpfcIClass operator & (const MpfcIClass& z, const mpfi_t& x);
MpfcIClass operator & (const MpfcIClass& z, const MpfiClass& x);
MpfcIClass operator & (const MpfcIClass& z, const cxsc::interval& x);
MpfcIClass operator & (const MpfcIClass& z, const mpfr_t& x);
MpfcIClass operator & (const MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass operator & (const MpfcIClass& z, const double& x);
MpfcIClass operator & (const MpfcIClass& z, const cxsc::real& x);
MpfcIClass operator & (const MpfcIClass& z, const int x);

MpfcIClass operator & (const cxsc::cinterval& x, const MpfcIClass& z);
MpfcIClass operator & (const MPFR::MpfcClass& x, const MpfcIClass& z);
MpfcIClass operator & (const cxsc::complex& x, const MpfcIClass& z);
MpfcIClass operator & (const mpfi_t& x, const MpfcIClass& z);
MpfcIClass operator & (const MpfiClass& x, const MpfcIClass& z);
MpfcIClass operator & (const cxsc::interval& x, const MpfcIClass& z);
MpfcIClass operator & (const mpfr_t& x, const MpfcIClass& z);
MpfcIClass operator & (const MPFR::MpfrClass& x, const MpfcIClass& z);
MpfcIClass operator & (const double& x, const MpfcIClass& z);
MpfcIClass operator & (const cxsc::real x, const MpfcIClass& z);
MpfcIClass operator & (const int& x, const MpfcIClass& z);
```

```
MpfcIClass & operator &= (MpfcIClass& z, const MpfcIClass& x);
MpfcIClass & operator &= (MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass & operator &= (MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass & operator &= (MpfcIClass& z, const cxsc::complex& x);
MpfcIClass & operator &= (MpfcIClass& z, const mpfi_t& x);
MpfcIClass & operator &= (MpfcIClass& z, const MpfiClass& x);
MpfcIClass & operator &= (MpfcIClass& z, const cxsc::interval& x);
MpfcIClass & operator &= (MpfcIClass& z, const mpfr_t& x);
MpfcIClass & operator &= (MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass & operator &= (MpfcIClass& z, const double& x);
MpfcIClass & operator &= (MpfcIClass& z, const cxsc::real& x);
MpfcIClass & operator &= (MpfcIClass& z, const int x);
```

Die Anweisung  $z \&= x;$  liefert an  $z$  den Durchschnitt  $(z \& x)$ , wobei die neue Präzision von  $z$  gleich dem Maximum der Präzisionen der ursprünglichen Operanden  $z$  und  $x$  ist. Auch hier wird erreicht, dass der Durchschnitt stets **rundungsfehlerfrei** an  $z$  zurückgegeben wird.

### Hinweis:

Wenn beide Operanden als Präzision die Current-Präzision besitzen, so wird der Durchschnitt stets **rundungsfehlerfrei** in der gleichen Current-Präzision zurückgegeben.

## 6.10. Konvexe Hülle

Berechnet wird die konvexe Hülle zweier komplexer Intervalle. Ist einer der nachfolgenden Operanden eine reelle oder komplexe Zahl, so ist diese jeweils als Punktintervall zu interpretieren. Die Präzisionen der Intervalle oder Zahlen können unterschiedlich sein und müssen mit der Current-Präzision nicht übereinstimmen. Die Präzision des zurückgegebenen Hüllen-Intervalls ist das Maximum der Präzisionen beider Operanden und muss mit der Current-Präzision **nicht** übereinstimmen. Dadurch wird die konvexe Hülle stets **rundungsfehlerfrei** berechnet.

```

MpfciClass operator | (const MpfciClass& z, const MpfciClass& w);
MpfciClass operator | (const MpfciClass& z, const cxsc::cinterval& x);
MpfciClass operator | (const MpfciClass& z, const MPFR::MpfcClass& x);
MpfciClass operator | (const MpfciClass& z, const cxsc::complex& x);
MpfciClass operator | (const MpfciClass& z, const mpfi_t& x);
MpfciClass operator | (const MpfciClass& z, const MpfiClass& x);
MpfciClass operator | (const MpfciClass& z, const cxsc::interval& x);
MpfciClass operator | (const MpfciClass& z, const mpfr_t& x);
MpfciClass operator | (const MpfciClass& z, const MPFR::MpfrClass& x);
MpfciClass operator | (const MpfciClass& z, const double& x);
MpfciClass operator | (const MpfciClass& z, const cxsc::real& x);
MpfciClass operator | (const MpfciClass& z, const int x);

MpfciClass operator | (const cxsc::cinterval& x, const MpfciClass& z);
MpfciClass operator | (const MPFR::MpfcClass& x, const MpfciClass& z);
MpfciClass operator | (const cxsc::complex& x, const MpfciClass& z);
MpfciClass operator | (const mpfi_t& x, const MpfciClass& z);
MpfciClass operator | (const MpfiClass& x, const MpfciClass& z);
MpfciClass operator | (const cxsc::interval& x, const MpfciClass& z);
MpfciClass operator | (const mpfr_t& x, const MpfciClass& z);
MpfciClass operator | (const MPFR::MpfrClass& x, const MpfciClass& z);
MpfciClass operator | (const double& x, const MpfciClass& z);
MpfciClass operator | (const cxsc::real& x, const MpfciClass& z);
MpfciClass operator | (const int x, const MpfciClass& z);

```

```

MpfciClass & operator |= (MpfciClass& z, const MpfciClass& x);
MpfciClass & operator |= (MpfciClass& z, const cxsc::cinterval& x);
MpfciClass & operator |= (MpfciClass& z, const MPFR::MpfcClass& x);
MpfciClass & operator |= (MpfciClass& z, const cxsc::complex& x);
MpfciClass & operator |= (MpfciClass& z, const mpfi_t& x);
MpfciClass & operator |= (MpfciClass& z, const MpfiClass& x);
MpfciClass & operator |= (MpfciClass& z, const cxsc::interval& x);
MpfciClass & operator |= (MpfciClass& z, const mpfr_t& x);
MpfciClass & operator |= (MpfciClass& z, const MPFR::MpfrClass& x);
MpfciClass & operator |= (MpfciClass& z, const double& x);
MpfciClass & operator |= (MpfciClass& z, const cxsc::real& x);
MpfciClass & operator |= (MpfciClass& z, const int x);

```

Die Anweisung `z |= x;` liefert an `z` die konvexe Hülle (`z | x`), wobei die neue Präzision von `z` gleich dem Maximum der Präzisionen der ursprünglichen Operanden `z` und `x` ist. Auch hier wird erreicht, dass die konvexe Hülle stets **rundungsfehlerfrei** an `z` zurückgegeben wird.

### Hinweis:

Wenn beide Operanden als Präzision die Current-Präzision besitzen, so wird die konvexe Hülle stets **rundungsfehlerfrei** in der gleichen Current-Präzision zurückgegeben.

## 6.11. Arithmetische Operatoren

Für alle arithmetischen Operationen mit komplexen Intervall-Operanden gilt:

Das exakte Ergebnis einer arithmetischen Operation mit komplexen Intervall-Operanden wird unabhängig von der Präzision dieser Operanden mit der vor-eingestellten Current-Präzision außer bei Multiplikation und Division optimal eingeschlossen.

Die Operatoren  $\odot =$ , mit  $\odot \in \{+, -, \cdot, /\}$ , bedeuten:  $u \odot = v \iff u = u \odot v$ . Dabei wird  $u \odot v$  durch  $u$  stets eingeschlossen, wobei  $u$  als Präzision die Current-Präzision erhält, d.h. die Präzision von  $u$  kann sich ändern.

### 6.11.1. Addition

```
MpfcIClass operator + (const MpfcIClass& z, const MpfcIClass& w);
MpfcIClass operator + (const MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass operator + (const MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass operator + (const MpfcIClass& z, const cxsc::complex& x);
MpfcIClass operator + (const MpfcIClass& z, const mpfi_t& x);
MpfcIClass operator + (const MpfcIClass& z, const MpfiClass& x);
MpfcIClass operator + (const MpfcIClass& z, const cxsc::interval& x);
MpfcIClass operator + (const MpfcIClass& z, const mpfr_t& x);
MpfcIClass operator + (const MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass operator + (const MpfcIClass& z, const double& x);
MpfcIClass operator + (const MpfcIClass& z, const cxsc::real& x);
MpfcIClass operator + (const MpfcIClass& z, const int x);
```

```
MpfcIClass operator + (const cxsc::cinterval& x, const MpfcIClass& z);
MpfcIClass operator + (const MPFR::MpfcClass& x, const MpfcIClass& z);
MpfcIClass operator + (const cxsc::complex& x, const MpfcIClass& z);
MpfcIClass operator + (const mpfi_t& x, const MpfcIClass& z);
MpfcIClass operator + (const MpfiClass& x, const MpfcIClass& z);
MpfcIClass operator + (const cxsc::interval& x, const MpfcIClass& z);
MpfcIClass operator + (const mpfr_t& x, const MpfcIClass& z);
MpfcIClass operator + (const MPFR::MpfrClass& x, const MpfcIClass& z);
MpfcIClass operator + (const double& x, const MpfcIClass& z);
MpfcIClass operator + (const cxsc::real& x, const MpfcIClass& z);
MpfcIClass operator + (const int x, const MpfcIClass& z);
```

```
MpfcIClass & operator += (MpfcIClass& z, const MpfcIClass& x);
MpfcIClass & operator += (MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass & operator += (MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass & operator += (MpfcIClass& z, const cxsc::complex& x);
MpfcIClass & operator += (MpfcIClass& z, const mpfi_t& x);
MpfcIClass & operator += (MpfcIClass& z, const MpfiClass& x);
MpfcIClass & operator += (MpfcIClass& z, const cxsc::interval& x);
MpfcIClass & operator += (MpfcIClass& z, const mpfr_t& x);
MpfcIClass & operator += (MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass & operator += (MpfcIClass& z, const double& x);
MpfcIClass & operator += (MpfcIClass& z, const cxsc::real& x);
MpfcIClass & operator += (MpfcIClass& z, const int x);
```

## 6.11.2. Subtraktion

Beachten Sie die Bemerkungen auf Seite 109 oben.

```
MpfcIClass operator - (const MpfcIClass& z, const MpfcIClass& w);
MpfcIClass operator - (const MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass operator - (const MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass operator - (const MpfcIClass& z, const cxsc::complex& x);
MpfcIClass operator - (const MpfcIClass& z, const mpfi_t& x);
MpfcIClass operator - (const MpfcIClass& z, const MpfiClass& x);
MpfcIClass operator - (const MpfcIClass& z, const cxsc::interval& x);
MpfcIClass operator - (const MpfcIClass& z, const mpfr_t& x);
MpfcIClass operator - (const MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass operator - (const MpfcIClass& z, const double& x);
MpfcIClass operator - (const MpfcIClass& z, const cxsc::real& x);
MpfcIClass operator - (const MpfcIClass& z, const int x);
```

```
MpfcIClass operator - (const cxsc::cinterval& x, const MpfcIClass& z);
MpfcIClass operator - (const MPFR::MpfcClass& x, const MpfcIClass& z);
MpfcIClass operator - (const cxsc::complex& x, const MpfcIClass& z);
MpfcIClass operator - (const mpfi_t& x, const MpfcIClass& z);
MpfcIClass operator - (const MpfiClass& x, const MpfcIClass& z);
MpfcIClass operator - (const cxsc::interval& x, const MpfcIClass& z);
MpfcIClass operator - (const mpfr_t& x, const MpfcIClass& z);
MpfcIClass operator - (const MPFR::MpfrClass& x, const MpfcIClass& z);
MpfcIClass operator - (const double& x, const MpfcIClass& z);
MpfcIClass operator - (const cxsc::real& x, const MpfcIClass& z);
MpfcIClass operator - (const int x, const MpfcIClass& z);
```

```
MpfcIClass & operator -= (MpfcIClass& z, const MpfcIClass& x);
MpfcIClass & operator -= (MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass & operator -= (MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass & operator -= (MpfcIClass& z, const cxsc::complex& x);
MpfcIClass & operator -= (MpfcIClass& z, const mpfi_t& x);
MpfcIClass & operator -= (MpfcIClass& z, const MpfiClass& x);
MpfcIClass & operator -= (MpfcIClass& z, const cxsc::interval& x);
MpfcIClass & operator -= (MpfcIClass& z, const mpfr_t& x);
MpfcIClass & operator -= (MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass & operator -= (MpfcIClass& z, const double& x);
MpfcIClass & operator -= (MpfcIClass& z, const cxsc::real& x);
MpfcIClass & operator -= (MpfcIClass& z, const int x);
```



### 6.11.3. Multiplikation

Beachten Sie die Bemerkungen auf Seite 109 oben.

```
Mpfciclass operator * (const Mpfciclass& z, const Mpfciclass& w);
Mpfciclass operator * (const Mpfciclass& z, const cxsc::cinterval& x);
Mpfciclass operator * (const Mpfciclass& z, const MPFR::Mpfciclass& x);
Mpfciclass operator * (const Mpfciclass& z, const cxsc::complex& x);
Mpfciclass operator * (const Mpfciclass& z, const mpfi_t& x);
Mpfciclass operator * (const Mpfciclass& z, const Mpfciclass& x);
Mpfciclass operator * (const Mpfciclass& z, const cxsc::interval& x);
Mpfciclass operator * (const Mpfciclass& z, const mpfr_t& x);
Mpfciclass operator * (const Mpfciclass& z, const MPFR::Mpfrclass& x);
Mpfciclass operator * (const Mpfciclass& z, const double& x);
Mpfciclass operator * (const Mpfciclass& z, const cxsc::real& x);
Mpfciclass operator * (const Mpfciclass& z, const int x);
```

```
Mpfciclass operator * (const cxsc::cinterval& x, const Mpfciclass& z);
Mpfciclass operator * (const MPFR::Mpfciclass& x, const Mpfciclass& z);
Mpfciclass operator * (const cxsc::complex& x, const Mpfciclass& z);
Mpfciclass operator * (const mpfi_t& x, const Mpfciclass& z);
Mpfciclass operator * (const Mpfciclass& x, const Mpfciclass& z);
Mpfciclass operator * (const cxsc::interval& x, const Mpfciclass& z);
Mpfciclass operator * (const mpfr_t& x, const Mpfciclass& z);
Mpfciclass operator * (const MPFR::Mpfrclass& x, const Mpfciclass& z);
Mpfciclass operator * (const double& x, const Mpfciclass& z);
Mpfciclass operator * (const cxsc::real& x, const Mpfciclass& z);
Mpfciclass operator * (const int x, const Mpfciclass& z);
```

```
Mpfciclass & operator *= (Mpfciclass& z, const Mpfciclass& x);
Mpfciclass & operator *= (Mpfciclass& z, const cxsc::cinterval& x);
Mpfciclass & operator *= (Mpfciclass& z, const MPFR::Mpfciclass& x);
Mpfciclass & operator *= (Mpfciclass& z, const cxsc::complex& x);
Mpfciclass & operator *= (Mpfciclass& z, const mpfi_t& x);
Mpfciclass & operator *= (Mpfciclass& z, const Mpfciclass& x);
Mpfciclass & operator *= (Mpfciclass& z, const cxsc::interval& x);
Mpfciclass & operator *= (Mpfciclass& z, const mpfr_t& x);
Mpfciclass & operator *= (Mpfciclass& z, const MPFR::Mpfrclass& x);
Mpfciclass & operator *= (Mpfciclass& z, const double& x);
Mpfciclass & operator *= (Mpfciclass& z, const cxsc::real& x);
Mpfciclass & operator *= (Mpfciclass& z, const int x);
```

## 6.11.4. Division

Beachten Sie die Bemerkungen auf Seite 109 oben.

```
MpfcIClass operator / (const MpfcIClass& z, const MpfcIClass& w);
MpfcIClass operator / (const MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass operator / (const MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass operator / (const MpfcIClass& z, const cxsc::complex& x);
MpfcIClass operator / (const MpfcIClass& z, const mpfi_t& x);
MpfcIClass operator / (const MpfcIClass& z, const MpfiClass& x);
MpfcIClass operator / (const MpfcIClass& z, const cxsc::interval& x);
MpfcIClass operator / (const MpfcIClass& z, const mpfr_t& x);
MpfcIClass operator / (const MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass operator / (const MpfcIClass& z, const double& x);
MpfcIClass operator / (const MpfcIClass& z, const cxsc::real& x);
MpfcIClass operator / (const MpfcIClass& z, const int x);
```

```
MpfcIClass operator / (const cxsc::cinterval& x, const MpfcIClass& z);
MpfcIClass operator / (const MPFR::MpfcClass& x, const MpfcIClass& z);
MpfcIClass operator / (const cxsc::complex& x, const MpfcIClass& z);
MpfcIClass operator / (const mpfi_t& x, const MpfcIClass& z);
MpfcIClass operator / (const MpfiClass& x, const MpfcIClass& z);
MpfcIClass operator / (const cxsc::interval& x, const MpfcIClass& z);
MpfcIClass operator / (const mpfr_t& x, const MpfcIClass& z);
MpfcIClass operator / (const MPFR::MpfrClass& x, const MpfcIClass& z);
MpfcIClass operator / (const double& x, const MpfcIClass& z);
MpfcIClass operator / (const cxsc::real& x, const MpfcIClass& z);
MpfcIClass operator / (const int x, const MpfcIClass& z);
```

```
MpfcIClass & operator /= (MpfcIClass& z, const MpfcIClass& x);
MpfcIClass & operator /= (MpfcIClass& z, const cxsc::cinterval& x);
MpfcIClass & operator /= (MpfcIClass& z, const MPFR::MpfcClass& x);
MpfcIClass & operator /= (MpfcIClass& z, const cxsc::complex& x);
MpfcIClass & operator /= (MpfcIClass& z, const mpfi_t& x);
MpfcIClass & operator /= (MpfcIClass& z, const MpfiClass& x);
MpfcIClass & operator /= (MpfcIClass& z, const cxsc::interval& x);
MpfcIClass & operator /= (MpfcIClass& z, const mpfr_t& x);
MpfcIClass & operator /= (MpfcIClass& z, const MPFR::MpfrClass& x);
MpfcIClass & operator /= (MpfcIClass& z, const double& x);
MpfcIClass & operator /= (MpfcIClass& z, const cxsc::real& x);
MpfcIClass & operator /= (MpfcIClass& z, const int x);
```



## 6.12. Mathematische Funktionen

### 6.12.1. Standard-Implementierung

Die Implementierung komplexwertiger Intervallfunktionen mit komplexen Intervallargumenten der Klasse `MpfcClass` erfolgt ganz analog zu den Funktionen der Klasse `MpfiClass`, vgl. dazu die Seite 64.

Unabhängig von der Präzision des Eingangsarguments werden die Real- und Imaginärteil-Intervalle des exakten Funktionswertes durch die entsprechenden Ergebnisintervalle in der voreingestellten Current-Präzision garantiert und nahezu optimal eingeschlossen.

Es gibt jedoch einige Funktionen, wie z.B. die komplexe Konjugation oder die Real- und Imaginärteil-Funktionen, bei denen man von dieser Standard-Implementierung abweichen sollte, d.h. die Präzision der Ergebnisintervalle sollte nicht mit der Current-Präzision, sondern mit der Präzision der Eingangsintervalle genau übereinstimmen. Im folgenden Abschnitt werden diese Funktionen kurz zusammengestellt.

### 6.12.2. Davon abweichende Funktionen

```
MpfiClass Re (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  wird das Realteil-Intervall  $\text{Re}(z) = x$  ohne Rundung in der Präzision von  $z$  zurückgegeben.

```
MpfiClass Im (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  wird das Imaginärteil-Intervall  $\text{Im}(z) = y$  ohne Rundung in der Präzision von  $z$  zurückgegeben.

```
MpfcClass Inf (const MpfcClass& z, PrecisionType prec);
```

Mit  $z = x + i \cdot y$  wird eine komplexe Zahl  $c$  in der Präzision `prec` zurückgegeben.  $\text{Re}(c)$  und  $\text{Im}(c)$  sind dabei die jeweils größten Maschinenzahlen, für die gilt:  $\text{Re}(c) \leq \text{Inf}(x)$  und  $\text{Im}(c) \leq \text{Inf}(y)$ . Wird `prec` nicht angegeben, so wird  $c$  ganz entsprechend in der Current-Präzision berechnet. Ist speziell `prec` die Präzision von  $z$ , so gilt in den beiden oberen Ungleichungen das Gleichheitszeichen.

```
MpfcClass Sup (const MpfcClass& z, PrecisionType prec);
```

Mit  $z = x + i \cdot y$  wird eine komplexe Zahl  $c$  in der Präzision `prec` geliefert.  $\text{Re}(c)$  und  $\text{Im}(c)$  sind dabei die jeweils kleinsten Maschinenzahlen, für die gilt:  $\text{Re}(c) \geq \text{Sup}(x)$  und  $\text{Im}(c) \geq \text{Sup}(y)$ . Wird `prec` nicht angegeben, so wird  $c$  ganz entsprechend in der Current-Präzision berechnet. Ist speziell `prec` die Präzision von  $z$ , so gilt in den beiden oberen Ungleichungen das Gleichheitszeichen.

```
MpfcClass mid(const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  wird eine komplexe Zahl  $c$  in der gleichen Präzision von  $z$  geliefert.  $\text{Re}(c)$  und  $\text{Im}(c)$  sind dabei die jeweiligen Mittelpunkte von  $x$  und  $y$ . Beachten Sie, dass bei möglichen späteren Rundungen von  $c$  diese Mittelpunktseigenschaft verloren gehen kann.

```
MpfcClass Blow(const MpfcClass& op1, const MPFR::MpfrClass& op2);
```

Der Rückgabewert ist ein mit `op2` aufgeblähtes Intervall `op1` mit gleicher Präzision. `Blow(...)` ist genauso definiert wie die gleichnamige Funktion in C-XSC.

```
void times2pown (MpfcClass& z, long int k)
```

Obige Funktion liefert mit dem Eingabewert `z` eine optimale Einschließung von  $z \cdot 2^k$  mit gleicher Präzision zurück, d.h. die mit  $2^k$  multiplizierten Real- und Imaginärteil-Intervalle werden optimal eingeschlossen. Solange kein Über- oder Unterlauf entsteht, wird  $z \cdot 2^k$  sogar **exakt**, d.h. rundungsfehlerfrei eingeschlossen.

```
void set_nan (MpfcClass& z);
```

Setzt `z` auf `([NaN,NaN], [NaN,NaN])`, wobei die Präzision von `z` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_inf (MpfcClass& z);
```

Setzt `z` auf `([-Inf,+Inf],[-Inf,+Inf])`, wobei die Präzision von `z` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
void set_zero (MpfcClass& z);
```

Setzt `z` auf `([0,0],[0,0])`, wobei die Präzision von `z` erhalten bleibt und deshalb mit der Current-Precision nicht übereinstimmen muss.

```
MpfcClass conj (const MpfcClass& z);
```

Mit  $z = x + i \cdot y$  ist der Rückgabewert von `conj(z)` gegeben durch  $x - i \cdot y$ , wobei die Präzision von `x` und `y` nicht geändert wird und daher mit der Current-Precision auch nicht übereinstimmen muss.

### 6.12.3. Elementarfunktionen

Tabelle 6.1.: Funktionen mit  $z = x + i \cdot y$ ,  $a, b$  vom Typ `MpfiClass`;  $x, y$  vom Typ `MpfiClass`

Funktion	Aufruf	Funktion	Aufruf
$\text{conj}(z) = x - i \cdot y$	<code>conj(z)</code>	$2^z$	<code>exp2(z)</code>
$\text{Re}(z) = x$	<code>Re(z)</code>	$10^z$	<code>exp10(z)</code>
$\text{Im}(z) = y$	<code>Im(z)</code>	$e^z$	<code>exp(z)</code>
$ z $	<code>abs(z)</code>	$e^z - 1$	<code>expm1(z)</code>
$\text{Arg}(z)$	<code>Arg(z)</code>	$e^{z^2}$	<code>expx2(z)</code>
$\text{arg}(z)$	<code>arg(z)</code>	$e^{z^2} - 1$	<code>expx2m1(z)</code>
$\text{arg}(1 + z)$	<code>argp1(z)</code>	$e^{-z^2}$	<code>expmx2(z)</code>
$z^2$	<code>sqr(z)</code>	$e^{-z^2} - 1$	<code>expmx2m1(z)</code>
$z^2 + a \cdot z + b$	<code>poly2(z, a, b)</code>	$z^n, n \in \mathbb{Z}$	<code>power_fast(z, n)</code>
$1/z$	<code>reci(z)</code>	$z^n, n \in \mathbb{Z}$	<code>power(z, n)</code>
$1/z^2$	<code>reci_z2(z)</code>	$z^p, p: \text{MpfiClass}$	<code>pow(z, p)</code>
$1/(1 + z^2)$	<code>reci_1pz2(z)</code>	$z^p, p: \text{MpfiClass}$	<code>pow_all(z, p)</code>
$1/(1 - z^2)$	<code>reci_1mz2(z)</code>	$z^w, w: \text{MpfiClass}$	<code>pow(z, w)</code>
$\sqrt{z}$	<code>sqrt(z)</code>	$\sin(z)$	<code>sin(z)</code>
$\sqrt{z}$	<code>sqrt_all(z)</code>	$\cos(z)$	<code>cos(z)</code>
$1/\sqrt{z}$	<code>sqrt_r(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z+1} - 1$	<code>sqrtp1m1(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{1+z^2}$	<code>sqrt1px2(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$\sqrt{1-z^2}$	<code>sqrt1mx2(z)</code>	$\arccos(z)$	<code>acos(z)</code>
$\sqrt{z^2-1}$	<code>sqrtx2m1(z)</code>	$\arctan(z)$	<code>atan(z)</code>
$\sqrt{-z^2-1}$	<code>sqrtnx2m1(z)</code>	$\text{arccot}(z)$	<code>acot(z)</code>
$\sqrt[n]{z}$	<code>sqrt(z, n)</code>	$\sinh(z)$	<code>sinh(z)</code>
$\sqrt[n]{z}$	<code>sqrt_all(z, n)</code>	$\cosh(z)$	<code>cosh(z)</code>
$\log(z)$	<code>Ln(z)</code>	$\tanh(z)$	<code>tanh(z)</code>
$\log(1 + z)$	<code>Lnp1(z)</code>	$\text{coth}(z)$	<code>coth(z)</code>

Fortsetzung auf der nächsten Seite

Fortsetzung der vorherigen Seite			
Funktion	Aufruf	Funktion	Aufruf
$\log(z)$	<code>ln(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$\log(1+z)$	<code>lnp1(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$\log_2(z)$	<code>log2(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$\log_{10}(z)$	<code>log10(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>

### Anmerkungen:

- Bei den Funktionen `conj(z)`, `Re(z)` und `Im(z)` werden die Ergebnisse in der Präzision von  $z$  zurückgegeben, die also nicht mit der Current-Präzision übereinstimmen muss.
- `abs(z)` kann mit einem zusätzlichen Präzisionsparameter `prec` aufgerufen werden. Dann wird das exakte, reelle Ergebnisintervall durch ein Intervall der Präzision `prec` garantiert eingeschlossen. Ohne `prec` wird das exakte, reelle Ergebnisintervall durch ein Intervall mit der Current-Präzision eingeschlossen.
- `Arg(z)` kann mit einem zusätzlichen Präzisionsparameter `prec` aufgerufen werden. Dann wird das exakte, reelle Argumentintervall, d.h. das exakte, reelle Ergebnisintervall, durch ein Intervall der Präzision `prec` garantiert eingeschlossen. Ohne `prec` wird das exakte, reelle Argumentintervall durch ein Intervall mit der Current-Präzision eingeschlossen. Enthält das komplexe Eingangsintervall  $z$  eine negative, reelle Zahl -auch auf dem Intervallrand!-, so erfolgt eine entsprechende Fehlermeldung. Nicht erlaubte Intervalle sind demnach  $z = [-2, -1] + i \cdot [-1, 0]$  oder  $z = [-1, -1] + i \cdot [0, 1]$  oder  $z = [-1, +1] + i \cdot [0, 0]$ .  
 $z = [-2, -1] + i \cdot [-2, -1]$  liefert: `Arg(z) = [-2.67795, -2.03444]`.
- `arg(z)` kann mit einem zusätzlichen Präzisionsparameter `prec` aufgerufen werden. Dann wird das exakte, reelle Argumentintervall, d.h. das exakte, reelle Ergebnisintervall, durch ein Intervall der Präzision `prec` garantiert eingeschlossen. Ohne `prec` wird das exakte, reelle Argumentintervall durch ein Intervall mit der Current-Präzision eingeschlossen. `arg(z)` liefert für jedes  $z \in \mathbb{IC}$  eine Einschließung des exakten Argumentintervalls, dabei gilt: `arg(z)  $\subset$   $[-\pi, 3\pi/2]$` . Hier einige Beispiele:
  - `arg([-2, -1] + i · [-2, -1]) = [-2.67795, -2.03444]`,
  - `arg([-2, +1] + i · [-2, +2]) = [-3.14160, +3.14160]`,
  - `arg([-0, +0] + i · [-0, +0]) = [-0.00000, +0.00000]`,
  - `arg([-2, -1] + i · [+0, +1]) = [+2.35619, +3.14160]`,
  - `arg([-2, -1] + i · [-1, -0]) = [-3.14160, -2.35619]`,
  - `arg([-2, -1] + i · [-0, +0]) = [+3.14159, +3.14160]`,
  - `arg([-0, +0] + i · [-1, +1]) = [-1.57080, +1.57080]`,
  - `arg([-1, +0] + i · [-1, +1]) = [+1.57079, +4.71239]`.
- `argp1(z)` berechnet `arg(1+z)` und kann mit einem zusätzlichen Präzisionsparameter `prec` aufgerufen werden. Dann wird das exakte, reelle Argumentintervall, d.h. das exakte, reelle Ergebnisintervall, durch ein Intervall der Präzision `prec` garantiert eingeschlossen. Ohne `prec` wird das exakte, reelle Argumentintervall durch ein Intervall mit der Current-Präzision eingeschlossen. `argp1(z)` liefert für jedes  $z \in \mathbb{IC}$  eine Einschließung des exakten Argumentintervalls, dabei gilt: `argp1(z)  $\subset$   $[-\pi, 3\pi/2]$` . Hier einige Beispiele:
  - `argp1([-3, -2] + i · [-2, -1]) = [-2.67795, -2.03444]`,
  - `argp1([-3, +0] + i · [-2, +2]) = [-3.14160, +3.14160]`,

```

argp1([-1, -1] + i * [-0, +0]) = [-0.00000, +0.00000],
argp1([-3, -2] + i * [+0, +1]) = [+2.35619, +3.14160],
argp1([-3, -2] + i * [-1, -0]) = [-3.14160, -2.35619],
argp1([-3, -2] + i * [-0, +0]) = [+3.14159, +3.14160],
argp1([-1, -1] + i * [-1, +1]) = [-1.57080, +1.57080],
argp1([-2, -1] + i * [-1, +1]) = [+1.57079, +4.71239].

```

6. Beim Aufruf von `sqr(z)`, mit  $z = X + i \cdot Y$ , wird der Realteil von `sqr(z)` anstelle von  $X, Y$  mithilfe der reellen Intervalle  $|X|, |Y|$  berechnet, weitere Einzelheiten findet man auf Seite 287. Mit  $z = [-3, 2] + i \cdot [-5, 3]$  liefert daher  $z \cdot z = [-31, 24] + i \cdot [-20, 30]$  eine erhebliche Überschätzung des Realteils von `sqr(z) = [-25, 9] + i \cdot [-20, 30]`.
7. `power_fast(z, n)` liefert für zu breite Intervalle  $z$  und für betragsmäßig zu große  $n$ -Werte erhebliche Überschätzungen von  $z^n$ . Für sehr schmale Intervalle werden auch bei größeren  $n$ -Werten brauchbare Einschließungen in recht kurzen Laufzeiten berechnet. Bei breiteren Intervallen  $z$  sollte man daher die Funktion `power(z, n)` benutzen, [54].  
Mit  $Z = [1, 1] + i[1, 1]$  erhält man für  $\alpha := \{y \in \mathbb{C} \mid y = z^8, z \in Z\}$  die Einschließungen:

$$\begin{aligned} \text{power\_fast}(Z, 8) &= ([1.599999e1, 1.600001e1], [-3.918870e-15, 1.029199e-14]) \\ \text{power}(Z, 8) &= ([1.599999e1, 1.600001e1], [-3.918870e-15, 1.029199e-14]). \end{aligned}$$

Mit  $Z = [1, 1.125] + i[1, 1.25]$  erhält man für  $\beta := \{y \in \mathbb{C} \mid y = z^8, z \in Z\}$  die Einschließungen:

$$\begin{aligned} \text{power\_fast}(Z, 8) &= ([1.012943e1, 6.397266e1], [-2.897500e1, 4.951985e1]) \\ \text{power}(Z, 8) &= ([1.599999e1, 5.839539e1], [-1.193387e1, 3.337647e1]). \end{aligned}$$

Das letzte Beispiel zeigt deutlich die mit `power(Z, 8)` berechnete bessere Einschließung von  $\beta$ .

8. Mit `Ln(z)` werden nur Funktionswerte des Hauptzweiges des natürlichen Logarithmus eingeschlossen. Die Funktion wird auch als *analytische* Funktion bezeichnet, da ein Rechteckintervall  $z$  den Verzweigungsschnitt nur von oben berühren darf. Im Gegensatz dazu darf bei der *nicht-analytische* Funktion `ln(z)` das Rechteckintervall  $z$  den Verzweigungsschnitt ganz im Innern enthalten oder diesen von oben oder unten berühren, vgl. die Abbildungen C.17, C.18 auf den Seiten 294 und 295.
9. Mit `pow(Z, P)` wird eine Rechteck-Einschließung der komplexen Menge

$$\{y \in \mathbb{C} \mid y = e^{p \cdot \ln(z)}, z \in Z : \text{MpfiClass}, p \in P : \text{MpfiClass}\}$$

berechnet, wobei nur die Funktionswerte des Hauptzweiges eingeschlossen werden, d.h. das Rechteckintervall  $Z$  darf den Verzweigungsschnitt nur von oben berühren und nicht in seinem Innern enthalten.  $\ln(z)$  bedeutet den Hauptwert des komplexen Logarithmus.

10. Mit `pow_all(Z, P)` und  $z = |z| \cdot e^{i\varphi}, -\pi/2 < \varphi < +3\pi/2$  wird eine Einschließung der komplexen Menge

$$\{y \in \mathbb{C} \mid y = e^{p \cdot \ln(z)} = e^{p \cdot \ln|z|} \cdot e^{ip(\varphi + 2\pi k)}, z \in Z : \text{MpfiClass}, p \in P : \text{MpfiClass}, k \in \mathbb{Z}\}$$

berechnet, wobei aber jetzt  $Z$  den Verzweigungsschnitt auch in seinem Innern enthalten darf. Die Einschließung besteht entweder aus einem einzelnen Rechteckintervall oder aus vier Rechtecken, die einen Kreisring optimal einschließen. Weitere Einzelheiten finde man ab Seite 367.

11. Mit `pow(Z,W)` wird eine Rechteck-Einschließung der komplexen Menge

$$\{y \in \mathbb{C} \mid y = e^{w \cdot \ln(z)}, z \in \mathbf{Z}:\text{Mpfciclass}, w \in \mathbf{W}:\text{Mpfciclass}\}$$

berechnet, wobei unter  $\ln(z)$  der Hauptwert des komplexen Logarithmus zu verstehen ist. Das Rechteckintervall  $\mathbf{Z}$  darf den Verzweigungsschnitt nur von oben berühren und nicht in seinem Innern enthalten.

12. Bei der Auswertung von  $z^2 + a \cdot z + b$  mit `poly2(..)` sollte für  $a$  möglichst nur ein Punktintervall gewählt werden, um eine optimale Einschließung zu gewährleisten, vgl. S. 291.
13. Die mit **Magenta** gedruckten Funktionsaufrufe, wie z.B. `exp2` und `expmx2`, beziehen sich auf Implementierungen, die mithilfe der *Globalen Optimierung* realisiert wurden. Dadurch erhält man zwar optimale Einschließungen der Wertemengen, aber bei zu breiten Eingangsintervallen oder bei zu groß gewählten Präzisionen muss mit größeren Laufzeiten gerechnet werden, vgl. Abschnitt 7.4.1.

### 6.12.4. Optimale, komplexe Intervall-Einschließungen

Ab Seite 243 wird ausführlich beschrieben, wie ein Algorithmus zu realisieren ist, mit dem man für eine vorgegebene holomorphe Funktion  $f(z)$ ,  $z \in \mathbb{C}$ , zu einem komplexen Rechteck-Intervall  $Z$  den zugehörigen Wertebereich  $W_Z := \{w \in \mathbb{C} \mid w = f(z), z \in Z\}$  durch  $F(Z)$  **optimal** einschließt. Die bisher realisierten Funktionen mit optimaler Einschließung sind in der Tabelle auf Seite 115 zusammengestellt.

Am Beispiel der Funktion  $f(z) = \sqrt{-z^2 - 1}$  soll hier noch ausdrücklich betont werden, dass zu einem vorgegebenen Intervall  $Z$  die berechnete Einschließung  $F_1(Z)$  bei naiver intervallmäßiger Auswertung des Funktionsterms  $\sqrt{-z^2 - 1}$  durch  $\text{sqrt}(-\text{sqr}(z) \diamond 1)$  i.a. viel größer ausfällt, als bei Auswertung der Funktion  $\text{sqr}(\text{mx}2\text{m}1(Z))$  aus der Tabelle auf Seite 115, welche die optimale Einschließung  $W_Z \subseteq F(Z) \subset F_1(Z)$  des Wertebereichs  $W_Z$  liefert, vgl. Seite 320.

In der Praxis müssen aber oft auch Einschließungen komplexer Intervallausdrücke berechnet werden, die aus den vier Grundoperationen zusammen mit den Intervallfunktionen aus der Tabelle von Seite 115 aufgebaut sind. Im Gegensatz zu einem reellen Intervallausdruck reicht es jetzt jedoch nicht aus, den komplexen Ausdruck äquivalent so umzuformen, dass dieser die Variable  $Z$  nur einmal enthält, um danach eine optimale Einschließung des Wertebereichs berechnen zu können, vgl. Seite 74. Als Beispiel betrachten wir wieder die obige Funktion  $f(z) = \sqrt{-z^2 - 1}$ , in deren Funktionsterm die Variable  $z$  zwar nur einmal vorkommt, die aber bei naiver Intervallauswertung nach Seite 320 **keine optimale** Einschließung liefert.

Wie kompliziert schon bei der scheinbar einfachen Funktion  $f(z) := e^{z^2} = u(x, y) + i \cdot v(x, y)$  die Realteilmfunktion  $u(x, y)$  über dem Eingangsintervall  $Z = [2, 4] + i \cdot [0, 3]$  ausfällt, zeigt die folgende mit *Mathematica* erstellte Abbildung 6.1

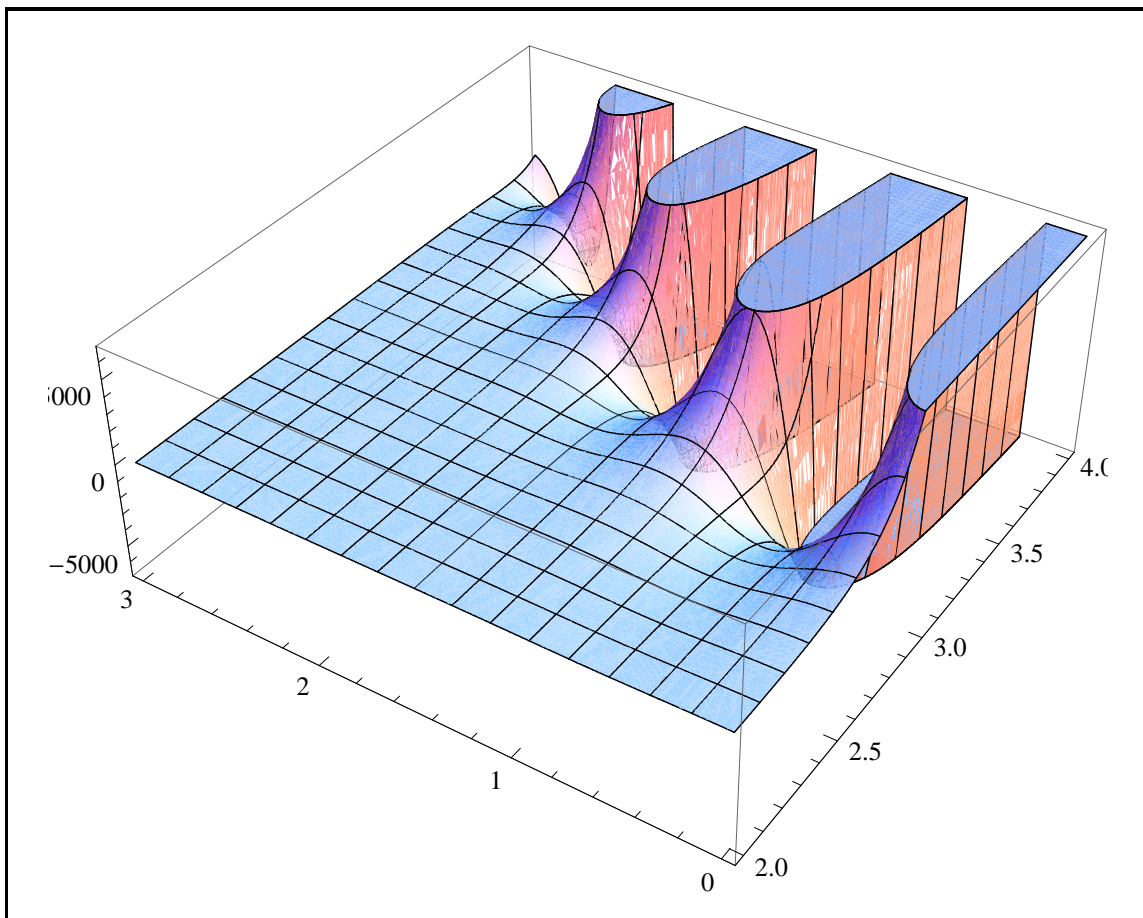


Abbildung 6.1.: Realteil-Funktion  $u(x, y)$  von  $f(z) := e^{z^2}$  über  $Z = [2, 4] + i \cdot [0, 3]$

Man erkennt mit wachsenden  $x, y$ -Werten aufeinanderfolgende Maxima und Minima, bedingt durch die  $2\pi$ -periodischen trigonometrischen Funktionen der komplexen Exponentialfunktion. Es wird daher sehr aufwendig oder sogar unmöglich sein, auf dem Rand eines gegebenen Eingangsintervalls  $Z$  die Lage des Maximums und Minimums, z.B. der Realteildfunktion  $u(x, y)$ , zu bestimmen. Bei komplizierten Intervallausdrücken wird es daher kaum möglich sein, eine optimale Einschließung des Wertebereichs  $W$  mit vertretbarem Aufwand zu berechnen. Wir fassen zusammen:

Um den Wertebereich  $W$  eines aus holomorphen Funktionen zusammengesetzten Intervallausdrucks möglichst optimal einzuschließen, sollte man auf möglichst viele Teilfunktionen zurückgreifen können, deren Wertebereiche selbst jeweils optimal eingeschlossen werden, vgl. dazu die Tabelle auf Seite 115.

Die Einschließung des Wertebereichs  $W$  eines komplexen Intervallausdrucks ist i.a. ein eigenständiges Problem, das nur in Ausnahmefällen optimal gelöst werden kann. Selbst wenn die Variable  $Z$  im Funktionsterm nur einmal vorkommt, wird bei naiver Intervallauswertung für  $W$  i.a. keine optimale Einschließung berechnet. Eine optimale Einschließung von  $W$  wird jedoch erreicht, wenn in den Real- und Imaginärteildfunktionen von  $f(z) = u(x, y) + i \cdot v(x, y)$  die Variablen  $x, y$  jeweils nur einmal vorkommen, vgl. z.B. die separable, komplexe Exponentialfunktion auf Seite 244.

Abschließend noch ein beruhigender Hinweis aus der numerischen Praxis:

Die erfolgreiche Einschließung einer komplexen Nullstelle  $z_0$ , z.B. nach dem Newton-Verfahren, liefert eine Folge von Einschließungsintervallen  $Z_k \ni z_0$ , deren Durchmesser für  $k \rightarrow \infty$  gegen Null streben, so dass die Intervalle  $Z_k$  selbst nicht wirklich optimal berechnet werden müssen, da die auftretenden Überschätzungen mit  $k \rightarrow \infty$  vernachlässigbar klein werden, vgl. Seite 155.



## 7. Anwendungen

### 7.1. Erste Nullstelle von $J_0(x)$

Das nachfolgende Programm liefert mit dem Intervall-Newton-Verfahren eine Einschließung der ersten Nullstelle  $x_1 = 2.4048\dots$  der Besselfunktion 1. Art  $J_0(x)$ , [20], [37]. In einer Umgebung von  $x_1$  benötigt man dazu Intervallfunktionen von  $J_0(x)$  und seiner ersten Ableitung  $J_0'(x) = -J_1(x)$ . Da beide Funktionen als Intervallfunktionen nicht zur Verfügung stehen, im Umgebungsintervall  $\alpha = [2, 3] \ni x_1$  jedoch monoton fallend bzw. wachsend sind, kann man  $J_0(x)$  und  $J_0'(x) = -J_1(x)$  in `MpfiClass f(const MpfiClass& x)` bzw. in `MpfiClass DERIV(const MpfiClass& x)` sehr einfach implementieren. Die Funktion `bool criter(const MpfiClass& x)` untersucht in `x` den Vorzeichenwechsel und die Monotonie von  $J_0(x)$ , wodurch die Eindeutigkeit der Nullstelle in `x` verifiziert ist. Zusätzlich wird untersucht, ob das Startintervall `x` die Bedingung `x ⊆ α` erfüllt.

```
1 // MPFR-13.cpp
2 // Newton-Verfahren zur Einschliessung der 1. Nst. x1 von J0(x);
3 // J0'(x) = -J1(x); x1 = 2.4048...;
4 // In [2,3] sind J0 und -J1(x) monoton fallend bzw. monoton wachsend,
5 // so dass deren Intervallfunktionen dort einfach zu implementieren sind.
6 // Als Start-Intervall sollte [2,3] gewaehlt werden.
7 #include "mpfiClass.hpp"
8
9 using namespace MPFR;
10 using namespace MPFI;
11 using namespace cxsc;
12 using namespace std;
13
14 MpfiClass f(const MpfiClass& x) // J0(x)
15 {
16     MpfrClass left ( J0(Sup(x),RoundDown) );
17     MpfrClass right ( J0(Inf(x),RoundUp) );
18     return MpfiClass(left, right);
19 }
20
21 MpfiClass DERIV(const MpfiClass& x) // -J1(x)
22 {
23     MpfrClass left ( J1(Sup(x),RoundDown) );
24     MpfrClass right ( J1(Inf(x),RoundUp) );
25     return -MpfiClass(left, right);
26 }
27
28 MpfiClass Start_Interval(void)
29 {
30     MpfrClass a, b;
31     cout << endl << "Left boundary point = ? "; cin >> a;
32     cout << "Right boundary point = ? "; cin >> b; cout << endl;
33     return MpfiClass(a,b);
34 }
35
36 bool criter(const MpfiClass& x) // Computing: J0(Inf(x))*f(Sup(x)) < 0
37 // and not 0 in J0'([x]) and x <= [2,3]?
38 {
39     MpfrClass a(2), b(3); // Start-interval <= [2,3] must be verified!
40     return ( Sup( f(MpfiClass(Inf(x))) * f(MpfiClass(Sup(x))) ) < 0 )
41     && !( 0 <= DERIV(x) ) && x <= MpfiClass(a,b);
```

```

42 }
43
44 int main(void)
45 {
46     PrecisionType prec = 6000;
47     MpfiClass::SetCurrPrecision(prec);
48     MpfiClass x, xOld, xMid, fxMid, fx, dfx;
49     x = Start_Interval();
50     cout << "Starting interval is: " << x << endl;
51
52     if (criter(x))
53     { // There is exactly one zero of f in the interval x
54         do
55         {
56             xOld = x;
57             cout << "Actual enclosure is " << x
58                 << ", Absolute diameter: " << diam(x) << endl;
59             xMid = MpfiClass(mid(x));
60             fxMid = f(xMid);
61             dfx = DERIV(x);
62             x = ( xMid - fxMid / dfx ) & x;
63         } while (x != xOld);
64         cout.precision(prec/3.321928095);
65         cout << "Final enclosure of the zero: " << x << endl
66             << "Absolute diameter = " << diam(x) << endl;
67         cout << "Correct decimal digits = " << common_decimals(x) << endl;
68         cout << "Enclosure of J0(x): " << f(x) << endl;
69         cout << "J0(Inf(x)) = " << f(MpfiClass(Inf(x))) << endl;
70     }
71     else
72         cout << "Criterion not satisfied!" << endl;
73
74     return 0;
75 }

```

Das Programm liefert nach Zeile 46 mit  $\text{prec} = 6000$  Bits  $\approx 1806$  Dezimalstellen die verkürzte Ausgabe

```

Left boundary point = ? 2.3
Right boundary point = ? 2.5

Starting interval is: [2.29999,2.50000]
Actual enclosure is [2.29999,2.50000], Absolute diameter: 2.00000e-1
Actual enclosure is [2.40464,2.40505], Absolute diameter: 3.99731e-4
Actual enclosure is [2.40482,2.40483], Absolute diameter: 3.20189e-9
Actual enclosure is [2.40482,2.40483], Absolute diameter: 1.03227e-19
Actual enclosure is [2.40482,2.40483], Absolute diameter: 5.36453e-41
Actual enclosure is [2.40482,2.40483], Absolute diameter: 7.24403e-84
Actual enclosure is [2.40482,2.40483], Absolute diameter: 6.60462e-170
Actual enclosure is [2.40482,2.40483], Absolute diameter: 2.74507e-342
Actual enclosure is [2.40482,2.40483], Absolute diameter: 2.37102e-687
Actual enclosure is [2.40482,2.40483], Absolute diameter: 8.84440e-1378
Actual enclosure is [2.40482,2.40483], Absolute diameter: 2.64293e-1806
Final enclosure of the zero:
[2.4048255576957727686216...1643831,2.4048255576957727686216...1643832]
Absolute diameter = 2.6429...4131431e-1806
Correct decimal digits = 1805

Enclosure of J0(x):
[-8.081224285204...20008648601e-1807, 5.6394916205927...970488619e-1807]

J0(Inf(x)) = [5.63949162059...488618e-1807,5.63949162059...488619e-1807]

```

**Anmerkungen:**

- Die Nullstelle  $x_1 = 2.40482\dots$  wird mit 1805 korrekten Dezimalstellen eingeschlossen. Anstelle von `prec = 6000` Bits kann in Zeile 46 auch eine sehr viel größere Präzision gewählt werden.
- Prinzipiell können ganz analog auch die Nullstellen  $x_3, x_5, x_7, \dots$  eingeschlossen werden, da  $J_0(x)$  und  $J'_0(x)$  in deren Umgebung das gleiche Monotonieverhalten aufweisen.

## 7.2. Einschließung reeller arithmetischer Ausdrücke

Das folgende Programm zeigt wesentliche Punkte, die zu beachten sind, wenn ein arithmetischer Ausdruck an einer speziellen Stelle  $x_0$  einzuschließen ist. Als Beispiel betrachten wir die Funktion

$$f(x) := \tan(x) - \sin(x) \equiv 2 \cdot \tan(x) \cdot \sin^2(x/2), \quad \text{mit: } x_0 = 2^{-16000}.$$

```
1 // MPFR-12.cpp
2 // Inclusion of f(x) := tan(x) - sin(x); x = 2^(-16000);
3 // f(x) = 2*tan(x)*sqr(sin(x/2));
4 #include "mpficlass.hpp"
5
6 using namespace MPFR;
7 using namespace MPFI;
8 using namespace cxsc;
9 using namespace std;
10
11 int main(void)
12 {
13     PrecisionType prec = 600, prec_old = prec;
14     MpfrClass::SetCurrPrecision(prec);
15     MpfiClass x(exp2(MpfiClass(-16000))), y; // x = 2^(-16000)
16     // Evaluating the difference expression
17     y = tan(x) - sin(x);
18     cout.precision(prec/3.321928095);
19     cout << "f(x) included by " << y << endl;
20     // Evaluation with fifty-fivefold Precision
21     MpfrClass::SetCurrPrecision(55*prec);
22     y = tan(x) - sin(x);
23     y.RoundPrecision(prec_old);
24     cout << "f(x) included by " << y << endl;
25     // Evaluating the simplified expression
26     MpfrClass::SetCurrPrecision(prec_old);
27     y = 2*tan(x)*sqr(sin(x/2)); // product expression
28     cout << "f(x) included by " << y << endl;
29
30     return 0;
31 }
```

Das Programm liefert die verkürzte Ausgabe

```
f(x) included by [0,2.394380862273962683...288266110563e-4997]
f(x) included by [1.81626048...4826128e-14450,1.81626048...4826129e-14450]
f(x) included by [1.81626048...4826128e-14450,1.81626048...4826129e-14450]
```

Anmerkungen:

- Die erste Einschließung ist wegen auftretender Auslöschung sehr grob, da die Differenz in der Nähe ihrer Nullstelle 0 auszuwerten ist.
- Eine optimale Einschließung dieser Differenz erhält man erst bei 55-facher Präzision.
- Die Einschließung des Produkts ist schon bei einfacher Präzision nahezu optimal.

Ist eine Summe oder Differenz in der Nähe einer Nullstelle auszuwerten, so muss zur Kompensation möglicher Auslöschungen eine mehrfache Präzision gewählt werden. Summen oder Differenzen sind daher möglichst in Produkte umzuformen!

## 7.3. Automatische Differentiation

Die Grundlagen der Automatischen Differentiation werden beschrieben in [20], so dass hier nur eine kurze Zusammenfassung angegeben wird.

Bei vielen wissenschaftlichen Anwendungen ist die Berechnung von Ableitungen komplizierter Funktionsausdrücke ein zentrales Problem. Beispiele sind die Einschließung der Nullstellen oder der relativen Extrema solcher Ausdrücke. Es gibt dabei drei verschiedene Methoden, um diese Ableitungen zu berechnen: Die *numerische* Differentiation, die *symbolische* Differentiation und die *automatische* Differentiation.

Die *numerische* Differentiation benutzt Differenzapproximationen zur Berechnung der Ableitungen, die *symbolische* Differentiation berechnet explizite Formeln für diese Ableitungen und die *automatische* Differentiation benutzt ebenfalls diese Formeln, die jedoch automatisch intervallmäßig ausgewertet werden, so dass garantierte Einschließungen dieser Ableitungen berechnet werden. Der Vorteil der *automatischen* Differentiation besteht darin, dass nur der aktuelle Funktionsausdruck selbst, nicht aber die Formeln für die oft sehr komplizierten Ableitungen dieser Ausdrücke explizit angegeben werden müssen.

In der in `mpficlass.cpp` implementierten Klasse `MPDerivType` sind alle notwendigen arithmetischen Operatoren und Funktionen definiert, um für einen entsprechenden Funktionsausdruck den Funktionswert selbst und die beiden ersten Ableitungen berechnen zu können. Das folgende Programm zeigt, wie man für  $f(x) := e^{\cos(x)} \cdot (1 + \sin(x))$  z.B. an der Stelle  $x_0 = \sqrt{2}$  die Einschließungen des Funktionswertes  $f(x_0)$  und der beiden Ableitungen  $f'(x_0)$ ,  $f''(x_0)$  mit der Präzision `prec = 7000`, d.h. also mit nahezu 2107 korrekten Dezimalstellen, sehr einfach berechnen kann.

```
1 // MPFR-14.cpp
2 // f(x) = exp(cos(x))*(1+sin(x));
3 // Berechnung der Einschliessungen von
4 // f(x0), f'(x0), f''(x0) fuer x0 = sqrt(2)
5 // mit einer Praezision von 7000 bits, d.h.
6 // mit nahezu 2107 Dezimalstellen
7
8 #include "mpficlass.hpp"
9
10 using namespace MPFR;
11 using namespace MPFI;
12 using namespace cxsc;
13 using namespace std;
14
15 MPDerivType f(const MPDerivType& x) // f(x)
16 {
17     return exp(cos(x))*(1+sin(x));
18 }
19
20 int main(void)
21 {
22     PrecisionType prec = 7000;
23     MpfrClass::SetCurrPrecision(prec);
24     MpfiClass x(2), fx, dfx, ddfx;
25     x = sqrt(x);
26     ddfEval(f, x, fx, dfx, ddfx);
27
28     cout.precision(x.GetPrecision()/3.321928095); // output format
29     cout << "x0 in " << x << endl;
30     cout << "fx in " << fx << endl;
31     cout << "dfx in " << dfx << endl;
32     cout << "ddfx in " << ddfx << endl;
33
34     return 0;
35 }
```

Das Programm liefert die hier auf nur wenige Dezimalstellen verkürzte Ausgabe

```

x0 ∈ [1.414213562373095048801688724... 88360, 1.414213562373095048801688724... 88361]
fx ∈ [2.323222110037337194992998377... 68783, 2.323222110037337194992998377... 68785]
dfx ∈ [-2.11253887115066778267776671... 41006, -2.11253887115066778267776671... 41003]
ddfx ∈ [0.389909375302151731543830546... 39325, 0.389909375302151731543830546... 39359]

```

### Anmerkungen:

1. Beachten Sie, dass  $x_0 = \sqrt{2}$  nicht darstellbar ist und daher oben durch das echte Intervall  $x = [1.4142... 88360, 1.4142... 88361]$  eingeschlossen werden muss.
2. Aber auch bei Punktintervallen  $x$  werden die berechneten Einschließungen von  $fx = f(x_0)$ ,  $dfx = f'(x_0)$ ,  $ddf_x = f''(x_0)$  wegen der internen Intervallauswertungen grundsätzlich leicht überschätzt sein. Wenn jedoch, wie z.B. im Fall  $x \approx 1.5\pi$ , stärkere Auslöschungen auftreten, so werden die berechneten Einschließungen entsprechend grob ausfallen. Die damit verbundene geringere Genauigkeit kann jedoch durch die Wahl einer höheren Präzision `prec` stets korrigiert werden. Mit dieser höheren Präzision darf natürlich der Werte des Eingangsintervalls  $x \approx 1.5\pi$  *nicht* geändert werden.

Mit den folgenden Funktionen können zu einer vordefinierten Funktion  $f$  vom Typ `MPDerivType` und zu einem gegebenen Eingangsintervall  $x$  vom Typ `MpfiClass` garantierte Einschließungen von  $f(x)$ ,  $f'(x)$  oder  $f''(x)$  berechnet werden:

```
void fEval ( MPddf_FctPtr f, MpfiClass x, MpfiClass& fx );
```

```
void dfEval ( MPddf_FctPtr f, MpfiClass x, MpfiClass& fx, MpfiClass& dfx );
```

```
void ddfEval( MPddf_FctPtr f, MpfiClass x, MpfiClass& fx, MpfiClass& dfx,
               MpfiClass& ddfx );
```

`MPddf_FctPtr f` ist dabei ein vordefinierter Zeiger auf eine Funktion  $f$  vom Typ `MPDerivType`.

Hat man mit den Deklarationen

```

MPDerivType x, fx;
MpfiClass y, fy, dfy, ddfy;

```

und mit den Anweisungen `y = 123.0; x = DerivVar(y); fx = f(x);` eine Funktion  $f$  vom Typ `MPDerivType` z.B. an der Stelle  $y = 123.0$  ausgewertet, so kann man mit den Funktionen

```
const MpfiClass fValue ( const MPDerivType& fx );
```

```
const MpfiClass dfValue ( const MPDerivType& fx );
```

```
const MpfiClass ddfValue( const MPDerivType& fx );
```

und den Anweisungen `fy = fValue(fx); dfy = dfValue(fx); ddfy = ddfValue(fx);` Einschließungen von  $f(x)$ ,  $f'(x)$  und  $f''(x)$  berechnen.

Die nachfolgende Funktion liefert mit einem Objekt vom Typ `MpfiClass` eine Variable vom Typ `MPDerivType`. Ein Objekt vom Typ `MPDerivType` ist ein Tripel von `MpfiClass`-Intervallen, welche den Funktionswert und die beiden ersten Ableitungen einschließen. Bei einer Variablen  $u$  vom Typ `MPDerivType` lautet daher das Tripel:  $u = (y, [1,1], [0,0])$ .

```
MPDerivType DerivVar( const MpfiClass& y );
```

Für  $y$  sind folgende Typen möglich: `MpfiClass`, `MpfrClass`, `interval`, `real`, `double`.

Mit der folgenden Funktion wird aus einem Objekt vom Typ `MpfiClass` eine Konstante vom Typ `MPDerivType` generiert. Bei einer Konstanten `c` vom Typ `MPDerivType` lautet daher mit der Anweisung `c = DerivConst(y)` das Tripel: `c = (y, [0,0], [0,0])`.

```
MPDerivType DerivConst( const MpfiClass& y);
```

Für `y` sind folgende Typen möglich: `MpfiClass`, `MpfrClass`, `interval`, `real`, `double`. Die beiden letzten Funktionen liefern `MPDerivType`-Objekte, d.h. die oben beschriebenen Tripel von `MpfiClass`-Objekten, jeweils in der vordefinierten Current-Precision.

Wertzuweisungen an ein Objekt vom Typ `MPDerivType` erfolgen mit Hilfe der Operatoren

```
MPDerivType& operator = ( const MPDerivType& u);
```

```
MPDerivType& operator = ( const MpfiClass& u);
```

```
MPDerivType& operator = ( const MpfrClass& u);
```

```
MPDerivType& operator = ( const interval& u);
```

```
MPDerivType& operator = ( const double& u);
```

```
MPDerivType& operator = ( const real& u);
```

```
MPDerivType& operator = ( int u);
```

Der linke Operand vom Typ `MPDerivType` wird dabei stets in der Current-Precision generiert. Beachten Sie, dass bis auf den 1. Operator stets **Konstanten** vom Typ `MPDerivType` erzeugt werden, d.h. mit `u = 1.25` erhält man das Tripel `( [1.25,1.25], [0,0], [0,0] )` von `MpfiClass`-Intervallen in der voreingestellten Current-Precision.

Mit den Funktionen

```
static void SetCurrPrecision (PrecisionType prec);
```

```
static const PrecisionType GetCurrPrecision ();
```

kann die Current-Precision mit `prec` gesetzt oder abgerufen werden.

Die Memberfunktion

```
void SetPrecision (PrecisionType prec);
```

setzt die Präzision des aktuellen Objekts vom Typ `MPDerivType` auf `prec` und löscht dabei das aktuelle Tripel der `MpfiClass`-Intervalle.

Die Memberfunktion

```
void RoundPrecision(PrecisionType prec);
```

setzt die Präzision des aktuellen Objekts vom Typ `MPDerivType` auf `prec` und rundet dabei das aktuelle Tripel der `MpfiClass`-Intervalle in die neue Präzision.

Mit den beiden folgenden Funktionen wird die Basis nur für die Ein- und Ausgabe gesetzt oder abgerufen.

```
static void MpfrClass::SetBase (int b);
```

```
static const int MpfrClass::GetBase ();
```

Der voreingestellte Standardwert ist `b = 10`.

In der nachfolgenden Tabelle sind alle Funktionen vom Typ `MPDerivType` aufgelistet, für die Einschließungen ihres Funktionswertes und ihrer beiden ersten Ableitungen berechnet werden können. Wie im Programm `MPFR-14.cpp` können diese Funktionen beliebig verschachtelt und mit den vier Grundoperationen verknüpft werden. Damit kann man fast beliebig komplizierte Funktionen generieren, um deren Funktionswerte und die beiden ersten Ableitungen in fast beliebiger Genauigkeit einzuschließen.

Tabelle 7.1.: Funktionen vom Typ `MPDerivType` zur Automatischen Differentiation

Funktion	Aufruf	Funktion	Aufruf
$x^2$	<code>sqr(x)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$x^n, n \in \mathbb{Z}$	<code>power(x,n)</code>	$\sin(x)$	<code>sin(x)</code>
$x^y$	<code>pow(x, y)</code>	$\cos(x)$	<code>cos(x)</code>
$1/x$	<code>reci(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt{x}$	<code>sqrt(x)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt[3]{x}, x \in \mathbb{R}$	<code>cbrt(x)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqrt(x,n)</code>	$\arccos(x)$	<code>acos(x)</code>
$1/\sqrt{x}$	<code>sqrt_r(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$\sqrt{x^2-1}$	<code>sqrtx2m1(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\sqrt{x+1}-1$	<code>sqrtp1m1(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$e^x$	<code>exp(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$2^x$	<code>exp2(x)</code>	$1/\sin(x)$	<code>csc(x)</code>
$10^x$	<code>exp10(x)</code>	$1/\cos(x)$	<code>sec(x)</code>
<i>Fortsetzung auf der nächsten Seite</i>			



Fortsetzung der vorherigen Seite			
Funktion	Aufruf	Funktion	Aufruf
$\ln(x)$	<code>ln(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$\ln(1+x)$	<code>lnp1(x)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$\log_2(x)$	<code>log2(x)</code>	$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>
$\ln(\sin(x))$	<code>ln_sin(x)</code>		

Wenn man bei der praktischen Anwendung der Funktionen aus obiger Tabelle neue Ausdrücke generiert, so sollte man unbedingt Teilausdrücke, die in obiger Tabelle bereits zu finden sind, auch **direkt** benutzen, um starke Überschätzungen bei den Einschließungen von Funktionswert und den beiden ersten Ableitungen zu vermeiden. Als Beispiel betrachten wir analog zum Programm MPFR-14.cpp auf Seite 125 den Ausdruck `g`

```
MPDerivType g(const MPDerivType& x) // g(x)
{
    return sqrt(1-x*x);
}
```

der unbedingt zu ersetzen ist durch den äquivalenten Ausdruck `f`

```
MPDerivType f(const MPDerivType& x) // f(x)
{
    return sqrt1mx2(x);
}
```

Für `g(x)` erhält man mit den Anweisungen

```
int main(void)
{
    PrecisionType prec = 100;
    MpfrClass::SetCurrPrecision(prec);
    MpfrClass r(1);
    r = pred(r);
    MpfiClass x(r), gx, dgx, ddgx;
    ddfEval(g, x, gx, dgx, ddgx);

    cout.precision(x.GetPrecision()/3.321928095); // output format
    cout << "x0 in " << x << endl;
    cout << "gx in " << gx << endl;
    cout << "dgx in " << dgx << endl;
    cout << "ddgx in " << ddgx << endl;

    return 0;
}
```

für das Argument `x0 = pred(1) = 1 - 2-100` die folgenden sehr groben Einschließungen:



## 7.4. Globale Optimierung

Bei der eindimensionalen globalen Optimierung werden für eine zweimal stetig differenzierbare Funktion  $f : \mathbb{R} \rightarrow \mathbb{R}$  über einem Intervall  $[x] \in \mathbb{R}$  alle Lösungen  $x^*$  berechnet von

$$(7.1) \quad \min_{x \in [x]} f(x).$$

Dabei werden die Minimumstellen  $x^*$  und der minimale Funktionswert  $f^* := f(x^*)$  garantiert eingeschlossen. Der verwendete Algorithmus basiert auf Hansen, [27], [28], [29], Ratscheck, [31] und Ratz, [59], weitere Einzelheiten findet man in [20]. Die Folgenden Abbildungen zeigen die

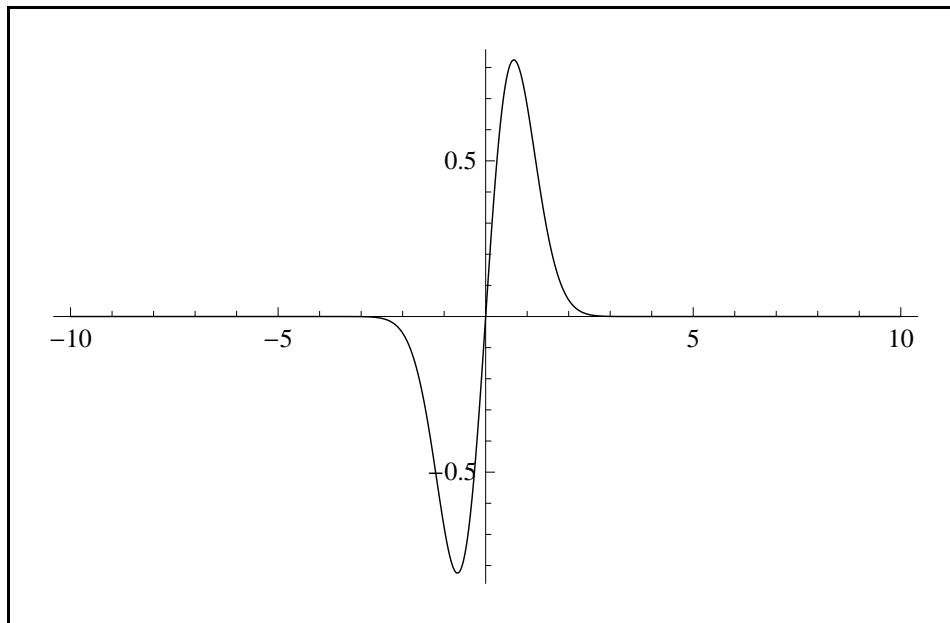


Abbildung 7.1.:  $f(x) := (x + \sin(x)) \cdot e^{-x^2}$ ,  $x \in [-10, 10]$

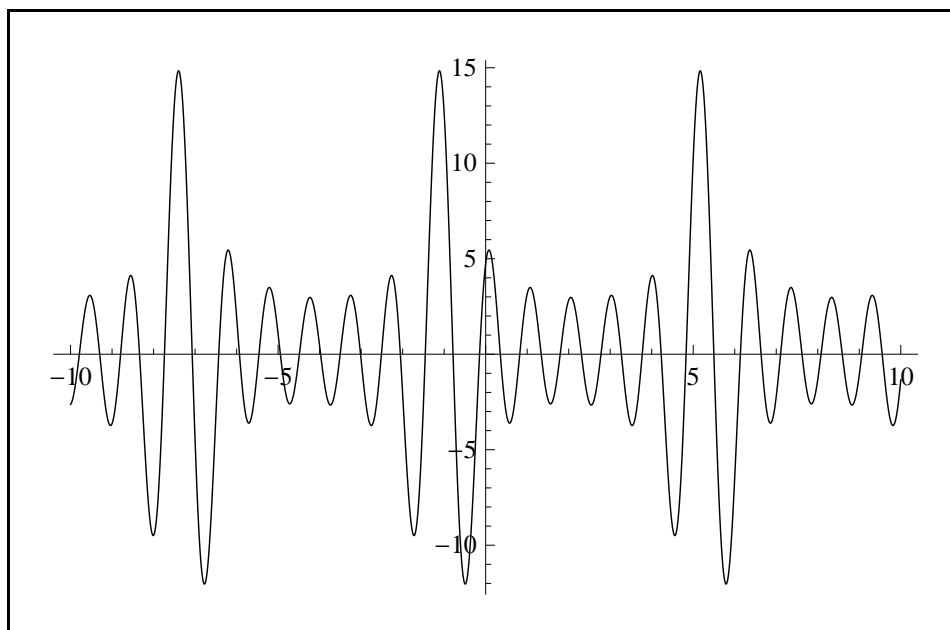


Abbildung 7.2.:  $g(x) := -\sum_{k=1}^5 k \cdot \sin((k+1) \cdot x + k)$ ,  $x \in [-10, 10]$

## Graphen der Funktionen

$$f(x) = (x + \sin(x)) \cdot e^{-x^2},$$

$$g(x) = - \sum_{k=1}^5 k \cdot \sin((k+1) \cdot x + k),$$

für die im Intervall  $[-10, 10]$  jeweils der oder die minimalen Funktionswerte  $f^*$  einzuschließen sind. Die Einschließung von  $f^*$  und der lokalen Minimumstellen  $x^*$  wird realisiert durch das nachfolgende Programm MPR-16.cpp:

```

1 //-----
2 // Program MPFR-16.cpp: Global optimization (one-dimensional)
3 // This program uses module 'MPgop1' to compute the global optimizers of the
4 // functions
5 //
6 //    f(x) = (x + sin(x)) * exp(-sqr(x)).
7 //
8 // and
9 //
10 //    g(x) = - sum ( k * sin((k+1)*x + k) )
11 //              k=1
12 //
13 // A starting interval and a tolerance must be entered.
14 //-----
15
16 #include <iostream>
17 #include "mpficlass.hpp"
18
19 using namespace cxsc;
20 using namespace std;
21 using namespace MPR;
22 using namespace MPFI;
23
24 MPDerivType f ( const MPDerivType &x )
25 { return (x + sin(x)) * exp(-sqr(x)); }
26
27 MPDerivType g ( const MPDerivType &x )
28 {
29     MPDerivType s;
30     int k;
31
32     s = DerivConst(0.0);
33     for (k = 1; k <= 5; k++)
34         s = s + k * sin ( (k+1) * x + k );
35     return -s;
36 }
37
38 //-----
39 // Function for printing and reading information to call the function
40 // 'AllGOpl'. It must be called with the function parameter 'f' and a
41 // string 'Name' containing a textual description of that function.
42 //-----
43 void compute ( MPddf_FctPtr f, char *Name )
44 {
45     MpfiClass SearchInterval, Minimum;
46     MpfrClass Tolerance;
47     MPivector Opti;
48     intvector Unique;
49     int NumberOfOptis, i, Error;
50     PrecisionType prec;
51
52     cout << "Computing all global minimizers of the function " << Name << endl;
53     cout << "Search interval      : "; cin >> SearchInterval;

```

```

54 cout << "Tolerance (relative) : "; cin >> Tolerance;
55 cout << endl;
56
57 if (Tolerance>1e-8) Tolerance = 1e-8;
58 // With the entered 'Tolerance' the corresponding output format is calculated:
59 prec = -expo(Tolerance) + 1;
60 if (prec < 53) prec = 53;
61 cout.precision(prec/3.321928095); // output format;
62
63 AllGOp1(f, SearchInterval, Tolerance, Opti, Unique, NumberOfOptis, Minimum, Error);
64
65 for (i = 1; i <= NumberOfOptis; i++) {
66     cout << Opti[i] << endl;
67     if (Unique[i])
68         cout << "encloses a locally unique candidate for a global minimizer!";
69     else
70         cout << "may contain a local or global minimizer!";
71     cout << endl;
72 }
73
74 if (NumberOfOptis != 0)
75     cout << endl << Minimum << endl
76         << "encloses the global minimum value!" << endl;
77 cout << endl << NumberOfOptis << " interval enclosure(s)" << endl;
78
79 if (Error)
80     cout << endl << AllGOp1ErrMsg(Error) << endl;
81 else if ( (NumberOfOptis == 1) && Unique[1] )
82     cout << endl << "We have validated that there is "
83         "a unique global optimizer!" << endl;
84 }
85
86 int main(void)
87 {
88     compute(f, (char*)" (x + SIN(x))*EXP(-x^2)");
89     cout << endl << endl;
90     compute(g, (char*)" -SUM(k*SIN((k+1)*x+k),k,1,5)");
91     return 0;
92 }

```

Mit dem Eingabeintervall  $[-10,10]$  und mit Tolerance (relative) :  $1e-34$  liefert das obige Programm für die Funktion  $f(x)$  die Ausgabe:

```

Search interval      : [-10,10]
Tolerance (relative) : 1e-34

```

```

[-6.795786600198815397288381602301209e-1,-6.795786600198815397288381602301207e-1]
encloses a locally unique candidate for a global minimizer!

```

```

[-8.242393984760766542477600525163374e-1,-8.242393984760766542477600525163372e-1]
encloses the global minimum value!

```

```

1 interval enclosure(s)

```

```

We have validated that there is a unique global optimizer!

```

Mit dem Eingabeintervall  $[-10,10]$  und mit Tolerance (relative) :  $1e-34$  liefert das obige Programm für die Funktion  $g(x)$  die Ausgabe:

```

Search interval      : [-10,10]
Tolerance (relative) : 1e-34

```

```
[-6.774576143438901030986146658809774,-6.774576143438901030986146658809773]
encloses a locally unique candidate for a global minimizer!
[-4.913908362593145540608598922507676e-1,-4.913908362593145540608598922507675e-1]
encloses a locally unique candidate for a global minimizer!
[5.791794470920271922864426874308238,5.791794470920271922864426874308239]
encloses a locally unique candidate for a global minimizer!

[-1.203124944216713894806863173737987e1,-1.203124944216713894806863173737986e1]
encloses the global minimum value!
```

3 interval enclosure(s)

### Anmerkungen:

1. Die Funktionen  $f(x), g(x)$  werden am Anfang des Programms als Funktionen vom Typ `MPDerivType` definiert, wobei alle Funktionen aus der Tabelle auf Seite 128 zur Verfügung stehen. Dabei sollte z.B.  $1/\sqrt{x}$  durch `sqrt_r(x)` und nicht mit dem Divisionsoperator `"/"` durch `1/sqrt(x)` definiert werden.
2. Durch Eingabe von `Tolerance (relative) : 1e-34` wird die relative Genauigkeit der zu berechnenden Einschließungen auf  $10^{-34}$  festgelegt. Damit wird die entsprechende Präzision für die internen Berechnungen festgelegt und zusätzlich das Ausgabeformat, das etwa 34 Dezimalstellen zur Verfügung stellt.
3. Die Funktion  $f(x)$  besitzt nur eine lokale Minimumstelle bei  $x^* = -0.679578\dots$  mit dem Minimum  $f^* = -0.82423\dots$ , wobei durch das Programm nachgewiesen wird, dass in der berechneten Einschließung von  $x^*$  auch keine weitere lokale Minimumstelle liegen kann.
4. Shubert's Funktion  $g(x)$ , [64], besitzt im Intervall  $[-10, +10]$  drei lokale Minimumstellen  $x^*$ , an denen das Minimum  $f^*$  von  $g(x)$  über  $[-10, 10]$  angenommen wird, vgl. die Abb. 7.2 auf Seite 131. Das Programm liefert Einschließungen der drei  $x^*$  und zeigt, dass jede Einschließung nur eine lokale Minimumstelle enthält. Mit `[-1.203...987e1, -1.203...986e1]` wird das globale Minimum  $f^*$  garantiert und nahezu optimal eingeschlossen.
5. Die Funktion  $r(x) := 2x^2 - (3 \cdot e^{-(a(x-b))^2})/c$  nimmt für  $a = 1\,000\,000\,000$ ,  $b = 0.0675$  und  $c = 100$  bei  $x^* = 0.067499\dots$  ein sehr schmales Minimum an, das bei üblichen Algorithmen ohne Intervallarithmetik leicht übersehen wird. Zeigen Sie, dass unser Programm mit `[0,10]` und `Tolerance (relative) : 1e-40000` das Minimum  $f^* = -2.08874999999999994\dots \cdot 10^{-2}$  problemlos und mit nahezu optimaler Genauigkeit einschließt.
6. Bei der Auswertung der zu untersuchenden Funktion  $f(x)$  mit der Differentiationsarithmetik kann es zu Fehlermeldungen kommen, wenn ein Intervallargument nicht mehr im Definitionsbereich einer verwendeten Elementarfunktion liegt oder wenn durch ein Intervall, das die Null enthält, zu dividieren ist. In diesen Fällen sollte man das Eingangsintervall in mehrere geeignete Teilintervalle unterteilen, um diese Fehlermeldungen zu vermeiden.

### 7.4.1. Optimale Einschließung des Wertebereichs analytischer Funktionen

Um den Wertebereich einer analytischen Funktion  $f(z) = u(x, y) + i \cdot v(x, y)$ ,  $z = x + i \cdot y \in \mathbb{C}$ , über einem Rechteck  $Z$  aus dem Definitionsbereich von  $f(z)$  optimal einzuschließen, genügt es, die Extrema der harmonischen Funktionen  $u(x, y)$ ,  $v(x, y)$  jeweils auf dem Rand von  $Z$  zu bestimmen. Bei den vergleichsweise einfachen Funktionen  $f(z) = 1/\sqrt{z}$  und  $f(z) = 1/(1+z^2)$  haben wir gesehen, dass die Berechnung der gesuchten Extrema auf dem Rand von  $Z$  jeweils zu aufwendigen Algorithmen führt, die zwar auch bei großen Präzisionen schnell ausführbar sind aber bei komplizierteren Funktionen kaum noch realisierbar sind, vergleichen Sie dazu z.B. die Abschnitte B.6, B.7.

Um eine optimale Einschließung z.B. für  $f(z) = e^{z^2} = u(x, y) + i \cdot v(x, y)$  zu realisieren, kann man die gesuchten Extrema von  $u(x, y) = e^{x^2-y^2} \cdot \cos(2xy)$  und  $v(x, y) = e^{x^2-y^2} \cdot \sin(2xy)$  durch globale Optimierung berechnen, indem man die jeweiligen Extrema auf den vier Rändern von  $Z$  für  $u(x, y)$  und  $v(x, y)$  einzeln bestimmt. Der Vorteil dieses Verfahrens ist die vergleichsweise einfache Realisierung des Algorithmus, der Nachteil ist aber die verhältnismäßig große Laufzeit, die z.B. bei einer Präzision von `prec = 4000` schon in der Größenordnung mehrerer Sekunden liegt. In der folgenden Klasse `f_expx2` werden u.a. die Realteil- und die Imaginärteilmfunktionen  $u, v$  von  $f(z) = e^{z^2}$  definiert, und mit der Funktion `Mpfciclass expx2(const Mpfciclass& z)` wird die optimale Einschließung des Wertebereichs von  $f(z) = e^{z^2}$  für  $z \in Z$  berechnet.

```
class f_expx2
{
    typedef MPDerivType T;          // T wird hier nur zur Abkuerzung eingefuehrt
    enum minORmax {mini, maxi};    // im Namensraum dieser Klasse f_expx2

    template <typename T1, typename T2>
    static T u(const T1& x, const T2& y)
    { // Real part u(x,y) of f(Z)
        T res = exp((x-y)*(x+y))*cos(2*x*y);
        if (mM==maxi) res= -res; // Es soll -Maximum berechnet werden
        return res;
    }
    static T u_x(const T& x)
    { // Funktion in x bei festgehaltenem y
        MpfrClass y(p);
        return u(x,y);
    }
    static T u_y(const T& y)
    { // Funktion in y bei festgehaltenem x
        MpfrClass x(p);
        return u(x,y);
    }
    template <typename T1, typename T2>
    static T v(const T1& x, const T2& y)
    { // Imaginary part v(x,y) of f(z)
        T res = exp((x-y)*(x+y))*sin(2*x*y);
        if (mM==maxi) res = -res; // Es soll -Maximum berechnet werden
        return res;
    }
    static T v_x(const T& x)
    { // Funktion in x bei festgehaltenem y
        MpfrClass y(p);
        return v(x,y); }
}
```

```

static T v_y(const T& y)
{ // Funktion in y bei festgehaltenem x
  MpfrClass x(p);
  return v(x,y);
}

static MpfiClass OPT(T (*f)(const T&), const MpfiClass& S,
                    MpfrClass fix, minORmax m)
{
  PrecisionType prec;
  MpfrClass Tolerance;
  MpfiClass minimum;
  MPivector Opti;
  intvector Unique;
  int      NumberOfOptis, Error;
  mM = m;
  p = fix;

  // Berechnung der 'Tolerance' aus der CurrentPrecision:
  prec = MPFI::MpfcicClass::GetCurrPrecision();
  Tolerance = 0.5;
  Tolerance = comp(Tolerance, -prec+1);

  AllGOpl(f, S, Tolerance, Opti, Unique, NumberOfOptis, minimum, Error);
  if (Error)
  {
    cout << endl << AllGOplErrMsg(Error) << endl;
    exit(1);
  }
  return minimum; // von f bzw -f ueber Intervall S (= eine Rechteckseite)
}

public:
static MpfcicClass EinschliessungUeberRechteck(MpfiClass X, MpfiClass Y)
{
  // Realteil u(x,y):
  MpfiClass u = OPT(u_x, X, Inf(Y), mini); // Einschliessung Minimum
  u|= -OPT(u_x, X, Inf(Y), maxi);         // Einschliessung Maximum
  u|= -OPT(u_x, X, Sup(Y), maxi);
  u|=  OPT(u_x, X, Sup(Y), mini);

  u|=  OPT(u_y, Y, Inf(X), mini);
  u|= -OPT(u_y, Y, Inf(X), maxi);
  u|= -OPT(u_y, Y, Sup(X), maxi);
  u|=  OPT(u_y, Y, Sup(X), mini);
  // Imaginaerteil v(x,y):
  MpfiClass v= OPT(v_x, X, Inf(Y), mini);
  v|= -OPT(v_x, X, Inf(Y), maxi);
  v|= -OPT(v_x, X, Sup(Y), maxi);
  v|=  OPT(v_x, X, Sup(Y), mini);
  v|=  OPT(v_y, Y, Inf(X), mini);
}

```



```

    v|= -OPT(v_y, Y, Inf(X), maxi);
    v|= -OPT(v_y, Y, Sup(X), maxi);
    v|=  OPT(v_y, Y, Sup(X), mini);
    return Mpfciclass(u,v);
}

Mpfciclass operator()(const Mpfciclass& Z)
{
    return EinschliessungUeberRechteck(Re(Z), Im(Z));
}

private:
    static MpfrClass p; // Wird mit festgehaltenem x- bzw. y-Wert belegt
    static minORmax mM; // AllGOp1 soll Minimum bzw. -Maximum berechnen
                        // -Maximum ist gerade Minimum von -f
}; // class f_expx2

// Definition der static Variablen der Klasse f_expx2
MpfrClass f_expx2::p(real(12), RoundNearest, 53);
f_expx2::minORmax f_expx2::mM (f_expx2::mini);

Mpfciclass expx2(const Mpfciclass& z)
// Intervallhuelle von f(z) ueber komplexem Intervall z
{
    int kx, ky, k;
    f_expx2 f; // Funktionsobjekt f der Klasse f_expx2 erzeugen
    PrecisionType prec = z.GetPrecision(),
        prec_old = MPFR::MpfrClass::GetCurrPrecision();
    if (prec<prec_old) prec = prec_old;
    MPFI::MpfiClass::SetCurrPrecision(prec);

    Mpfciclass res;
    MpfiClass x(Re(z)), y(Im(z)); // Realteil-/Imaginaerteilintervall
    kx = expo(Inf(x));
    k = expo(Sup(x));
    if (k > kx) kx = k;
    ky = expo(Inf(y));
    k = expo(Sup(y));
    if (k > ky) ky = k;
    if (kx+ky > 10000)
    {
        std::cerr << "Mpfciclass expx2(const Mpfciclass& z);
            Values of the interval bounds too great!"
            << std::endl;
        exit(1);
    }
    res = f( Mpfciclass(x,y) );
    res.RoundPrecision(prec_old);
    MPFI::MpfiClass::SetCurrPrecision(prec_old);
    return res;
}

```

Um eine neue Funktion  $f(z)$  mit dem Namen `abcd` durch *Globale Optimierung* zu realisieren, sind folgende Schritte durchzuführen:

1. Kopieren des obigen kompletten Codes aus der Datei `mpfciclass.cpp` in ein neues Arbeitsverzeichnis
2. Ersetzen aller Worte `exp2` durch `abcd`
3. Definitionen von  $u(x, y)$  und  $v(x, y)$  in den beiden Zeilen der Klasse `f_abcd` realisieren, dazu vorher eventuell die *Mathematica*-Funktion `ComplexExpand[...]` benutzen.
4. In der Funktion `Mpfciclass abcd(const Mpfciclass& z)` die Zeilen `kx=expo(Inf(x));` bis `exit(1); }` entfernen, die zur Vermeidung eines Überlaufs bei der Funktion  $e^{z^2}$  notwendig waren, FERTIG!

Mit der obigen Funktion `Mpfciclass exp2(const Mpfciclass& Z)` wird für das komplexe Intervall  $Z = [0.125, 0.35] + i \cdot [1, 3.25]$  mit  $f(z) := e^{+z^2}$ ,  $z \in \mathbb{C}$ , eine optimale und garantierte Einschließung des Wertebereichs  $W_{f,Z} := \{f(z) \mid z \in Z\}$  berechnet. Zum Vergleich wird bei gleicher Präzision `prec = 74` der gleiche Wertebereich  $W_{f,Z}$  durch die Intervallauswertung des komplexwertigen Ausdrucks `exp(sqr(Z))` eingeschlossen. Man erhält die Ergebnisse:

$$W_{f,Z} \subseteq ([-1.4918459347e - 6, 3.6205612011e - 1], [1.9076213940e - 5, 1.8774573850e - 1]),$$

$$W_{f,Z} \subset ([-2.1216070916e - 2, 3.7943155613e - 1], [6.5006537073e - 6, 3.9160562669e - 1]).$$

Die Einschließung von  $W_{f,Z}$  durch den Intervallausdruck `exp(sqr(Z))` ist deutlich größer als die optimale Einschließung mit Hilfe der *Globalen Optimierung*.

Wir zeigen jetzt, dass bei viel schmalere Intervallen, z.B.  $Z^* = [0.125, 0.126] + i \cdot [1, 1.01]$ , die Einschließung von  $W_{f,Z^*}$  durch die Intervallauswertung des Ausdrucks `exp(sqr(Z*))` i.a. sehr viel besser gelingt, man erhält die Ergebnisse

$$W_{f,Z^*} \subseteq ([3.5452724316e - 1, 3.6205612011e - 1], [9.1495293392e - 2, 9.3195421280e - 2]),$$

$$W_{f,Z^*} \subset ([3.5443826798e - 1, 3.6214700764e - 1], [9.0608448859e - 2, 9.4107266993e - 2]).$$

Bei hinreichend schmalen Intervallen  $Z^*$  sollte man also insbesondere bei hohen Präzisionen die Auswertung von `exp(sqr(Z*))` aus Laufzeitgründen dem Verfahren der *Globalen Optimierung* vorziehen, da dann bei naiver Intervallauswertung die Intervallüberschätzungen praktisch keine Rolle mehr spielen!

## 7.5. Reelle, eindimensionale Taylor Arithmetik

Die eindimensionale Taylorarithmetik liefert Einschließungen aller Taylorkoeffizienten und Ableitungen einer beliebig oft differenzierbaren Funktion  $f : D_f \rightarrow \mathbb{R}$  einer Variablen an der Stelle  $x_0 \in D_f \subseteq \mathbb{R}$ . Dabei kann  $x_0$  durch ein Intervall eingeschlossen werden, wenn  $x_0$  im verwendeten Zahlenraster nicht darstellbar ist. Die in **C-XSC** implementierte Taylorarithmetik kombiniert die symbolische Differentiation mit der Intervallrechnung und liefert dadurch sehr effektiv auch für Taylorkoeffizienten sehr hoher Ordnungen enge Einschließungen, wenn die interne Präzision hinreichend groß gewählt wird.  $f(x)$  besitzt für den Entwicklungspunkt  $x_0$  die Taylorreihe

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} \cdot (x - x_0)^k, \quad (f)_k := \frac{f^{(k)}(x_0)}{k!},$$

wobei mit  $(f)_k$  der Taylorkoeffizient der Ordnung  $k$  bezeichnet wird. Die Taylorarithmetik wurde entwickelt in den 60ziger Jahren von Moore, [50], Rall, [57], Griewank, [25], Corliss, [19] und vielen anderen. Eine Zusammenfassung der wichtigsten Rekursionsformeln zur Berechnung der  $(f)_k$  und weitere Anwendungen findet man in [11]. In diesem Abschnitt werden die theoretischen Grundlagen nicht weiter behandelt. Beschrieben werden nur die zur Verfügung stehenden Werkzeuge, mit denen in der Programmierumgebung **C-XSC** die Taylorkoeffizienten und Ableitungen in fast beliebiger Präzision berechnet werden können.

Ein Objekt der Klasse `MPitaylor` besitzt einen Vektor `tay1` vom Typ `MPivector` mit den Indices von 0 bis  $p$ , wobei  $0 \leq p \leq 500$  die gewählte Ordnung der zu berechnenden Taylorkoeffizienten ist. Die Taylorkoeffizienten  $(f)_k$  werden dabei durch die Komponenten vom Typ `MpfiClass` dieses Vektors eingeschlossen. Der maximale Wert  $p = 500$  kann bei Bedarf auch größer gewählt werden. Mit  $p = 0$  erhält man lediglich eine Einschließung für  $f(z_0) = (f)_0$ . Die Präzision kann für alle Intervall-Taylorkoeffizienten nach Bedarf einheitlich festgelegt werden. Der Vektor `tay1` einer unabhängigen Variablen `x` vom Typ `MPitaylor` hat die Form

$$\text{tay1} = (\mathbf{x}, [1], [0], [0], [0], [0], [0]), \quad \text{hier mit } p = 6,$$

und der entsprechende Vektor einer Konstanten `c` vom Typ `MpfiClass` ist gegeben durch

$$\text{tay1} = (\mathbf{c}, [0], [0], [0], [0], [0]), \quad \text{hier mit } p = 5.$$

Werden die Vektoren für Summe, Differenz, Produkt und Quotient dieser Objekte berechnet, so müssen die Ordnungen  $p$  aller dieser Objekte übereinstimmen. Für alle Elementarfunktionen der nachfolgenden Tabelle gibt es entsprechende Formeln, [11], mit denen alle Taylorkoeffizienten  $(f)_k$ ,  $k = 0, 1, \dots, p$ , z.B für  $f(x) = \sin(x)$ , berechnet werden können. Mithilfe dieser Funktionen und den vier Grundoperationen lassen sich dann beliebige Ausdrücke  $T(x)$  definieren, für die an Entwicklungspunkten  $x_0$  alle Taylorkoeffizienten oder Ableitungen bis zur Ordnung  $p$  berechnet werden können, falls  $T(x)$  an der Stelle  $x_0$  mindestens  $p$ -mal differenzierbar ist. Wählt man für  $x_0$  ein einschließendes Intervall  $z \ni x_0$ , z.B. weil  $x_0$  im Zahlenraster nicht darstellbar ist, so ist die berechnete Einschließung für  $(T)_k$  auch die Einschließung aller Taylorkoeffizienten der Ordnung  $k$  für jedes andere  $\hat{x}_0 \in z$ , mit  $\hat{x}_0 \neq x_0$ . Um möglichst gute Einschließungen der Taylorkoeffizienten zu berechnen, sollte  $z$  daher möglichst eng gewählt werden. Bei hoher Ordnung  $p$  muss mit einer hinreichend großen Präzision `prec` gerechnet werden, um die bekannten Intervallüberschätzungen bei der internen Intervallauswertung klein zu halten.

## 7.5.1. Elementarfunktionen

Tabelle 7.2.: Funktionen vom Typ MPitaylor zur Taylor-Arithmetik

Funktion	Aufruf	Funktion	Aufruf
$x^2$	<code>sqr(x)</code>	$\cos(x)$	<code>cos(x)</code>
$\sqrt{x}$	<code>sqrt(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt[n]{x}$ , $n = 2, 3, \dots$	<code>sqrt(x, n)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x)</code>	$\arccos(x)$	<code>acos(x)</code>
$\sqrt{x^2-1}$	<code>sqrtox2m1(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{x+1}-1$	<code>sqrtp1m1(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$x^y$	<code>pow(x, y)</code>	$\sinh(x)$	<code>sinh(x)</code>
$x^n$ , $n \in \mathbb{Z}$	<code>power(x, n)</code>	$\cosh(x)$	<code>cosh(x)</code>
$e^x$	<code>exp(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$2^x$	<code>exp2(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$\log_2(x)$	<code>log2(x)</code>	$1/\sin(x)$	<code>csc(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$1/\cos(x)$	<code>sec(x)</code>
$\ln(1+x)$	<code>lnp1(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$\ln(\sin(x))$	<code>ln_sin(x)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$\ln(\cos(x))$	<code>ln_cos(x)</code>	$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>
$\sin(x)$	<code>sin(x)</code>	$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>

## 7.5.2. Konstruktoren / Destruktoren

### 7.5.2.1. Konstruktoren

```
MPitaylor ();
```

Der Default-Konstruktor legt ein neues MPitaylor-Objekt in der Current-Precision an. Aufruf: `MPitaylor f;`  
Es wird kein Vektor `tay1` und damit auch keine Ordnung `p` dieses Vektors initialisiert.

```
MPitaylor(const MPitaylor& g, PrecisionType prec);
```

Mit dem Aufruf: `MPitaylor f(g, prec);` legt der Copy-Konstruktor mithilfe von `g` ein neues MPitaylor-Objekt `f` an, wobei die Komponenten des Vektors von `f` die Komponenten des Vektors von `g` in der neuen Präzision `prec` einschließen. Wird `prec` nicht gesetzt, so werden die Komponenten des Vektors von `g` in der voreingestellten Current-Precision eingeschlossen.

Alle folgenden Konstruktoren erzeugen ein MPitaylor-Objekt `f`. Der Vektor `tay1` enthält  $p + 1$  Komponenten vom Typ `MpfiClass` mit der Präzision `prec`.

Aufruf z.B.: `MPitaylor f(5, MpfiClass(1)/10, 400);` dies erzeugt den Vektor `([9.99999e-2, 1.00001e-1], [1.00000, 1.00000], [0,0], [0,0], [0,0], [0,0])`, wobei hier die Dezimalstellen aus Platzgründen verkürzt angegeben sind. Die 1. Komponente ist eine Einschließung des nicht-darstellbaren Entwicklungspunktes  $x = 0.1$  mit ca.  $400/3.321928095 \approx 120$  Dezimalstellen.

Wird beim Konstruktoraufruf `prec` nicht angegeben, so erhalten alle Komponenten von `tay1` die voreingestellte Current-Precision. Falls diese mit `SetCurrPrecision` nicht gesetzt wurde, erhalten die Komponenten die Präzision `prec = 53`.

```
MPitaylor (int p, const MpfiClass& x, PrecisionType prec);  
MPitaylor (int p, const interval& x, PrecisionType prec);  
MPitaylor (int p, const MpfrClass& x, PrecisionType prec);  
MPitaylor (int p, const real& x, PrecisionType prec);  
MPitaylor (int p, const double& x, PrecisionType prec);  
MPitaylor (int p, int x, PrecisionType prec);
```

### 7.5.2.2. Destruktor

```
~MPitaylor(){};
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

### 7.5.3. Zuweisungs-Operatoren

Unabhängig von der Präzision des rechten Operanden `op` erhalten bei allen folgenden Zuweisungsoperatoren alle Vektorkomponenten des linken Operanden stets die Current-Precision und schließen die `op`-Werte optimal ein.

```
MPitaylor operator = (const MPitaylor& op);  
MPitaylor operator = (const MpfiClass& op);  
MPitaylor operator = (const interval& op);  
MPitaylor operator = (const MpfrClass& op);  
MPitaylor operator = (const real& op);  
MPitaylor operator = (const double& op);  
MPitaylor operator = (int op);
```

Nur beim ersten Operator wird die Ordnung des linken Operanden gleich der Ordnung von `op`, sonst bleibt die Ordnung des linken Operanden erhalten.

## 7.5.4. Arithmetische Operatoren

Für alle arithmetischen Operationen mit Operanden vom Typ `MPitaylor` gilt:

Die exakten Ergebnisse für die Komponenten des Vektors `tayl` werden unabhängig von der Präzision der `tayl`-Komponenten der beiden rechten Operanden in der vor-eingestellten Current-Präzision optimal eingeschlossen.

Sind die rechten Operanden beide vom Typ `MPitaylor`, so müssen deren Ordnungen übereinstimmen, und der Ergebnisvektor `tayl` übernimmt diese Ordnung. Sind die Ordnungen der beiden rechten Operanden verschieden, so erfolgt eine entsprechende Fehlermeldung.

Ist einer der beiden rechten Operanden nicht vom Typ `MPitaylor`, so sind für diesen Operanden nur die folgenden Typen zugelassen:

`MpfiClass`, `interval`, `MpfrClass`, `real`, `double`, `int`.

Alle Operanden dieser Typen werden dabei stets als Konstanten interpretiert!

Für die Addition existieren beispielsweise die Operatoren:

```
MPitaylor operator + (const MPitaylor& op1, const MPitaylor& op2);  
MPitaylor operator + (const MPitaylor& op1, int op2);  
MPitaylor operator + (double& op1, const MPitaylor& op2);
```

## 7.5.5. Unabhängige Variablen und Konstanten

Zur Definition von unabhängigen Variablen und Konstanten siehe Seite 139. Unabhängige Variablen werden generiert durch die Funktion

```
MPitaylor var_MPitaylor(int ord, const MpfiClass& x, PrecisionType prec);
```

Für den 2. Parameter stehen folgende Datentypen zur Verfügung:

`MpfiClass`, `MpfrClass`, `interval`, `real`, `double`, `int`.

Konstanten werden generiert durch die Funktion

```
MPitaylor const_MPitaylor(int ord, const MpfiClass& x, PrecisionType prec);
```

Für den 2. Parameter stehen folgende Datentypen zur Verfügung:

`MpfiClass`, `MpfrClass`, `interval`, `real`, `double`, `int`.

Für beide obigen Funktionen gilt:

Die Vektorkomponenten des Rückgabewertes erhalten die Präzision `prec`. Wird dieser 3. Parameter nicht gesetzt, so erhalten die Vektorkomponenten die Current-Präzision. Die Ordnung des Rückgabevektors `tayl` wird durch den 1. Parameter `ord` festgelegt.

## 7.5.6. Precision Handling

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision in `MpfrClass` auf `prec`. Diese Current-Precision wird dann automatisch auch in `class MPitaylor` benutzt. Wird die Current-Precision nicht gesetzt, so wird in beiden Klassen mit der Default-Precision von 53 Bits gerechnet. Das Setzen der Current-Precision hat auf die Präzision der bis dahin benutzten Objekte **keinerlei** Einfluss.

**PrecisionType GetPrecision () const;**

Diese Memberfunktion gibt die Präzision des aktuellen Objekts in Bits zurück. Als Beispiel entsprechen dabei 302 Bits  $302/\log_2(10) \approx 91$  Dezimalstellen.

**void SetPrecision (PrecisionType prec);**

Diese Memberfunktion setzt die Präzision der aktuellen Komponenten von `tay1` auf `prec`. Die Komponentenwerte bleiben dabei nicht erhalten.

**void RoundPrecision (PrecisionType prec);**

Diese Memberfunktion schließt die Komponenten des aktuellen Vektors `tay1` mit der neuen Präzision `prec` ein. Sollte die Präzision der ursprünglichen Komponenten größer sein als `prec`, so erhält man eine gröbere Einschließung. Ist die Präzision jedoch kleiner als `prec`, so werden die restlichen binären Stellen mit Nullen aufgefüllt, so dass die ursprünglichen Werte der Taylor-Komponenten erhalten bleiben.

### 7.5.7. Zugriff auf die Vektorkomponenten

**MPivector get\_all\_coef(const MPitaylor& x);**

Zurückgegeben wird ein Vektor vom Typ `MPivector` der Dimension  $p + 1$ , welcher alle  $p + 1$  Taylor-Koeffizienten des Objekts `x` enthält.

**MpfiClass get\_j\_coef(const MPitaylor& x, int j);**

Zurückgegeben wird ein Intervall vom Typ `MpfiClass`, welches den  $j$ -ten Taylorkoeffizienten einschließt.  $j \leq p$  muss erfüllt sein.

**MPivector get\_all\_derivative(const MPitaylor& x);**

Zurückgegeben wird ein Vektor vom Typ `MPivector` der Dimension  $p + 1$ , welcher alle  $p + 1$  Ableitungen bezüglich des Objekts `x` enthält. Beachten Sie, dass `x` in seinem Vektor `tay1` nur die Taylorkoeffizienten enthält. Die Ableitungen müssen daher aus diesen Taylor-Koeffizienten neu berechnet werden.

**MpfiClass get\_j\_derivative(const MPitaylor& x, int j);**

Zurückgegeben wird ein Intervall vom Typ `MpfiClass`, welches die  $j$ -te Ableitung einschließt.  $j \leq p$  muss erfüllt sein.

**int get\_order(const MPitaylor& x);**

Zurückgegeben wird die maximale Ordnung  $p$  der zu berechnenden Taylorkoeffizienten. Die Dimension des Vektors `tay1` ist damit  $p + 1$ .

### 7.5.8. Ausgabe

**void print\_MPitaylor(const MPitaylor& x);**

Gibt alle  $p + 1$  Taylorkoeffizienten des Objekts `x` auf dem Bildschirm aus.

**void print\_MPderivatives(const MPitaylor& x);**

Gibt alle  $p + 1$  Ableitungen des Objekts `x` auf dem Bildschirm aus.

**static void MpfrClass::SetBase (int b);**

Setzt die Basis  $2 \leq b \leq 36$  NUR für das Ausgabeformat. Gerechnet wird stets im Binärformat!

**static const int MpfrClass::GetBase ();**

Liefert die Basis  $2 \leq b \leq 36$  NUR für das Ausgabeformat. Gerechnet wird stets im Binärformat!

## 7.5.9. Programme

Die Anwendung der Taylorarithmetik zeigt das nachfolgende Programm MPFR-17.cpp, in dem für ein Polynom 4. Grades alle Taylorkoeffizienten und Ableitungen an einer gewählten Stelle  $x_0$  berechnet und auf dem Bildschirm ausgegeben werden.

```

1 //=====
2 // MPFR-17.cpp; Example: Taylor Arithmetic;
3 //-----
4 // Inclusions of the Taylor coefficients and derivatives up to order p=5 for
5 //      P(x) = 2x^4 + x^3 + 4x^2 - 3x + 2
6 // at points of expansion x0 included by the interval z.
7 // A starting interval including the point of expansion x0 must be entered.
8 //=====
9
10 #include "mpficlass.hpp"
11 #include <iostream>      // Input | output
12
13 using namespace cxsc;
14 using namespace MPFI;
15
16 int main()
17 {
18     PrecisionType prec=100;
19     MPitaylor::SetCurrPrecision(prec);
20     int p = 5;      // Order of expansion
21     MpfiClass z;   // Interval to include the point of expansion
22     MPitaylor P;   // Default constructor
23
24     while(1) {
25         cout << endl << "Inclusion of point of expansion x0; [x0,x0] = ? ";
26         cin >> z;
27         MPitaylor x(p,z); // Constructor
28         P = (((2.0*x + 1) // Polynomial of order 4.
29             *x + 4)
30             *x - 3.0)
31             *x + 2.0;
32         cout.precision(P.GetPrecision()/3.321928095);
33         print_MPitaylor(P); // Output of Taylor coefficients.
34
35         MPivector derivative(0,p); // interval vector with
36         for(int i=0; i<=p; i++)
37             derivative[i] = get_j_derivative(P,i); // Derivatives
38         cout << "Inclusions of the derivatives up to order 5 " << endl;
39         for(int i=0; i<=p; i++) // Output of the derivatives
40             cout << i<<"th derivative: " << derivative[i] << endl;
41     }
42 } // main

```

Das Programm liefert die Ausgabe:

```

Inclusion of point of expansion x0; [x0,x0] = ? [2,2]
Output MPitaylor of order 5
i = 0 component: [5.200000000000000000000000000000e1,5.200000000000000000000000000000e1]
i = 1 component: [8.900000000000000000000000000000e1,8.900000000000000000000000000000e1]
i = 2 component: [5.800000000000000000000000000000e1,5.800000000000000000000000000000e1]
i = 3 component: [1.700000000000000000000000000000e1,1.700000000000000000000000000000e1]
i = 4 component: [2.000000000000000000000000000000,2.000000000000000000000000000000]
i = 5 component: [0,-0]

Inclusions of the derivatives up to order 5
0th derivative: [5.200000000000000000000000000000e1,5.200000000000000000000000000000e1]
1th derivative: [8.900000000000000000000000000000e1,8.900000000000000000000000000000e1]
2th derivative: [1.160000000000000000000000000000e2,1.160000000000000000000000000000e2]
3th derivative: [1.020000000000000000000000000000e2,1.020000000000000000000000000000e2]
4th derivative: [4.800000000000000000000000000000e1,4.800000000000000000000000000000e1]
5th derivative: [0,-0]

```







## 7.6. Komplexe, eindimensionale Taylor Arithmetik

Die komplexe, eindimensionale Taylorarithmetik liefert Einschließungen aller Taylorkoeffizienten und Ableitungen einer analytischen Funktion  $f : D_f \rightarrow \mathbb{C}$  einer Variablen an der Stelle  $z_0 \in D_f \subseteq \mathbb{C}$ . Dabei kann  $z_0 = x_0 + i \cdot y_0$  durch ein Intervall eingeschlossen werden, wenn  $z_0$  im verwendeten Zahlenraster nicht darstellbar ist. Die in **C-XSC** implementierte Taylorarithmetik kombiniert die symbolische Differentiation mit der Intervallrechnung und liefert dadurch sehr effektiv auch für Taylorkoeffizienten sehr hoher Ordnungen enge Einschließungen, wenn die interne Präzision hinreichend groß gewählt wird.  $f(z)$  besitzt für den Entwicklungspunkt  $z_0$  die Taylorreihe

$$f(z) = \sum_{k=0}^{\infty} \frac{f^{(k)}(z_0)}{k!} \cdot (z - z_0)^k, \quad (f)_k := \frac{f^{(k)}(z_0)}{k!} \in \mathbb{C},$$

wobei mit  $(f)_k$  der Taylorkoeffizient der Ordnung  $k$  bezeichnet wird. Literaturhinweise findet man im Abschnitt 7.5 auf Seite 139. Da die Rekursionsformeln zur Berechnung der  $(f)_k$  für beliebig oft differenzierbare reelle und für die entsprechenden analytischen Funktionen identisch sind, können die entsprechenden reellen Algorithmen direkt für die komplexen Berechnungen übernommen werden. Man muss lediglich darauf achten, dass die den Entwicklungspunkt  $z_0$  einschließenden Intervalle im Definitionsbereich  $D_f$  liegen. In diesem Abschnitt werden die theoretischen Grundlagen nicht weiter behandelt. Beschrieben werden nur die zur Verfügung stehenden Werkzeuge, mit denen in der Programmierumgebung **C-XSC** die Taylorkoeffizienten und Ableitungen in fast beliebiger Präzision berechnet werden können.

Ein Objekt der Klasse `MPcitaylor` besitzt einen Vektor `tay1` vom Typ `MPcivector` mit den Indices von 0 bis  $p$ , wobei  $0 \leq p \leq 500$  die gewählte Ordnung der zu berechnenden Taylorkoeffizienten ist. Die Taylorkoeffizienten  $(f)_k$  werden dabei durch die Komponenten vom Typ `MpfciClass` dieses Vektors eingeschlossen. Der maximale Wert  $p = 500$  kann bei Bedarf auch noch größer gewählt werden. Mit  $p = 0$  erhält man lediglich eine Einschließung für  $f(z_0) = (f)_0$ . Die Präzision kann für alle Intervall-Taylorkoeffizienten nach Bedarf einheitlich festgelegt werden. Der Vektor `tay1` einer unabhängigen Variablen<sup>1</sup>  $z$  vom Typ `MPcitaylor` hat die Form

$$\text{tay1} = (z, [1], [0], [0], [0], [0], [0]), \quad \text{hier mit } p = 6,$$

und der entsprechende Vektor einer Konstanten  $c$  vom Typ `MpfciClass` ist gegeben durch

$$\text{tay1} = (c, [0], [0], [0], [0], [0]), \quad \text{hier mit } p = 5.$$

Werden die Vektoren für Summe, Differenz, Produkt und Quotient dieser Objekte berechnet, so müssen die Ordnungen  $p$  aller dieser Objekte übereinstimmen. Für alle Elementarfunktionen der nachfolgenden Tabelle gibt es entsprechende Formeln, [11], mit denen alle Taylorkoeffizienten  $(f)_k$ ,  $k = 0, 1, \dots, p$ , z.B für  $f(z) = \sin(z)$ , berechnet werden können. Mithilfe dieser Funktionen und den vier Grundoperationen lassen sich dann beliebige Ausdrücke  $T(z)$  definieren, für die an Entwicklungspunkten  $z_0 \in \mathbb{C}$  alle Taylorkoeffizienten oder Ableitungen bis zur Ordnung  $p$  garantiert eingeschlossen werden können, falls  $T(z)$  an der Stelle  $z_0$  analytisch ist. Wählt man für  $z_0$  ein einschließendes Intervall  $[z] \ni z_0$ , z.B. weil  $z_0$  im Zahlenraster nicht darstellbar ist, so ist die berechnete Einschließung für  $(T)_k$  auch die Einschließung aller Taylorkoeffizienten der Ordnung  $k$  für jedes andere  $\hat{z}_0 \in [z]$ , mit  $\hat{z}_0 \neq z_0$ . Um möglichst gute Einschließungen der Taylorkoeffizienten zu berechnen, sollte  $[z]$  daher möglichst eng gewählt werden. Bei höherer Ordnung  $p$  muss oft mit einer hinreichend großen Präzision `prec` gerechnet werden, um die bekannten Intervallüberschätzungen bei der internen Intervallauswertung klein zu halten.

<sup>1</sup>Beachten Sie, dass  $z$  hier vom Typ `MPcitaylor` ist, während  $z \in \mathbb{C}$  eine komplexe Variable bezeichnet.

### 7.6.1. Elementarfunktionen

Tabelle 7.3.: Funktionen vom Typ `MPcittaylor` zur Taylor-Arithmetik,  $z$  vom Typ `MPcittaylor`,  $w$  vom Typ `Mpfciclass`;

Funktion	Aufruf	Funktion	Aufruf
$z^2$	<code>sqr(z)</code>	$\log_{10}(z)$	<code>log10(z)</code>
$\sqrt{z}$	<code>sqr(z)</code>	$\ln(1+z)$	<code>lnp1(z)</code>
$\sqrt[n]{z}$ , $n = 2, 3, \dots$	<code>sqr(z,n)</code>	$\sin(z)$	<code>sin(z)</code>
$\sqrt{1+z^2}$	<code>sqrtp1x2(z)</code>	$\cos(z)$	<code>cos(z)</code>
$\sqrt{1-z^2}$	<code>sqrtp1mx2(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z^2-1}$	<code>sqrtpx2m1(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{z+1}-1$	<code>sqrtp1m1(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$z^w$	<code>pow(z,w)</code>	$\arccos(z)$	<code>acos(z)</code>
$z^n$ , $n \in \mathbb{Z}$	<code>power(z,n)</code>	$\arctan(z)$	<code>atan(z)</code>
$e^z$	<code>exp(z)</code>	$\operatorname{arccot}(z)$	<code>acot(z)</code>
$2^z$	<code>exp2(z)</code>	$\sinh(z)$	<code>sinh(z)</code>
$10^z$	<code>exp10(z)</code>	$\cosh(z)$	<code>cosh(z)</code>
$e^z - 1$	<code>expm1(z)</code>	$\tanh(z)$	<code>tanh(z)</code>
$e^{z^2}$	<code>expx2(z)</code>	$\operatorname{coth}(z)$	<code>coth(z)</code>
$e^{z^2} - 1$	<code>expx2m1(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$e^{-z^2}$	<code>expmx2(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$e^{-z^2} - 1$	<code>expmx2m1(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$\ln(z)$	<code>ln(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>
$\log_2(z)$	<code>log2(z)</code>		

## 7.6.2. Konstruktoren / Destruktoren

### 7.6.2.1. Konstruktoren

```
MPcitaylor ( );
```

Der Default-Konstruktor legt ein neues MPcitaylor-Objekt in der Current-Precision an. Aufruf: `MPcitaylor f;`  
Es wird kein Vektor `tayl` und damit auch keine Ordnung `p` dieses Vektors initialisiert.

```
MPcitaylor (const MPcitaylor& g, PrecisionType prec);
```

Mit dem Aufruf: `MPcitaylor f(g, prec);` legt der Copy-Konstruktor mithilfe von `g` ein neues MPcitaylor-Objekt `f` an, wobei die Komponenten des Vektors von `f` die Komponenten des Vektors von `g` in der neuen Präzision `prec` einschließen. Wird `prec` nicht gesetzt, so werden die Komponenten des Vektors von `g` in der voreingestellten Current-Precision eingeschlossen.

Alle folgenden Konstruktoren erzeugen ein MPcitaylor-Objekt `f`. Der zugehörige Vektor `tayl` enthält `p + 1` Komponenten vom Typ `Mpfciclass` mit der Präzision `prec`.

Der Aufruf: `MPcitaylor f(5, Mpfciclass( Mpfciclass(1)/10, Mpfciclass(-1) ), 400);` erzeugt den Vektor

`tayl = {( [9.999e-2, 1.001e-1], [-1, -1] ), ([1, 1], [0, 0]), (0, 0), (0, 0), (0, 0), (0, 0) }`, mit  $5 + 1 = 6$  komplexen Komponenten vom Typ `Mpfciclass`, wobei hier die Dezimalstellen aus Platzgründen verkürzt angegeben sind. Die 1. Komponente ist eine Einschließung des nicht-darstellbaren komplexen Entwicklungspunktes  $z_0 = 0.1 - i$  mit ca.  $400/3.321928095 \approx 120$  Dezimalstellen.

Wird beim Konstruktoraufruf `prec` nicht angegeben, so erhalten alle Komponenten von `tayl` die voreingestellte Current-Precision. Falls diese mit `SetCurrPrecision` nicht gesetzt wurde, erhalten die Komponenten die Präzision `prec = 53`.

```
MPcitaylor (int p, const Mpfciclass& x, PrecisionType prec);  
MPcitaylor (int p, const cinterval& x, PrecisionType prec);  
MPcitaylor (int p, const Mpfciclass& x, PrecisionType prec);  
MPcitaylor (int p, const complex& x, PrecisionType prec);  
MPcitaylor (int p, const Mpfciclass& x, PrecisionType prec);  
MPcitaylor (int p, const interval& x, PrecisionType prec);  
MPcitaylor (int p, const MpfrClass& x, PrecisionType prec);  
MPcitaylor (int p, const real& x, PrecisionType prec);  
MPcitaylor (int p, const double& x, PrecisionType prec);  
MPcitaylor (int p, int x, PrecisionType prec);
```

### 7.6.2.2. Destruktor

```
~MPcitaylor () { };
```

Der Speicher für das Objekt wird freigegeben. Dieser Destruktor muss nicht explizit aufgerufen werden!

### 7.6.3. Zuweisungs-Operatoren

Unabhängig von der Präzision des rechten Operanden `op` erhalten bei allen folgenden Zuweisungsoperatoren alle Vektorkomponenten des linken Operanden stets die Current-Precision und schließen die `op`-Werte optimal ein.

```
MPcitaylor operator = (const MPcitaylor& op);
MPcitaylor operator = (const MpfcClass& op);
MPcitaylor operator = (const cinterval& op);
MPcitaylor operator = (const MpfcClass& op);
MPcitaylor operator = (const complex& op);
MPcitaylor operator = (const MpfiClass& op);
MPcitaylor operator = (const interval& op);
MPcitaylor operator = (const MpfrClass& op);
MPcitaylor operator = (const real& op);
MPcitaylor operator = (const double& op);
MPcitaylor operator = (int op);
```

Nur beim ersten Operator wird die Ordnung des linken Operanden gleich der Ordnung von `op`, sonst bleibt die Ordnung des linken Operanden erhalten.

### 7.6.4. Arithmetische Operatoren

Für alle arithmetischen Operationen mit Operanden vom Typ `MPcitaylor` gilt:

Die exakten Ergebnisse für die Komponenten des Vektors `tayl` werden unabhängig von der Präzision der `tayl`-Komponenten der beiden rechten Operanden in der voreingestellten Current-Präzision optimal eingeschlossen.

Sind die rechten Operanden beide vom Typ `MPcitaylor`, so müssen deren Ordnungen übereinstimmen, und der Ergebnisvektor `tayl` übernimmt diese Ordnung. Sind die Ordnungen der beiden rechten Operanden verschieden, so erfolgt eine entsprechende Fehlermeldung.

Ist einer der beiden rechten Operanden nicht vom Typ `MPcitaylor`, so sind für diesen Operanden nur die folgenden Typen zugelassen:

`MpfcClass, cinterval, MpfcClass, complex,`  
`MpfiClass, interval, MpfrClass, real, double, int.`

Alle Operanden dieser Typen werden dabei stets als Konstanten interpretiert!

Für die Addition existieren beispielsweise die Operatoren:

```
MPcitaylor operator + (const MPcitaylor& op1, const MPcitaylor& op2);
MPcitaylor operator + (const MPcitaylor& op1, int op2);
MPcitaylor operator + (complex& op1, const MPcitaylor& op2);
```

### 7.6.5. Unabhängige Variablen und Konstanten

Zur Definition von unabhängigen Variablen und Konstanten siehe Seite 147. Unabhängige Variablen werden generiert durch die Funktion

```
MPcitaylor var_MPcitaylor(int ord, const MpfcClass& x, PrecisionType prec);
```

Für den 2. Parameter stehen folgende Datentypen zur Verfügung:

`MpfcClass, cinterval, MpfcClass, complex,`  
`MpfiClass, interval, MpfrClass, real, double, int.`

Konstanten werden generiert durch die Funktion

```
MPcitaylor const_MPcitaylor(int ord, const MpfiClass& x, PrecisionType prec);
```

Für den 2. Parameter stehen folgende Datentypen zur Verfügung:

```
Mpfciclass, cinterval, MpfcClass, complex,  
MpfiClass, interval, MpfrClass, real, double, int.
```

Für beide obigen Funktionen gilt:

Die Vektorkomponenten des Rückgabewertes erhalten die Präzision `prec`. Wird dieser 3. Parameter nicht gesetzt, so erhalten die Vektorkomponenten die Current-Präzision. Die Ordnung des Rückgabevektors `tayl` wird durch den 1. Parameter `ord` festgelegt.

### 7.6.6. Precision Handling

```
static const PrecisionType GetCurrPrecision ();
```

Gibt die aktuelle Current-Precision in Bits zurück.

```
static void SetCurrPrecision (PrecisionType prec);
```

Setzt die Current-Precision in `MpfrClass` auf `prec`. Diese Current-Precision wird dann automatisch auch in `class MPcitaylor` benutzt. Wird die Current-Precision nicht gesetzt, so wird in beiden Klassen mit der Default-Precision von 53 Bits gerechnet. Das Setzen der Current-Precision hat auf die Präzision der bis dahin benutzten Objekte **keinerlei** Einfluss.

```
PrecisionType GetPrecision () const;
```

Diese Memberfunktion gibt die Präzision des aktuellen Objekts in Bits zurück. Als Beispiel entsprechen dabei 302 Bits  $302/\log_2(10) \approx 91$  Dezimalstellen.

```
void SetPrecision (PrecisionType prec);
```

Diese Memberfunktion setzt die Präzision der aktuellen Komponenten von `tayl` auf `prec`. Die Komponentenwerte bleiben dabei nicht erhalten.

```
void RoundPrecision (PrecisionType prec);
```

Diese Memberfunktion schließt die Komponenten des aktuellen Vektors `tayl` mit der neuen Präzision `prec` ein. Sollte die Präzision der ursprünglichen Komponenten größer sein als `prec`, so erhält man eine gröbere Einschließung. Ist die Präzision jedoch kleiner als `prec`, so werden die restlichen binären Stellen mit Nullen aufgefüllt, so dass die ursprünglichen Werte der Taylor-Komponenten erhalten bleiben.

### 7.6.7. Zugriff auf die Vektorkomponenten

```
MPcivector get_all_coef(const MPcitaylor& x);
```

Zurückgegeben wird ein Vektor vom Typ `MPcivector` der Dimension `p + 1`, welcher alle `p + 1` Taylor-Koeffizienten des Objekts `x` enthält.

```
Mpfciclass get_j_coef(const MPcitaylor& x, int j);
```

Zurückgegeben wird ein komplexes Intervall vom Typ `Mpfciclass`, welches den `j`-ten Taylorkoeffizienten einschließt.  $j \leq p$  muss erfüllt sein.

```
MPcivector get_all_derivative(const MPcivector& x);
```

Zurückgegeben wird ein Vektor vom Typ `MPcivector` der Dimension  $p + 1$ , welcher alle  $p + 1$  Ableitungen bezüglich des Objekts `x` enthält. Beachten Sie, dass `x` in seinem Vektor `tay1` nur die Taylorkoeffizienten enthält. Die Ableitungen müssen daher aus diesen Taylor-Koeffizienten neu berechnet werden.

```
Mpfciclass get_j_derivative(const MPcivector& x, int j);
```

Zurückgegeben wird ein komplexes Intervall vom Typ `Mpfciclass`, welches die  $j$ -te Ableitung einschließt.  $j \leq p$  muss erfüllt sein.

```
int get_order(const MPcivector& x);
```

Zurückgegeben wird die maximale Ordnung  $p$  der zu berechnenden Taylorkoeffizienten. Die Dimension des Vektors `tay1` ist damit  $p + 1$ .

### 7.6.8. Ausgabe

```
void print_MPcivector(const MPcivector& x);
```

Gibt die Einschließungen aller  $p + 1$  komplexen Taylorkoeffizienten des Objekts `x` auf dem Bildschirm aus.

```
void print_MPderivatives(const MPcivector& x);
```

Gibt die Einschließungen aller  $p + 1$  komplexen Ableitungen des Objekts `x` auf dem Bildschirm aus.

```
static void MpfrClass::SetBase (int b);
```

Setzt die Basis  $2 \leq b \leq 36$  NUR für das Ausgabeformat. Gerechnet wird stets im Binärformat!

```
static const int MpfrClass::GetBase ();
```

Liefert die Basis  $2 \leq b \leq 36$  NUR für das Ausgabeformat. Gerechnet wird stets im Binärformat!

### 7.6.9. Programme

Die Anwendung der komplexen Taylorarithmetik zeigt das folgende Programm `MPFR-19.cpp`, in dem für ein komplexes Polynom 4. Grades alle Taylorkoeffizienten und Ableitungen an einer gewählten Stelle  $z_0$  eingeschlossen und auf dem Bildschirm ausgegeben werden.



```

1 //=====
2 // MPFR-19.cpp; Example: Complex Taylor Arithmetic;
3 // Enclosures of the Taylor coefficients and derivatives up to order p=5 for
4 //       $P(z) = 2*z^4 + z^3 + (4+i)*z^2 - 3z + 2-i$ 
5 // at points of expansion  $z_0=x+i*y$  included by an interval  $[z]$ .
6 // A starting interval including the point of expansion  $z_0$  must be entered.
7 //=====
8
9 #include "mpfciclass.hpp"
10 #include <iostream>      // Input | output
11
12 using namespace cxsc;
13 using namespace MPFI;
14
15 int main() {
16     PrecisionType prec=70;
17     MPCitaylor::SetCurrPrecision(prec);
18     int p = 5;          // Order of expansion
19     Mpfciclass z;      // Interval to include the point of expansion
20     MPCitaylor P;      // Default constructor
21
22     while(1) {
23         cout << endl << "Inclusion of point of expansion x0; ([x,x],[y,y]) = ? ";
24         cin >> z;
25         MPCitaylor x(p,z); // Constructor
26
27         P = (((2.0*x + 1) // Polynomial of order 4.
28             *x + complex(4,1))
29             *x - 3.0)
30             *x + complex(2,-1);
31         cout.precision(P.GetPrecision()/3.321928095);
32         print_MPCitaylor(P); // Output of Taylor coefficients.
33
34         MPcivector derivative(0,p); // interval vector with
35                                     // components from index 0 up to the order p.
36         for(int i=0; i<=p; i++)
37             derivative[i] = get_j_derivative(P,i); // Derivatives
38         cout << "Enclosures of the derivatives up to order 5 " << endl;
39         for(int i=0; i<=p; i++) // Output of the derivatives
40             cout << i << "th derivative: " << derivative[i] << endl; }
41 } // main

```

Das Programm liefert die Ausgabe:

```

Inclusion of point of expansion x0; ([x,x],[y,y]) = ? ([2,2],[-3,-3])
Output MPCitaylor of order 5
i = 0 Taylor coefficient: ([-2.9600000000000000000000000000000000 e2, -2.9600000000000000000000000000000000 e2],
[1.8600000000000000000000000000000000 e2, 1.8600000000000000000000000000000000 e2])
i = 1 Taylor coefficient: ([-3.6400000000000000000000000000000000 e2, -3.6400000000000000000000000000000000 e2],
[-1.2800000000000000000000000000000000 e2, -1.2800000000000000000000000000000000 e2])
i = 2 Taylor coefficient: ([-5.0000000000000000000000000000000000 e1, -5.0000000000000000000000000000000000 e1],
[-1.5200000000000000000000000000000000 e2, -1.5200000000000000000000000000000000 e2])
i = 3 Taylor coefficient: ([1.7000000000000000000000000000000000 e1, 1.7000000000000000000000000000000000 e1],
[-2.4000000000000000000000000000000000 e1, -2.4000000000000000000000000000000000 e1])
i = 4 Taylor coefficient: ([2.0000000000000000000000000000000000, 2.0000000000000000000000000000000000], [0,-0])
i = 5 Taylor coefficient: ([0,-0], [0,-0])

Enclosures of the derivatives up to order 5
0th derivative: ([-2.9600000000000000000000000000000000 e2, -2.9600000000000000000000000000000000 e2],
[1.8600000000000000000000000000000000 e2, 1.8600000000000000000000000000000000 e2])
1th derivative: ([-3.6400000000000000000000000000000000 e2, -3.6400000000000000000000000000000000 e2],
[-1.2800000000000000000000000000000000 e2, -1.2800000000000000000000000000000000 e2])
2th derivative: ([-1.0000000000000000000000000000000000 e2, -1.0000000000000000000000000000000000 e2],
[-3.0400000000000000000000000000000000 e2, -3.0400000000000000000000000000000000 e2])
3th derivative: ([1.0200000000000000000000000000000000 e2, 1.0200000000000000000000000000000000 e2],
[-1.4400000000000000000000000000000000 e2, -1.4400000000000000000000000000000000 e2])
4th derivative: ([4.8000000000000000000000000000000000 e1, 4.8000000000000000000000000000000000 e1], [0,-0])
5th derivative: ([0,-0], [0,-0])

```

Mit Hilfe der Identität

$$T_1(z) := \cot(4z) \equiv \frac{\cot^4 z - 6 \cot^2 z + 1}{4 \cot^3 z - 4 \cot z} =: T_2(z)$$

zeigt das folgende Programm MPFR-20, dass bei ungünstiger Wahl des Funktionsterms  $T_2(z)$  und zu kleiner Präzision `prec=200` z.B. die 5. Ableitung an der Stelle  $z_0 = 2^{-200} + i \cdot 2^{-100}$  viel zu grob eingeschlossen wird, während mit  $T_1(z)$  die gleiche Ableitung bei gleicher Präzision fast optimal eingeschlossen wird.

```

1 //=====
2 //-----
3 // MPFR-20.cpp; Example: Complex Taylor Arithmetic;
4 //-----
5 // To get best enclosures of the derivatives (Taylor coefficients) the
6 // choice of the optimal function term is essential:
7 // cot(4z) = (cot^4(z) - 6*cot^2(z) + 1)/(4*cot^3(z)-4*cot(z));
8 //-----
9 //=====
10
11 #include "mpfciclass.hpp"
12 #include <iostream> // Input | output
13
14 using namespace cxsc;
15 using namespace MPFI;
16
17 int main()
18 {
19     PrecisionType prec=200;
20     MPCitaylor::SetCurrPrecision(prec);
21     int p = 5; // Order of expansion
22     MpfiClass re(exp2(-200)), im(exp2(-100));
23     MpfiClass z(re,im);
24     MPCitaylor y; // Default constructor
25     MPCitaylor x(p,z); // x: Variable of type MPCitaylor
26
27     y = cot(4*x);
28     cout.precision(y.GetPrecision()/3.321928095);
29     cout << "5-th derivative: " << get_j_derivative(y,5) << endl;
30
31     y = (power(cot(x),4)-6*power(cot(x),2)+1)/(4*power(cot(x),3)-4*cot(x));
32     cout << "5-th derivative: " << get_j_derivative(y,5) << endl;
33 } // main

```

Das Programm liefert die Ausgabe:

```

5-th derivative: ([1.24485467066429788755372235910734834530373386967273106999850e182,
1.24485467066429788755372235910734834530373386967273106999887e182],
[5.89210309421305536602374145428967847398995568287751662106821e152,
5.89210309421305536602374145428967847398995568287751662107157e152])
5-th derivative: ([-1.14181716389186769287770038422182448327311962913481632129061e775,
1.09960973422552681253501467712448598966940730546911198376147e775],
[-1.14168379364897743492204414456301343922534958133554464565727e775,
1.14168379364897743492204414456157393246606487206729495198471e775])

```

### Anmerkungen:

1. Erst mit `prec = 400` erhält man mit Hilfe von  $T_2(z)$  eine Einschließung der 5. Ableitung mit hinreichend vielen übereinstimmenden Dezimalstellen.
2. Ein wesentlicher Vorteil einer Taylorarithmetik mit beliebiger Präzision besteht aber darin, dass man bei ungeschickter Wahl des Funktionsterms immer noch brauchbare Einschließungen erhält, wenn die Präzision nur hinreichend groß gewählt wird.

## 7.7. Nullstellen komplexer Ausdrücke

In der Dissertation von W. Krämer wird ein Algorithmus zur Einschließung von Nullstellen komplexwertiger Ausdrücke angegeben, [38],[13]. Benötigt wird eine möglichst gute Nullstellenapproximation, die mithilfe eines vereinfachten Newton-Verfahrens weiter verbessert wird, um danach die eindeutig bestimmte Nullstelle garantiert einzuschließen. *Eindeutig* bedeutet hier, dass das einschließende Intervall **genau eine** Nullstelle enthält. Mit dem folgenden Programm können mit geeigneten Näherungen die drei einfachen Nullstellen  $z_1 = 2 + i$ ,  $z_2 = 3 - i$ ,  $z_3 = 4i$  der Funktion

$$f(z) = \arctan((z - a) \cdot \ln(z^2 - 5 \cdot z + 8 + i))$$

garantiert eingeschlossen werden, wobei zur Realisierung von  $z_3 = 4i$  der Parameter  $a = 4i$  zu setzen ist. Im Programm muss die 1. Ableitung

$$f'(z) = \frac{\ln(z^2 - 5 \cdot z + 8 + i) + \frac{(z - a) \cdot (2 \cdot z - 5)}{z^2 - 5 \cdot z + 8 + i}}{1 + (z - a)^2 \cdot (\ln(z^2 - 5 \cdot z + 8 + i))^2}$$

für das Newton-Verfahren **nicht** explizit definiert werden, da  $f(z)$  und  $f'(z)$  mit Hilfe der durch komplexe Taylor-Arithmetik definierten Funktion

```
MPcittaylor F(const MPcittaylor Z)
{
    MpfcClass a(0,4);
    return atan( (Z-a)*ln(sqr(Z) - 5*Z + MpfcClass(8,1)) );
}
```

direkt eingeschlossen werden können.

```
1 //=====
2 //-----
3 // MPFR-11.cpp; Complex (interval) Newton method using complex interval
4 // Automatic Taylor Arithmetic (ATA) to compute first derivatives
5 //-----
6 //=====
7
8 #include "mpfciclass.hpp"
9
10 using namespace MPFR;
11 using namespace MPFI;
12 using namespace cxsc;
13 using namespace std;
14
15 // The zeros of the function F() are to be enclosed. Zeros are:
16 // a (see parameter in function definition), 2+i, and 3-i.
17
18 MPcittaylor F(const MPcittaylor& Z) // Definition of F(z) and F'(z) via ATA
19 {
20     MpfcClass a(0,4); // The parameter a is set to 0 + 4i
21     return atan( (Z-a)*ln(sqr(Z) - 5*Z + MpfcClass(8,1)) );
22 };
23
24
25 int main(void)
26 {
27     PrecisionType prec= 99; // number of bits to be used in arithmetic operations
28     MPcittaylor::SetCurrPrecision(prec);
29     cout << "Current-precision: about " << Round(prec/3.321928095)
30         << " decimals" << endl;
31     prec= prec/3.321928095; // precision in decimal digits
32     cout.precision(prec); // perform output with prec decimal digits
```

```

33
34 MpfciClass Y, Z0, Zb, Z, D, CH, Zi, C; // complex intervals
35 MpfcClass approx; // complex point
36 int k, kmax(9); // iteration counter, maximum number of iteration steps
37 MpfrClass eps= MpfrClass(0.125); // used for epsilon inflation
38
39 while (true) {
40     cout << "The zeros of function F are: 4i, 2+i, and 3-i" << endl;
41     cout << "Approximation of a complex zero
42         ( try e.g. (0.1,4.3) or (-0.2,3.8) for the zero 4i )? ";
43     cin >> approx;
44     Zi= MpfcClass(approx);
45     k= 0;
46
47     MPCitaylor FD; // function values and values of derivative using ATA
48     MpfcClass FZi; // function values
49     MpfcClass DFZi; // values of first derivative
50     do { // try to improve approximation
51         k++;
52         cout << "k: " << k << endl;
53         Z0 = Zi;
54         FD = F( MPCitaylor(1, Zi) ); // compute function value
55                                     // and first derivative via ATA
56         FZi = get_j_derivative(FD,0); // take function value
57         DFZi = get_j_derivative(FD,1); // take first derivative
58
59         Zi= MpfcClass( mid(Zi - FZi/DFZi) ); // perform Newton step
60         Z = MpfcClass( MpfiClass(Re(Sup(Z0)),Re(Sup(Zi))),
61                       MpfiClass(Im(Sup(Z0)),Im(Sup(Zi))) );
62     } while (common_decimals(Re(Z))<prec/3 && common_decimals(Im(Z))<prec/3
63             && k<kmax );
64     cout << "Improved approximation of complex zero: " << endl << Sup(Zi) << endl;
65
66     C = MpfcClass(1/Sup(DFZi)); // C is a point interval
67     FD = F( MPCitaylor(1, Zi) ); // compute function value and first derivative
68     FZi = get_j_derivative(FD,0); // take function value
69     Y = -C*FZi;
70     Z = Y; // Y,C are the initial values of the interval iteration (verification)
71     k = 0;
72
73     do
74     {
75         k++;
76         cout << " k: " << k << endl;
77         Zb = Blow(Z, eps); // epsilon inflation of Z
78         CH = Zi | Zi+Zb; // convex hull
79
80         FD = F( MPCitaylor(1, CH) ); // function value and 1. derivative via ATA
81
82         D = 1 - C*get_j_derivative(FD,1);
83         Z = Y + D*Zb;
84     } while ( (k<kmax) && !in(Z,Zb) );
85
86     if ( in(Z,Zb) )
87     { // inclusion verified
88         cout << "Verification successful, inclusion of a simple zero: "
89             << endl << Zi + Z << endl;
90     }
91     else cout << "Verification failed" << endl;
92 }
93
94 return 0;
95 }

```

Das Programm liefert die Ausgabe:

```
Current-precision: about 30 decimals
The zeros of function F are: 4i, 2+i, and 3-i
Approximation of a complex zero
  ( try e.g. (0.1,4.3) or (-0.2,3.8) for the zero 4i )? (-0.2,3.8)
k: 1
k: 2
k: 3
k: 4
k: 5
k: 6
Improved approximation of complex zero:
(4.3951815917420156712814208664e-35, 4.000000000000000000000000000000)
  k: 1
  k: 2
Verification successful, inclusion of a simple zero:
([-8.3560996162686228977764531973e-64,8.3560996162686228977764531973e-64],
[3.999999999999999999999999999999,4.000000000000000000000000000001])
```

#### Anmerkungen:

- Mit den Näherungen  $\zeta_1 = 1.9 + 0.9i$ ,  $\zeta_2 = 2.9 - 0.9i$ ,  $\zeta_3 = 0.1 + 3.9i$  erhält man z.B. mit der Präzision  $\text{prec} = 3000$  Bits  $\approx 3000/3.321928095 \approx 903$  Dezimalstellen nahezu optimale Einschließungen der drei Nullstellen  $z_1 = 2 + i$ ,  $z_2 = 3 - i$ ,  $z_3 = 4i$ , wobei der Realteil 0 von  $z_3$  natürlich nur vergleichsweise grob eingeschlossen werden kann.
- In Zeile 50 wird das vereinfachte Newton-Verfahren so lange angewandt, bis zwei aufeinanderfolgende Nullstellen-Approximationen in mindestens einem Drittel ihrer dezimalen Präzisionsstellen übereinstimmen. Danach erfolgt ab Zeile 73 der Einschließungstest.

## 7.8. Einschließung aller Nullstellen nichtlinearer Funktionen

Die Nullstellenberechnung für differenzierbare, nichtlineare Funktionen ist seit Generationen in der Numerik ein klassisches Arbeitsgebiet. In den meisten Fällen geht man aus von einer Nullstellennäherung und versucht dann, diese Näherung mit iterativen Methoden (Newtonverfahren, Halley's Methode, Fixpunktiteration, Bisektion, ...) nach Möglichkeit weiter zu verbessern. Bis in die 70er Jahre des 20. Jahrhunderts war man dabei der Auffassung, dass es keinen Algorithmus geben kann, mit dem *alle* Nullstellen einer nichtlineare Funktion über einem Intervall mit garantierten Fehlerschranken berechnet werden können, wobei zusätzlich auch noch die Existenz und Eindeutigkeit dieser Nullstellen nachzuweisen ist.

Die Aufgabe besteht darin, zu einer gegebenen, stetig differenzierbaren Funktion  $f : \mathbb{R} \rightarrow \mathbb{R}$  und  $x \in [x_1, x_2] \subseteq \mathbb{R}$  **alle** ihre Nullstellen möglichst eng einzuschließen und die Existenz und Eindeutigkeit dieser Nullstellen nachzuweisen. Den theoretischen Hintergrund findet man in [20], [50], [5], [29]. In diesem Abschnitt werden nur die zur Verfügung stehenden Hilfsmittel beschrieben, mit denen die Nullstellen mit praktisch beliebiger Genauigkeit eingeschlossen werden können. Das zentrale Werkzeug ist dabei die in `mpficlass.cpp` definierte Funktion

```
void AllZeros ( MPddf_FctPtr f,
                MpfiClass   SearchInterval,
                MpfrClass   Tolerance,
                MPivector&  Zero,
                intvector&  Unique,
                int&        NumberOfZeros,
                int&        Err,
                int         MaxNumberOfZeros )
```

### Anmerkungen:

1.  $f$  : Funktion vom Typ `MPDerivType`. Zur Definition siehe z.B. Seite 125 und die Tabelle auf Seite 128.
2. `SearchInterval` : Das Intervall, in dem alle Nullstellen von  $f$  zu berechnen sind.
3. `Tolerance` : Die relative Genauigkeit für die Einschließung der Nullstellen. Mit `Tolerance` wird die Präzision bestimmt, mit der intern gerechnet werden muss, um die vorgegebene relative Genauigkeit zu erreichen.
4. `Zero` : Dieser Vektor liefert die Einschließungen aller Nullstellenkandidaten.
5. `Unique` : `Unique[i]==0` : Existenz und Eindeutigkeit sind nicht garantiert, d.h `Zero[i]` kann entweder keine Nullstelle oder eine oder mehrere einfache oder mehrfache Nullstellen enthalten.  
`Unique[i]==1` : In `Zero[i]` sind Existenz und Eindeutigkeit garantiert.
6. `NumberOfZeros` : Anzahl der gefundenen Nullstellenkandidaten, wobei aber die Eindeutigkeit nicht bei allen Einschließungen garantiert sein muss. Bei symmetrischer Lage der Nullstellen in `SearchInterval` kann die gleiche Nullstelle im Vektor `Zero` in verschiedenen Komponenten mehrfach eingeschlossen und entsprechend gezählt werden. Um diese Symmetrie zu vermeiden, kann bei einigen Programmversionen (Funktion `compute(...)`) das eingegebene Startintervall nach Bedarf nach rechts oder links zum `SearchInterval` aufgebläht werden, so dass dann in nahezu allen Fällen die gleiche Nullstelle nur durch eine einzige Vektorkomponente `Zero[i]` eingeschlossen wird. Durch diese möglichen Aufblähungen können Nullstellen außerhalb von `SearchInterval` eingefangen werden, die aber anschließend automatisch wieder aussortiert werden.

7. Da ein Nullstellenkandidat `Zero[k]` mit `Unique[k] = 0` keine Nullstelle einschließen muss, wird bei einigen Programmversionen (Funktion `compute(...)`) durch einen nochmaligen Aufruf von `AllZeros(f,Zero[k],Tolerance,ZeroV,UNIQUE,NOZ,Error)` getestet, ob `NOZ = 0` erfüllt ist. In diesem Fall ist garantiert, dass `Zero[k]` keine Nullstelle enthält. In den Vektoren `Zero` und `Unique` werden dann die Komponenten `Zero[k]`, `Unique[k]` entsprechend aussortiert.
8. Um die Anzahl der Nullstellenkandidaten weiter zu minimieren, betrachten wir den Fall, dass sich zwei aufeinanderfolgende Nullstellenkandidaten `Zero[k]` und `Zero[k+1]` überlappen. Unter der Bedingung `(Unique[k]==Unique[k+1] && Unique[k]==0)` kann deren Vereinigung `tmp = Zero[k] ∪ Zero[k+1]` keine oder mehrere einfache oder mehrfache Nullstellen einschließen. Nach dem folgenden Aufruf
 

```
AllZeros(f,tmp,Tolerance,ZeroV,UNIQUE,NOZ,Error)
```

 wird dann die Bedingung `(NOZ==1 && UNIQUE[1]==1)` getestet. Ist diese erfüllt, so setzt man `Zero[k]=ZeroV[1]` und `Unique[k]=1`, und in den Vektoren `Zero` und `Unique` werden die Komponenten `Zero[k+1]` und `Unique[k+1]` gelöscht. Ist obige Bedingung nicht erfüllt, so setzt man `Zero[k]=tmp` und löscht in den Vektoren `Zero` und `Unique` ebenfalls die Komponenten `Zero[k+1]` und `Unique[k+1]`. In jedem Fall reduziert sich damit die Anzahl der Nullstellen-Kandidaten um eins.  
 Wenn bei den beiden sich überlappenden Nullstellenkandidaten `Zero[k]` und `Zero[k+1]` die obige Bedingung `(Unique[k]==Unique[k+1] && Unique[k]==0)` nicht erfüllt ist, so wird nach dem Aufruf von `AllZeros(f,tmp,Tolerance,ZeroV,UNIQUE,NOZ,Error)` die Bedingung `(NOZ==1 && UNIQUE[1]==1)` getestet. Und nur, wenn diese erfüllt ist, setzt man `Zero[k]=ZeroV[1]` und `Unique[k]=1`, und zusätzlich werden in den Vektoren `Zero` und `Unique` die Komponenten `Zero[k+1]` und `Unique[k+1]` gelöscht.  
 Durch diese Vorgehensweise wird bei einigen Programmversionen (Funktion `compute()`) die Anzahl der Nullstellenkandidaten minimiert. Den vollständigen Programmcode findet man z.B. ab Seite ....
9. **Err** : Bestimmt die Ausgabe einer möglichen Fehlermeldung.

Der in `AllZeros()` realisierte Algorithmus ist sehr leistungsfähig, und ohne das Auftreten einer Fehlermeldungen werden alle Nullstellen aus dem Startintervall erkannt. Mehrfache Nullstellen werden zwar erkannt und eingeschlossen, die Eindeutigkeit der Nullstelle(n) kann in dieser Einschließung jedoch grundsätzlich nicht nachgewiesen werden. Für diese Einschließungen erfolgt die Meldung *may contain a zero (not verified unique)!* Die gleiche Meldung kann aber auch dann erfolgen, wenn der Kandidat keine Nullstelle enthält oder wenn er eine oder mehrere dicht zusammenliegende einfache Nullstellen enthält, deren Eindeutigkeit, z.B. wegen fast waagerechter Tangente, nicht mehr nachgewiesen werden kann.

Wird bei einer einfachen Nullstelle die Eindeutigkeit nachgewiesen, so erfolgt für diese Nullstelleneinschließung die Meldung *encloses a locally unique zero!*, d.h. außer dieser einfachen Nullstelle gibt es in diese Einschließung keine weitere einfache oder mehrfache Nullstelle.

Die bisherigen Anmerkungen beziehen sich auf die Anwendung der Funktion `AllZeros(...)` bzw. auf eine möglichst optimale Ausgabe der Nullstelleneinschließungen, d.h. auf eine möglichst minimale aber vollständige Anzahl der Nullstellenkandidaten. Die folgenden Anmerkungen beziehen sich auf eine möglichst geschickte Darstellung der zu untersuchenden Funktion, deren Nullstellen zu bestimmen sind. Dabei zeigt sich, dass bei komplizierten Funktionen deren optimale Darstellung oft ein eigenständiges Problem ist.

#### Anmerkungen:

1. Wenn ein Funktionsterm beispielsweise den Ausdruck  $\sqrt{1 + [x] \cdot [x]}$  enthält, so sollte dieser durch  $\sqrt{1 + \text{sqr}([x])}$  ersetzt werden, da sonst mit  $[x] = [-2, 2]$  der Radikand  $1 + [x] \cdot [x]$  die Null enthält und die Quadratwurzel damit eine Fehlermeldung erzeugt. Um zusätzlich einen vorzeitigen Überlauf zu vermeiden, benutzt man am besten die Funktion `sqr1px2(x)`.

2. Ein Funktionsterm sollte die unabhängige Variable  $x$  nur in minimaler Anzahl enthalten. Daher sollte z.B.  $(x^2 - 1)/(x - 1)$  durch den äquivalenten Ausdruck  $x + 1$  ersetzt werden.
3. Ein Polynom  $(1) P_N(x) := a_0 + a_1 \cdot x + a_2 \cdot x^2 + \dots + a_N \cdot x^N$  sollte nur nach dem Horner-Schema

```

r = a[N]*x + a[N-1];
for (int i=N-2; i>=0; i--)
    r = r*x + a[i];

```

mit  $P_N(x) = r$  ausgewertet werden, da das Horner-Schema weniger Rechenoperationen benötigt und die Potenzen  $x^i$  bei der Intervallauswertung zu große Überschätzungen verursachen. Bei aufwendigen Polynomen kann die Darstellung (1) bei der Nullstellenberechnung mit `AllZeros(...)` zum Programmabbruch führen. Dies zeigt, dass die geschickte Darstellung der auszuwertenden Funktion für die erfolgreiche Nullstellenberechnung oft entscheidend ist.

4. Bei komplizierten Nullstellenproblemen (z.B. fast waagerechte oder senkrechte Nullstellen-Tangenten, oder sehr dicht benachbarte einfache Nullstellen) führt oft auch eine geschickte Wahl des Funktionsterms nicht zum Ziel. Der Grund ist dann oft, dass im Funktionsterm die unabhängige Variable  $x$  zu häufig vorkommt und damit bei der Intervallauswertung zu große Überschätzungen verursacht, die auch mit noch so großer Präzision nicht beseitigt werden kann. Abhilfe schafft bei nicht zu breiten Intervallen  $[x]$  bei der Auswertung von  $f[x]$  dann oft die Anwendung der zentralen Form (*centered form*), die für differenzierbare Funktionen  $f$  gegeben ist durch [20],[6],[56],[30]

$$(7.2) \quad f(x) \in f(c) + f'([x]) \cdot ([x] - c), \quad c, x \in [x].$$

Wählt man  $c := \text{mid}([x])$ , so spricht man von der Mittelwertform (*mean-value form*). Zu beachten ist, dass die Güte der Einschließung von  $f(x)$ ,  $x \in [x]$ , nur von  $f'([x])$  abhängt, so dass auch  $f'([x])$  wieder mit Hilfe der Mittelwertform dargestellt werden sollte und damit auch  $f''([x])$  benötigt wird. Die höheren Ableitungen, etwa bis zur Ordnung  $k$ , können dabei mit Hilfe der implementierten Taylorarithmetik bequem berechnet werden.

5. Bei einigen Problemstellungen, z.B. bei der Bestimmung aller relativen Extrema eine zweimal differenzierbaren Funktion  $f(x)$ , benötigt man alle Nullstellen von  $f'(x)$ . Mithilfe der Taylorarithmetik ist es grundsätzlich kein Problem, zu einer  $n$ -mal differenzierbaren Funktion  $f(x)$  die  $n$ -te Ableitung  $f^{(n)}(x)$  zu berechnen und dann mit `AllZeros()` alle Nullstellen von  $f^{(n)}(x)$  für  $x \in [x]$  zu bestimmen. Bei einigen Programmversionen (Funktion `compute()`) kann dies durch die entsprechende interaktive Eingabe von  $n$  realisiert werden.

Die folgenden vier Beispielprogramme zeigen die geschilderten, verschiedenen Anwendungsmöglichkeiten zur Nullstellenbestimmung, dabei werden auch die Grenzbereiche angesprochen, in denen die gesuchten Nullstellen nicht oder nicht mehr vollständig angegeben werden können. Das erste Programm `MPFR-15.cpp` besitzt die einfachste Funktion `void compute(...)`, in der keine Mittelwertform, keine Intervall Aufblähung und keine optimierte Kandidatenausgabe zur Anwendung kommt.

Das folgende Programm `MPFR-15.cpp` berechnet für  $f(x) = e^{-3x} - \sin^3(x)$  die Einschließungen aller auftretenden einfachen Nullstellen im vorgegebenen Startintervall  $[0, 20]$ . Die relative Genauigkeit für die einzuschließenden Nullstellen muss eingegeben werden. Mit z.B. `Tolerance = 1e-400` wird die interne Präzision auf `prec = -1.5 * expo(Tolerance) + 53  $\approx$  2045` festgesetzt, so dass Einschließungen mit etwa  $2045/3.321928095 \approx 616$  korrekten Dezimalstellen berechnet werden. Falls `prec` kleiner ist als die `CurrentPrecision`, wird `prec` auf die `CurrentPrecision` gesetzt.



```

1 // Programm MPFR-15.cpp zur Einschliessung aller Nullstellen
2 // einer Funktion f vom Typ MPDerivType
3 #include <iostream>
4 #include "mpficlass.hpp"
5 #include <stacksz.hpp>
6
7 using namespace cxsc;
8 using namespace std;
9 using namespace MPFR;
10 using namespace MPFI;
11
12 MPDerivType f (const MPDerivType& x)
13 {
14     // return exp(-30.0*sqr(x)) - power(sin(x), 3); // Tolerance = 1e-4000 works
15     return exp(-3.0*x) - power(sin(x), 3);
16 }
17
18 //-----
19 // Function for prompting and reading information to call the function
20 // 'AllZeros()'. This function must be called with the function 'f'
21 // and a string 'Name' containing a textual description of that function.
22 //-----
23
24 void compute ( MPddf_FctPtr f, const char *Name )
25 {
26     MpfiClass   SearchInterval;
27     MpfrClass   Tolerance;
28     MPivector   Zero;
29     intvector   Unique;
30     int         NumberOfZeros, i, Error;
31     PrecisionType prec;
32
33     cout << "Computing all zeros of the function " << Name << endl;
34     cout << "Search interval      : ";
35     cin >> SearchInterval;
36     cout << "Tolerance (relative) : ";
37     cin >> Tolerance;
38     cout << endl;
39
40     if (Tolerance>1e-8) Tolerance = 1e-8;
41     // With the entered 'Tolerance' the corresponding output format is calculated:
42     prec = -expo(Tolerance) + 1;
43     if (prec < 53) prec = 53;
44     cout.precision(prec/3.321928095); // Only for the output format;
45
46     AllZeros(f, SearchInterval, Tolerance, Zero, Unique, NumberOfZeros, Error);
47
48     if (Error) cout << endl << AllZerosErrMsg(Error) << endl;
49     else
50     {
51         for ( i = 1; i <= NumberOfZeros; i++)
52         {
53             cout << Zero[i] << endl;
54             if (Unique[i])
55                 cout << "encloses a locally unique zero!" << endl;
56             else
57                 cout << "may contain a zero (not verified unique)!" << endl;
58         }
59         cout << "Calculated " << NumberOfZeros << " interval enclosure(s)" << endl;
60     }
61 } // compute()

```

```

62 int main(void)
63 {
64     int base;
65     MpfrClass::SetBase (10);
66     base = MpfrClass::GetBase();
67     cout << "Base = " << base << endl;
68
69     compute(f, "EXP(-3x) - POWER(SIN(x),3)");
70
71     return 0;
72 }

```

Beim Programmstart erhält man den folgenden Ausdruck:

```
Base = 10
```

```
Computing all zeros of the function EXP(-3x) - POWER(SIN(x),3)
```

```
Search interval      : [0,20]
```

```
Tolerance (relative) : 1e-34
```

```
[5.885327439818610774324520457029036e-1,5.885327439818610774324520457029037e-1]
```

```
encloses a locally unique zero!
```

```
[3.096363932410646115625840849904024,3.096363932410646115625840849904025]
```

```
encloses a locally unique zero!
```

```
[6.285049273382586533848300969784156,6.285049273382586533848300969784157]
```

```
encloses a locally unique zero!
```

```
[9.424697254738521219115865183918111,9.424697254738521219115865183918112]
```

```
encloses a locally unique zero!
```

```
[1.256637410168936767682201604251579e1,1.256637410168936767682201604251580e1]
```

```
encloses a locally unique zero!
```

```
[1.570796311724721594229036828468796e1,1.570796311724721594229036828468797e1]
```

```
encloses a locally unique zero!
```

```
[1.884955592805117152444424965806380e1,1.884955592805117152444424965806381e1]
```

```
encloses a locally unique zero!
```

```
Calculated 7 interval enclosure(s)
```

### Anmerkungen:

1. Wir erhalten also 7 eindeutig bestimmte Nullstellen, die mit einer Genauigkeit von etwa 221 Dezimalstellen eingeschlossen werden, und der verwendete Algorithmus garantiert wegen der nachgewiesenen Eindeutigkeit, dass es im Innern von  $[0, 20]$  keine weiteren Nullstellen geben kann. Die obige Ausgabe ist auf 34 Dezimalstellen beschränkt.
2. Die Funktion  $g(x) = -\sin^3(x)$  besitzt bei  $x_k = k \cdot \pi$ ,  $k \in \mathbb{Z}$ , nur dreifache Nullstellen mit waagerechter Tangente. Solche mehrfachen Nullstellen werden vom Algorithmus zwar eingeschlossen, die Eindeutigkeit kann aber nicht nachgewiesen werden. Der Nachweis der Eindeutigkeit gelingt nur bei einfachen Nullstellen!
3. Addiert man zu  $g(x) = -\sin^3(x)$  den kleinen positiven Funktionswert  $e^{-3x}$ , so werden die mehrfachen Nullstellen von  $g(x)$  zu einfachen Nullstellen von  $f(x) = e^{-3x} - \sin^3(x)$ , die mit obigem Programm erfolgreich eingeschlossen werden.
4. Wählt man in  $f(x) = e^{-3x} - \sin^3(x)$  den für  $x > 0$  noch viel kleineren Summanden  $e^{-30x^2}$ , so besitzt  $h(x) := e^{-30x^2} - \sin^3(x)$  ebenfalls einfache Nullstellen, jetzt jedoch mit nahezu waagerechten Tangenten, und der Nachweis der Eindeutigkeit dieser Nullstellen in ihren Einschließungen gelingt in  $[0, +20]$  nur mit der viel kleineren Toleranz  $10^{-2000}$ , wodurch aber erhebliche Laufzeiten erforderlich werden.

- $u(x) = \sin^2(x) - 2^{-20000}$  besitzt in  $[3,4]$  zwei sehr eng benachbarte, einfache Nullstellen. Der Algorithmus findet jedoch mit `Tolerance = 1e-34` nur *eine* Einschließung mit der Meldung *may contain a zero (not verified unique)!* Mit dem sehr viel kleineren Wert `Tolerance = 1e-2000` wird eine Einschließung wegen zu vieler interner rekursiver Funktionsaufrufe jedoch nicht gefunden, (Speicherzugriffsfehler). Die gleiche Fehlermeldung erhält man aber auch mit der sehr viel kleineren Fehlerschranke `Tolerance = 1e-40000`, vgl. dazu auch das nächste Beispiel und Seite 170.
- $v(x) = \cos(x) \cdot \cosh(x) + 1$  besitzt in  $[1e9, 1.00000001e9]$  drei einfache Nullstellen mit fast senkrechten Tangenten. Mit `Tolerance = 1e-34` erhält man das Ergebnis:

```
[1.000000000993400903293511449821674e9,1.000000000993400903293511449821675e9]
may contain a zero (not verified unique)!
[1.000000004134993556883304688284317e9,1.000000004134993556883304688284318e9]
may contain a zero (not verified unique)!
[1.000000007276586210473097926746961e9,1.000000007276586210473097926746962e9]
may contain a zero (not verified unique)!
Calculated 3 different interval enclosure(s)
```

Mit `Tolerance = 1e-4000` erhält man die Fehlermeldung: **Speicherzugriffsfehler**  
und mit `Tolerance = 1e-6600` erhält man das Ergebnis:

```
[1.0000000009934...92174077065973028e9,1.0000000009934...92174077065973029e9]
encloses a locally unique zero!
[1.0000000041349...87813976055331239e9,1.0000000041349...87813976055331240e9]
encloses a locally unique zero!
[1.0000000072765...83453875044689451e9,1.0000000072765...83453875044689452e9]
encloses a locally unique zero!

Calculated 3 different interval enclosure(s)
```

Bei einfachen Nullstellen mit nahezu senkrechten Tangenten erhält man also bei diesem Beispiel den Beweis der Eindeutigkeit erst mit hinreichend kleinen relativen Fehlerschranken. Es ist sogar möglich, mit `Tolerance = 1e-2000000` die drei Nullstellen mit jeweils etwa 3000000 korrekten Dezimalstellen einzuschließen. Mit wachsender Stellenzahl wächst die Laufzeit jedoch erheblich.

Die nächsten drei Beispielprogramme verlangen innerhalb der Funktion `void compute(...)` die folgenden Eingaben:

- n-th derivative of the function, n = ?*, d.h. die Ordnung  $n \geq 0$  der Ableitung von  $f(x)$ , deren Nullstellen zu berechnen sind. Die Funktion  $f(x)$  selbst ist jeweils vor der Funktion `void compute(...)` zu definieren.
- Start intervall [x1,x2] = ?*, d.h. das Intervall in der Form  $[1.234e-500, 1.23e-490]$ , in dem die Nullstellen einzuschließen sind.
- The precision in bits of the start interval ?* Die Eingabe von 53 ist meist völlig ausreichend. Wenn jedoch das Startintervall schon sehr schmal ist und noch weiter verbessert werden soll, so muss seine Präzision in Bits eingegeben werden, d.h. besteht z.B. die Intervallunterschranke aus 100 Dezimalstellen, so sollten wenigstens  $100 \cdot 3.321928095 \approx 332$  Bits eingegeben werden. Innerhalb der Funktion `void AllZeros(...)` wird dann mit einer `CurrentPrecision` von wenigstens 332 Bits gerechnet.
- Number k of derivatives for mean value forms (e.g. 1 <= k <= 10), k = ?* Hier ist die Anzahl  $k \geq 1$  der Ableitungen anzugeben, die bei der benutzten Mittelwertform zur Anwendung kommen sollen.

5. *No inflation:  $i=0$ , inflation to left:  $i=1$ , inflation to right:  $i=2$ ;  $i = ?$*  Bei symmetrischer Lage der Nullstellen im Startintervall können einige Nullstellenkandidaten mehrfach gezählt werden. Um dies zu vermeiden, kann das Startintervall wahlweise nach rechts oder links minimal aufgebläht werden. Dadurch möglicherweise eingefangene Nullstellenkandidaten werden jedoch anschließend wieder aussortiert.
6. *(Relative) Tolerance for the enclosures (e.g.  $1e-20$ ) = ?* Sollen die Nullstellen mit z.B. 20 korrekten Dezimalstellen eingeschlossen werden, so sollte man hier `1e-20` eingeben. Innerhalb von `void AllZeros(...)` wird dann die `CurrentPrecision` auf mindestens  $-1.5 \cdot \text{expo}(\text{Tolerance}) + 53$  gesetzt. Sollte dieser Wert immer noch zu klein sein, so kann bei der nächsten Eingabeaufforderung ein noch größerer Wert in Bits eingegeben werden.
7. *How many decimals should have the output format = ?* Hier kann die Anzahl der Dezimalstellen im Ausgabeformat festgelegt werden.

Das folgende Programm MPFR-21 berechnet für ein gegebenes Startintervall  $[x]$  alle Nullstellen der  $n$ -ten Ableitung  $f^{(n)}(x)$  der recht komplizierten Funktion

$$f(x) := \cosh(x-1) - \operatorname{erf}((x+3)/(3+\cos(9x))) + \exp(-\operatorname{sqr}(x)) \cdot \sin(5.5 \cdot x - 2), \quad x \in [x].$$

Versuchen Sie (lieber nicht!) etwa  $f^{(5)}(x)$  mit Papier und Bleistift zu berechnen! Löst man diese Aufgabe mit einem Algebrasystem, so wird man von seitenlangen Funktionstermen erschlagen, und es ist völlig ausgeschlossen, für die Nullstellenbestimmung den optimalen Funktionsterm zu bestimmen. Aber mit `AllZeros(...)` und mithilfe der Mittelwertform (mea-value form,  $k=8$ ) lassen sich im Intervall  $[-14, 14]$  alle 823 einfache Nullstellen erfolgreich einschließen!

```

1 //=====
2 //  MPFR-21.cpp
3 //  Calculating enclosures of all zeros of the n-th derivative of the
4 //  basic expression
5 //=====
6
7 #include <iostream>
8 #include "mpficlass.hpp"
9
10 using namespace cxsc;
11 using namespace std;
12 using namespace MPFR;
13 using namespace MPFI;
14
15 namespace{ int n; } // n-th derivative of the basic expression
16                 // is used in function definition of f
17 namespace{ int k; } // derivatives up to this order are used
18                 // by the centered forms
19 class BasicExpression
20 {
21 public:
22   MPitaylor operator()(const MPitaylor& x)
23   { // Evaluation of expression:
24     MPitaylor r = cosh(x-1)-erf((x+3)/(3+cos(9*x))) + exp(-sqr(x))*sin(5.5*x-2);
25     // MPitaylor r = 1-cos((x-0.992)*(x-0.993)); // Test for many candidates
26     // MpfrClass t(0.5);
27     // MPitaylor r = sqr(sin(x)) - comp(t,-19999);
28     return r;
29   } //operator
30
31 } basic_f; // class BasicExpression
32           // with one global object basic_f (used in function f)
33 MPDerivType f (const MPDerivType& xx)
34 {
35   MPitaylor x(:,n+:,k, fValue(xx)); // independent variable x

```

```

36 MPitaylor xmid(n+k, mid(fValue(xx))); // independent variable xmid
37 MPitaylor rinf;
38 MPitaylor rsup;
39
40 MPitaylor rx = basic_f(x); // function evaluated with interval argument
41
42 MprClass::SetCurrPrecision(MprClass::GetCurrPrecision()+10);
43 MPitaylor rmidx = basic_f(xmid); // function evaluated with
44 // point interval argument
45 MprClass::SetCurrPrecision(MprClass::GetCurrPrecision()-10);
46
47 MpfiClass y = fValue(xx) - mid(fValue(xx)); // used in mean-value forms
48
49 MpfiClass g[k+1]; // g[0]= f, g[1]= df, g[2]= ddf, ...
50
51 g[k] = get_j_derivative(rx,n+k); // highest derivative
52 for(int i = k-1; i>=0; i--)
53     g[i] = get_j_derivative(rmidx,n+i) + g[i+1]*y; // central form
54
55 MpfiClass dummy;
56 return MPDerivType(g[0], g[1], dummy);
57 } // MPDerivType f
58
59 //-----
60 // Function for prompting and reading information to call the function
61 // 'AllZeros()'. This function must be called with the function 'f'
62 // and a string 'Name' containing a textual description of that function.
63 //-----
64
65 void compute ( MPddf_FctPtr f, const char *Name )
66 {
67     MpfiClass    SearchInterval;
68     MprClass     Tolerance,
69                 SupStart, InfStart, EpsBlow;
70     MPivector    Zero, ZeroV;
71     intvector    Unique, UNIQUE;
72     int          NumberOfZeros, NOZ, i, Error;
73     long int     K;
74     PrecisionType prec;
75     string str;
76
77     cout << endl << " *****" << endl;
78     cout << " ***" << endl;
79     cout << " *** Computing all zeros of the "
80             << "n-th derivative of the function " << Name << endl;
81     cout << " *** n = ? "; cin >> n;
82     cout << " *** The zeros of the " << n << "-th derivative of the function "
83             << endl;
84     cout << " *** " << Name << " are to be calculated." << endl;
85     cout << " *** Start intervall [x1,x2] = ? "; cin >> str;
86     cout << " *** The precision in bits of the start interval must be entered."
87             << endl;
88     cout << " *** (e.g. 10 decimals require a precision of 10*3.3219280 bits) ? ";
89     cin >> prec;
90     if (prec<53) prec = 53;
91     MpfiClass Start(str,prec); // All zeros in Start are to be enclosed.
92     set_Mpfi(SearchInterval, Start, Start.GetPrecision());
93     cout << " *** Number k of derivatives for mean-value forms (e.g. 1<=k<=10),
94             if (k<1) k=1;
95     cout << " *** k for mean value forms: " << ::k << endl;
96     cout << " *** To avoid a possible symmetric location of the zeros in " << endl;
97     cout << " *** the start interval a small inflation can be realized. " << endl;
98
99     do {

```

```

100 cout << "   *** No inflation: i=0, inflation to left: i=1,
101           inflation to right: i=2; i = ? "; cin >> i;
102 } while (i<0 || i>3);
103
104 if (i>0)
105 {
106     K = -0.86*SearchInterval.GetPrecision();
107     EpsBlow = 1;
108     times2pown(EpsBlow, K);
109     MpfiClass::SetCurrPrecision(SearchInterval.GetPrecision());
110     if (i==1)
111         SearchInterval = Blow_L(SearchInterval, EpsBlow);
112     else // i=2
113         SearchInterval = Blow_R(SearchInterval, EpsBlow);
114 }
115 if (i==0)
116     cout << "   *** SearchInterval = " << SearchInterval << endl;
117 else cout << "   *** Blown up SearchInterval = " << SearchInterval << endl;
118 cout << "   *** (Relative) Tolerance for the enclosures (e.g. 1e-20) = ? ";
119 cin >> str;
120 Tolerance = string2Mpfr(str);
121 cout << "   *** Tolerance = " << Tolerance << endl;
122 cout << "   *** With Tolerance the CurrentPrecision is calculated in" << endl;
123 cout << "   *** AllZeros(...) by: -1.5*expo(Tolerance)+53 = "
124         << -1.5*expo(Tolerance)+53 << endl;
125 cout << "   *** In some cases this value can be too small.
126           So a greater value in bits" << endl;
127 cout << "   *** for the CurrentPrecision can be necessary: (e.g. 300) = ? ";
128 cin >> prec;
129 if (prec < 53) prec = 53;
130 MpfrClass::SetCurrPrecision(prec);
131 cout << "   *** Minimum CurrentPrecision in bits: "
132         << MpfrClass::GetCurrPrecision() << endl;
133 prec = -expo(Tolerance)/3.321928095 + 1;
134 cout << "   *** With respect to Tolerance each of the zero enclosures " << endl;
135 cout << "   *** should have " << prec << " common decimals." << endl;
136 cout << "   *** How many decimals should have the output format = ? ";
137 cin >> prec;
138 if (prec < 10) prec = 10;
139 cout.precision(prec); // only for the output format
140 cout << endl << "   *** Computing all zeros of the " << n
141         << " -th derivative of " << endl;
142 cout << "   " << Name << endl;
143 cout << "   *** please wait some moments:" << endl;
144
145 AllZeros(f, SearchInterval, Tolerance, Zero, Unique, NumberOfZeros, Error);
146
147 cout << endl << "   *** All zero enclosures are calculated ... "
148         << endl << endl;
149 if (Error)
150     cout << endl << AllZerosErrMsg(Error) << endl;
151 else
152 {
153     SupStart = Sup(Start); InfStart = Inf(Start);
154     if (i==2) // Inflation to right
155         // In Zero and Unique those zero enclosures are sorted out
156         // with Inf(Zero[k])>Sup(Start).
157         { i = 0;
158           for (int k=1; k<=NumberOfZeros; k++)
159               {
160                 if (Inf(Zero[k-i])<=SupStart)
161                     {
162                         Zero[k-i] = Zero[k];
163                         Unique[k-i] = Unique[k];

```

```

164     }
165     else i++;
166 }
167 }
168 else
169 {
170     if (i==1) // Inf(SearchInterval) < Inf(Start)
171         // In Zero and Unique those zero enclosures are sorted out
172         // with Sup(Zero[k])<Inf(Start).
173         { i = 0;
174         for (int k=1; k<=NumberOfZeros; k++)
175             {
176                 if (Sup(Zero[k-i])>=InfStart)
177                     {
178                         Zero[k-i] = Zero[k];
179                         Unique[k-i] = Unique[k];
180                     }
181                 else i++;
182             }
183     }
184 }
185 NumberOfZeros = NumberOfZeros - i;
186 // Candidates Zero[k] with Unique[k] = 0 are sorted out if it is guaranted
187 // that no zero is enclosed. However, it is possible that the remaining
188 // candidates, with Unique[k] = 0, enclose NO zero!
189
190 i=0;
191 for (int k=1; k<=NumberOfZeros; k++)
192 {
193     if (Unique[k] == 0)
194     {
195         AllZeros(f, Zero[k], Tolerance, ZeroV, UNIQUE, NOZ, Error);
196         if (NOZ==0) // No zero found
197             i++; // Elimination of Zero[k] including no zero
198         else
199             {
200                 Zero[k-i] = Zero[k];
201                 Unique[k-i] = Unique[k];
202             }
203     }
204     else
205     {
206         Zero[k-i] = Zero[k];
207         Unique[k-i] = Unique[k];
208     }
209 }
210 NumberOfZeros = NumberOfZeros - i;
211 // Now successive and overlapping enclosures are tested and
212 // under certain conditions one of these are sorted out:
213 i=0;
214 MpfClass tmp;
215 for (int k=1; k<=NumberOfZeros-1; k++)
216 {
217     if (!Disjoint(Zero[k-i], Zero[k-i+1])) // intersection is not empty
218     {
219         tmp = Zero[k-i] | Zero[k-i+1];
220         if (Unique[k-i]==Unique[k+1-i] && Unique[k-i]==0)
221             {
222                 AllZeros(f, tmp, Tolerance, ZeroV, UNIQUE, NOZ, Error);
223                 if (NOZ==1 && UNIQUE[1]==1)
224                     {
225                         Zero[k-i] = ZeroV[1]; Unique[k-i] = 1;
226                     }
227                 else

```

```

228     Zero[k-i] = tmp;
229     MPivec_del(Zero, k-i+1);
230     intvec_del(Unique, k-i+1);
231     NumberOfZeros--;
232     i++;
233 }
234 else
235 {
236     AllZeros(f, tmp, Tolerance, ZeroV, UNIQUE, NOZ, Error);
237     if (NOZ==1 && UNIQUE[1]==1)
238     {
239         Zero[k-i] = ZeroV[1];
240         Unique[k-i] = 1;
241         MPivec_del(Zero, k-i+1);
242         intvec_del(Unique, k-i+1); NumberOfZeros--;
243         i++;
244     }
245 }
246 } // !Disjoint( )
247 } // for (...)
248 cout << " *****" << endl << endl;
249 for (int i=1; i<=NumberOfZeros; i++)
250 {
251     cout << Zero[i] << endl;
252     if (Unique[i])
253         cout << "encloses a locally unique zero!" << endl;
254     else
255         cout << "may contain a zero (not verified unique)!" << endl;
256 }
257 cout << endl << "Calculated " << NumberOfZeros << " interval enclosure(s)"
258     << endl;
259 bool empty(true);
260 for (int i=1; i<=NumberOfZeros; i++)
261 {
262     for (int j=i+1; j<=NumberOfZeros; j++) // pairwise empty intersections?
263         empty = empty && Disjoint(Zero[i], Zero[j]);
264 }
265 cout << boolalpha << "Different enclosures do NOT overlap: "
266     << empty << endl;
267 bool allUnique(true);
268 for (int i=1; i<=NumberOfZeros; i++)
269     allUnique = allUnique && Unique[i];
270 cout << "Every enclosure contains a locally unique zero: "
271     << allUnique << endl;
272 cout.precision(Start.GetPrecision()/3.321928095);
273 if (empty && allUnique)
274 {
275     cout << "There is/are exactly " << NumberOfZeros
276         << " zero(s) in the interval " << endl;
277     cout << " " << Start << endl;
278 }
279 cout << "The zeros of the *** " << n
280     << "-th derivative *** of the basic expression have been computed."
281     << endl << endl;
282 }
283 return;
284 } // compute
285
286 int main(void) {
287     compute(f, "cosh(x-1)-erf((x+3)/(3+cos(9*x)) + exp(-sqr(x))*sin(5.5*x-2))");
288     return 0;
289 }

```



Mit den folgenden Eingaben

$$(7.3) \quad n = 20, [-14, 14], 53, k = 8, i = 0, 1e - 20, 152, 20;$$

erhält man die Ausgabe, wobei hier aus Platzgründen nur die letzte Einschließung angegeben wird:

```
*** Computing all zeros of the 20-th derivative of
    cosh(x-1)-erf((x+3)/(3+cos(9*x))) + exp(-sqr(x))*sin(5.5*x-2))
*** please wait some moments:
.
.
.
[1.39863986165066814584144317220e1,1.39863986165066814584144317221e1]
encloses a locally unique zero!

Calculated 823 interval enclosure(s)
Different enclosures do NOT overlap: true
Every enclosure contains a locally unique zero: true
There is/are exactly 823 zero(s) in the interval
[-1.4000000000000000e1,1.4000000000000000e1]
The zeros of the *** 20-th derivative ***
of the basic expression have been computed.
```

Die 823 Nullstelleneinschließungen werden intern mit etwa 45 korrekten Dezimalstellen berechnet und oben mit 30 Dezimalstellen ausgegeben. Will man die letzte Nullstelle mit z.B. 5000 Dezimalstellen berechnen, so benutzt man die obige Einschließung und startet das Programm MPFR-21 mit den Eingaben

$$n = 20, [1.39863986165066814584144317220e1, 1.39863986165066814584144317221e1], \\ 100, k = 8, i = 0, 1e - 5000, 24967, 5000;$$

und erhält so die Einschließung der letzten Nullstelle mit 5000 korrekten Dezimalziffern:

```
[1.398639861650668...818636349627e1,1.398639861650668...818636349628e1]
encloses a locally unique zero!
```

In einer `for`-Schleife kann man mit Hilfe des Vektors `Zero`, der die 823 Einschließungen mit 30 korrekten Dezimalstellen enthält, die Funktion `AllZeros(...)` für jeden Kandidaten `Zero[k]` mit den obigen Eingabedaten nochmals aufrufen und damit alle 823 Einschließungen mit 5000 korrekten Dezimalziffern berechnen. Ganz analog können die 823 Einschließungen dann auch mit einer Genauigkeit von z.B. 20 000 Dezimalstellen eingeschlossen werden.

Es sei noch angemerkt, dass auch alle 1045 Nullstellen von  $f^{(25)}$  mit den Eingaben

$$n = 25, [-14, 14], 53, k = 8, i = 0, 1e - 20, 152, 20;$$

erfolgreich eingeschlossen werden.

Für die Funktion  $f(x) := 1 - \cos((x - 0.992) \cdot (x - 0.993))$  erhält man mit den Eingabedaten

$$n = 0, [0.75, 1], 53, k = 1, i = 0, 1e - 20, 152, 20;$$

die folgenden Einschließungen der beiden doppelten Nullstellen:

```
[9.919999999999999288e-1,9.919999999999999291e-1]
may contain a zero (not verified unique)!
[9.929999999999999377e-1,9.929999999999999379e-1]
may contain a zero (not verified unique)!
```

Wählt man statt  $1e-20$  den `Tolerance`-Wert  $1e-100$ , so erhält man optimal nur die zwei Nullstellenkandidaten, wenn die obige Eingabe 152 durch 1000 ersetzt wird. Bei zu kleinen Werten erhält man zusätzliche Nullstellenkandidaten mit *may contain a zero (not verified unique)!*. Das Beispiel zeigt, wie nützlich es ist, dass man in Zeile 128, nach der Eingabe von `Tolerance`, den Mindestwert der `CurrentPrecision` in `AllZeros(...)` weiter vergrößern kann.

Die Funktion  $f(x) := \sin^2(x) - 2^{-20000}$ , vgl. Seite 163, besitzt in [3.1,3.2] zwei einfache aber sehr dicht benachbarte Nullstellen mit fast waagerechter Tangente, die auch jetzt z.B. mit den Eingabedaten

$$n = 0, [3.1415\dots3751, 3.1415\dots3752], 180, k = 1, i = 0, 1e - 25000, 124650, 50;$$

wegen der Fehlermeldung *Speicherzugriffsfehler* nicht getrennt werden können. Mit den Eingabedaten

$$n = 0, [3.1, 3.2], 53, k = 1, i = 0, 1e - 30, 202, 30;$$

erhält man lediglich die Einschließung

[3.14159265358979323846264338327, 3.14159265358979323846264338329]  
**may contain a zero (not verified unique)!**

Das prinzipiell einfache Beispiel zeigt, dass auch die Funktion `AllZeros(...)` nicht alle Nullstellenprobleme lösen kann, aber vorhandene Nullstellen werden grundsätzlich nicht übersehen, wenn das Programm ohne Fehlermeldung beendet wird. Das obige spezielle Problem lässt sich natürlich sehr einfach lösen, wenn  $f(x)$  als Produkt  $f(x) = (\sin(x) + 2^{-10000}) \cdot (\sin(x) - 2^{-10000})$  geschrieben wird und anschließend die einfache Nullstelle jedes einzelnen Faktors separat berechnet wird. Auch hier ist die Wahl des günstigsten Funktionsterms wieder entscheidend.

Mit dem folgenden Programm `MPFR-22.cpp` werden die Nullstellen von  $w_{20}^{(n)}(x)$  berechnet für  $0 \leq n \leq 20$ , wobei das Wilkinson-Polynom  $w_{20}(x)$  definiert ist durch [66]

$$w_{20}(x) := (x - 1) \cdot (x - 2) \cdot \dots \cdot (x - 20) = \sum_{i=0}^{20} a_i \cdot x^i, \quad a_{20} = 1, \quad a_{19} = -210, \dots$$

Da die Funktionswerte  $w_{20}(x)$  selbst und damit auch Lage der Nullstellen  $x_i = i$ ,  $i = 1, 2, \dots, 20$ , sehr empfindlich ist bezüglich nur kleiner Änderungen der ganzzahligen Koeffizienten  $a_i$ , muss zunächst sichergestellt werden, dass die  $a_i$  im Rechner exakt gespeichert werden. Dies wird im folgenden Quellcode realisiert von Zeile 34 bis 65. Der Quellcode ab Zeile 80 stimmt mit dem Quellcode von `MPFR-21.cpp`, dort ab Zeile 33, genau überein und wird hier aus Platzgründen weggelassen.

```

1 //=====
2 //  MPFR-22.cpp
3 //  Calculating enclosures of all zeros of the n-th derivative of the
4 //  basic expression
5 //  Wilkinson polynomial of degree 20
6 //=====
7
8 #include <iostream>
9 #include "mpficlass.hpp"
10
11 using namespace cxsc;
12 using namespace std;
13 using namespace MPFR;
14 using namespace MPFI;
15
16 namespace{ int n; } // n-th derivative of the basic expression
17                  // is used in function definition of f
18 namespace{ int k; } // derivatives up to this order are used
19
20 class BasicExpression {
21 #define MP MpfiClass
22 public:
23     static const int degree = 20;
24     MpfiClass a[degree+1]; // coeff.: a[degree] ... a[0]
25

```

```

26 public:
27 BasicExpression() // constructor
28 {
29     PrecisionType oldPrec(MP::GetCurrPrecision());
30     MP::SetCurrPrecision(57);
31     for (int i=degree; i>=0; i--) a[i] = 0;
32
33     // Wilkinson polynomial of degree 20
34     // Initialization of the coefficients:
35     a[20] = MP( 1);
36     a[19] = MP(-210);
37     a[18] = MP( 20615);
38     a[17] = MP(-1256850.0);
39     a[16] = MP( 53327946.0);
40     a[15] = MP(-1672280820.0);
41     a[14] = MP( 40171771630.0);
42     a[13] = MP(" -756111184500.0 ");
43     a[12] = MP(" 11310276995381.0 ");
44     a[11] = MP(" -135585182899530.0 ");
45     a[10] = MP(" 1307535010540395.0 ");
46     a[9] = MP(" -10142299865511450.0 ");
47     a[8] = MP(" 63030812099294896.0 ");
48     a[7] = MP(" -311333643161390640.0 ");
49     a[6] = MP(" 1206647803780373360.0 ");
50     a[5] = MP(" -3599979517947607200.0 ");
51     a[4] = MP(" 8037811822645051776.0 ");
52     a[3] = MP(" -12870931245150988800.0 ");
53     a[2] = MP(" 13803759753640704000.0 ");
54     a[1] = MP(" -8752948036761600000.0 ");
55     a[0] = MP(" 2432902008176640000.0 ");
56
57     for (int i=degree; i>=0; i--) // check whether coeff are exactly repr.
58     {
59         if (diam(a[i])!=0)
60         {
61             cout << " *** coeff a[" << i << "] not representable ,
62              increase the CurrenPrecision!" << endl;
63             exit(1);
64         }
65     }
66     MP::SetCurrPrecision(oldPrec); // restore precision setting
67     return;
68 } // constructor
69
70 MPitaylor operator()(const MPitaylor& x)
71 {
72     // Compute Taylor coefficients of basic expression r(x) via Horner's scheme
73     MPitaylor r = a[degree]*x + a[degree-1]; //expanded form of polynomial is used
74     for (int i=degree-2; i>=0; i--)
75         r = r*x+a[i];
76     return r;
77 } // operator
78
79 } basic_f; // class with one global object basic_f (used in function f)

```

Mit den Eingabedaten

$$n = 0, [1, 20], 53, k = 4, i = 0, 1e - 30, 202, 30;$$

erhält man die Ausgabe der Nullstelleneinschließungen mit 30 korrekten Dezimalstellen:



Das Tschebyscheff-Polynom  $T_N(x)$  vom Grad  $N = 40$  ist definiert durch

$$T_{40}(x) := \cos(40 \cdot \arccos(x)) = \sum_{k=0}^{20} a_{2k} \cdot x^{2k}; \quad a_0 = 1, \quad a_2 = -800, \dots, \quad a_{40} = 549755813888.$$

Mit folgendem Programm `MPFR-23.cpp` werden für die  $n$ -te Ableitung  $T_{40}^{(n)}(x)$ , mit  $0 \leq n \leq 40$ , alle Nullstellen für  $x \in [-1, 1]$  berechnet. Dazu muss im folgenden Quellcode zunächst ab Zeile 31 bis 67 dafür gesorgt werden, dass die ganzzahligen Koeffizienten  $a_{2k}$  exakt gespeichert werden. Die Abbildung 7.3 zeigt den zur  $y$ -Achse symmetrischen Graphen von  $T_{40}(x)$  für  $x \in [0, 1]$ . Man erkennt, dass sich die einfachen Nullstellen am Intervallrand zusammendrängen und die  $x$ -Achse fast senkrecht schneiden. Das numerische Problem besteht also darin, auch noch diese Nullstellen erfolgreich einzuschließen.

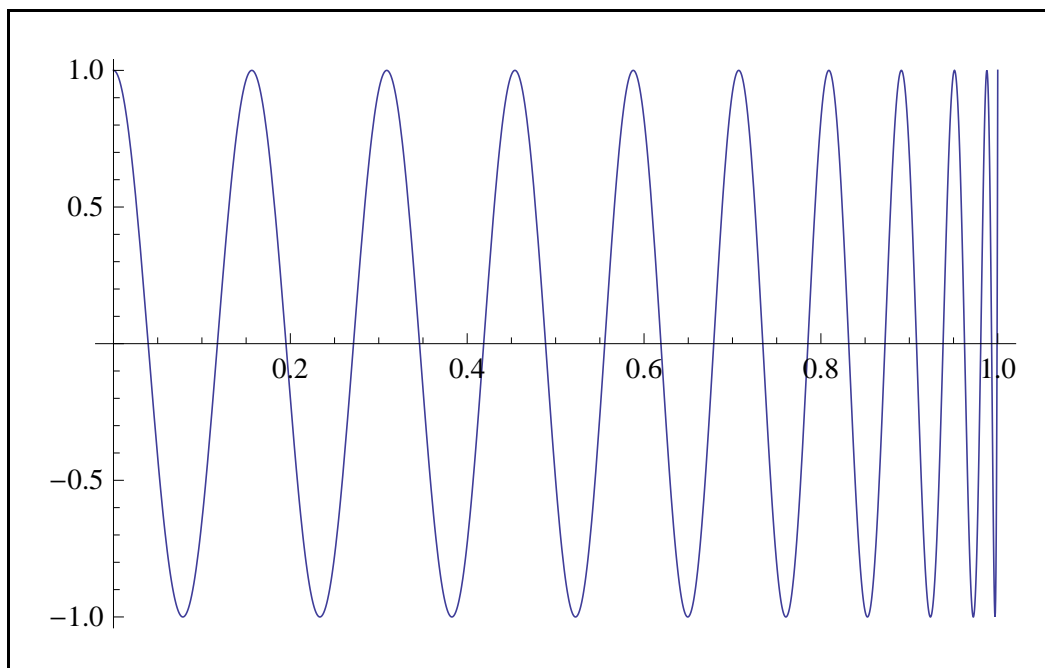


Abbildung 7.3.: Tschebyscheff-Polynom  $T_{40}(x)$ ,  $x \in [0, +1]$

```

1
2 //=====
3 //  MPFR-23.cpp
4 //  Calculating enclosures of all zeros of the n-th derivative of the
5 //  basic expression
6 //  Tschebyscheff polynomial of degree 40
7 //=====
8
9 #include <iostream>
10 #include "mpficlass.hpp"
11
12 using namespace cxsc;
13 using namespace std;
14 using namespace MPFR;
15 using namespace MPFI;
16
17 namespace{ int n; } // n-th derivative of the basic expression
18 // is used in function definition of f
19 namespace{ int k; } // derivatives up to this order are used
20
21 class BasicExpression {
22 #define MP MpfiClass

```

```

23 public:
24     static const int degree = 40;
25     MpfiClass a[degree+1]; // coeff. a[degree] ... a[0]
26
27 public:
28     BasicExpression() // constructor
29     {
30         PrecisionType oldPrec(MP::GetCurrPrecision());
31         MP::SetCurrPrecision(53);
32         for (int i=degree; i>=0; i--) a[i]= 0;
33
34         // Tschebyscheff polynomial of degree 40
35         // Initialization of the coefficients:
36         a[40]= MP(" 549755813888.0");
37         a[38]= MP(" -5497558138880.0");
38         a[36]= MP(" 25426206392320.0");
39         a[34]= MP(" -72155450572800.0");
40         a[32]= MP(" 140552804761600.0");
41         a[30]= MP(" -199183403319296.0");
42         a[28]= MP(" 212364657950720.0");
43         a[26]= MP(" -173752901959680.0");
44         a[24]= MP(" 110292369408000.0");
45         a[22]= MP(" -54553214976000.0");
46         a[20]= MP(" 21002987765760.0");
47         a[18]= MP(" -6254808268800.0");
48         a[16]= MP(" 1424085811200.0");
49         a[14]= MP(" -243433472000.0");
50         a[12]= MP(" 30429184000.0");
51         a[10]= MP(" -2677768192.0");
52         a[8]= MP(" 156900480.0");
53         a[6]= MP(" -5617920.0");
54         a[4]= MP(" 106400.0");
55         a[2]= MP(" -800.0");
56         a[0]= MP(" 1.0");
57
58         for (int i=degree; i>=0; i--) // check whether coeff are exactly repr.
59         {
60             if (diam(a[i])!=0)
61             {
62                 cout << "   *** coeff a[" << i << "] not representable ,
63                     increase the CurrenPrecision!"
64                     << endl;
65                 exit(1);
66             }
67         }
68         MP::SetCurrPrecision(oldPrec); // restore precision setting
69         return;
70     } //constructor
71
72     MPitaylor operator()(const MPitaylor& x)
73     { // Compute Taylor coefficients of basic expression r(x) via Horner's scheme
74         MPitaylor r = a[degree]*x + a[degree-1]; // expanded form of polynomial
75         for (int i=degree-2; i>=0; i--)
76             r = r*x+a[i];
77         // MPitaylor r = cos(degree*acos(x)); // polynomial via cos and acos
78         return r;
79     } //operator
80
81 } basic_f; // class with one global object basic_f (used in function f)

```

Der Rest des Programmcodes stimmt überein mit dem Code des Programms MPFR-21.cpp auf Seite 164 und wird hier aus Platzgründen nicht angegeben.

Mit den Eingabedaten

$$n = 0, [-1, +1], 53, k = 5, i = 0, 1e - 30, 202, 30;$$

erhält man die Einschließungen der folgenden 40 einfachen Nullstellen:

**[-9.99229036240722934737126260342e-1, -9.99229036240722934737126260341e-1]  
encloses a locally unique zero!**

**[-9.93068456954926295637437247811e-1, -9.93068456954926295637437247810e-1]  
encloses a locally unique zero!**

**[-9.80785280403230449126182236135e-1, -9.80785280403230449126182236134e-1]  
encloses a locally unique zero!**

**[-9.62455236453647287630266405186e-1, -9.62455236453647287630266405185e-1]  
encloses a locally unique zero!**

**[-9.38191335922484134452339726687e-1, -9.38191335922484134452339726686e-1]  
encloses a locally unique zero!**

**[-9.08143173825081299258085836572e-1, -9.08143173825081299258085836571e-1]  
encloses a locally unique zero!**

**[-8.72496007072797114525161099223e-1, -8.72496007072797114525161099222e-1]  
encloses a locally unique zero!**

**[-8.31469612302545237078788377618e-1, -8.31469612302545237078788377617e-1]  
encloses a locally unique zero!**

**[-7.85316930880744927470340278948e-1, -7.85316930880744927470340278947e-1]  
encloses a locally unique zero!**

**[-7.34322509435685535636126222188e-1, -7.34322509435685535636126222186e-1]  
encloses a locally unique zero!**

**[-6.78800745532941741393855554242e-1, -6.78800745532941741393855554241e-1]  
encloses a locally unique zero!**

**[-6.19093949309833986941560856247e-1, -6.19093949309833986941560856246e-1]  
encloses a locally unique zero!**

**[-5.55570233019602224742830813949e-1, -5.55570233019602224742830813948e-1]  
encloses a locally unique zero!**

**[-4.88621241496954947420190887839e-1, -4.88621241496954947420190887838e-1]  
encloses a locally unique zero!**

**[-4.18659737537428086675565205122e-1, -4.18659737537428086675565205121e-1]  
encloses a locally unique zero!**

**[-3.46117057077492976468214994929e-1, -3.46117057077492976468214994928e-1]  
encloses a locally unique zero!**

**[-2.71440449865074253343787401296e-1, -2.71440449865074253343787401295e-1]  
encloses a locally unique zero!**

**[-1.95090322016128267848284868478e-1, -1.95090322016128267848284868477e-1]  
encloses a locally unique zero!**

**[-1.17537397457837644105568266841e-1, -1.17537397457837644105568266840e-1]  
encloses a locally unique zero!**

**[-3.92598157590686090208033637984e-2, -3.92598157590686090208033637983e-2]  
encloses a locally unique zero!**

**[3.92598157590686090208033637983e-2, 3.92598157590686090208033637984e-2]  
encloses a locally unique zero!**

**[1.17537397457837644105568266840e-1, 1.17537397457837644105568266841e-1]  
encloses a locally unique zero!**

**[1.95090322016128267848284868477e-1, 1.95090322016128267848284868478e-1]  
encloses a locally unique zero!**

**[2.71440449865074253343787401295e-1, 2.71440449865074253343787401296e-1]  
encloses a locally unique zero!**

**[3.46117057077492976468214994928e-1, 3.46117057077492976468214994929e-1]  
encloses a locally unique zero!**

**[4.18659737537428086675565205121e-1, 4.18659737537428086675565205122e-1]  
encloses a locally unique zero!**

```

[4.88621241496954947420190887838e-1,4.88621241496954947420190887839e-1]
encloses a locally unique zero!
[5.55570233019602224742830813948e-1,5.55570233019602224742830813949e-1]
encloses a locally unique zero!
[6.19093949309833986941560856246e-1,6.19093949309833986941560856247e-1]
encloses a locally unique zero!
[6.7880074553294174139385554241e-1,6.7880074553294174139385554242e-1]
encloses a locally unique zero!
[7.34322509435685535636126222186e-1,7.34322509435685535636126222188e-1]
encloses a locally unique zero!
[7.85316930880744927470340278947e-1,7.85316930880744927470340278948e-1]
encloses a locally unique zero!
[8.31469612302545237078788377617e-1,8.31469612302545237078788377618e-1]
encloses a locally unique zero!
[8.72496007072797114525161099222e-1,8.72496007072797114525161099223e-1]
encloses a locally unique zero!
[9.08143173825081299258085836571e-1,9.08143173825081299258085836572e-1]
encloses a locally unique zero!
[9.38191335922484134452339726686e-1,9.38191335922484134452339726687e-1]
encloses a locally unique zero!
[9.62455236453647287630266405185e-1,9.62455236453647287630266405186e-1]
encloses a locally unique zero!
[9.80785280403230449126182236134e-1,9.80785280403230449126182236135e-1]
encloses a locally unique zero!
[9.93068456954926295637437247810e-1,9.93068456954926295637437247811e-1]
encloses a locally unique zero!
[9.99229036240722934737126260341e-1,9.99229036240722934737126260342e-1]
encloses a locally unique zero!

```

```

Calculated 40 interval enclosure(s)
Different enclosures do NOT overlap: true
Every enclosure contains a locally unique zero: true
There is/are exactly 40 zero(s) in the interval
[-1.000000000000000,1.000000000000000]
The zeros of the *** 0-th derivative *** of the basic expression
have been computed.

```

Mit den folgenden Eingabedaten

$$n = 0, [-1, +1], 53, k = 9, i = 0, 1e - 300, 1547, 30;$$

werden die 40 einfachen Nullstellen mit einer Genauigkeit von 300 Dezimalstellen eingeschlossen. Mit noch höherer relativer Genauigkeit, z.B.  $1e-500$ , können die Einschließungen jedoch nicht mehr berechnet werden. Will man z.B. die letzte Nullstelle mit noch größerer Genauigkeit, z.B. mit 10000 Dezimalstellen, einschließen, so berechnet man alle Einschließung in  $[-1, 1]$  zunächst mit **300** Dezimalstellen und wählt dann die Eingabedaten für die letzte Nullstelle:

$$n = 0, [9.99229036...2009244248e - 1, 9.99229036...2009244248e - 1], \\ 53, k = 9, i = 0, 1e - 10000, 50000, 30;$$

Nach dem neuen Programmstart erhält man dann die gewünschte hohe Genauigkeit. Dies lässt sich (in einer **for**-Schleife) mit jeder anderen Nullstelle analog durchführen. Natürlich wird man bei diesem Beispiel auch wieder an die Grenzen der Funktion `AllZeros(...)` stoßen, wenn nur der Polynomgrad hinreichend erhöht wird. Es bleibt dem Anwender überlassen zu testen, bis zu welcher Genauigkeit die Nullstellen eingeschlossen werden können, wenn man an Stelle der Polynomdarstellung die alternative Definition  $T_{40}(x) := \cos(40 \cdot \arccos(x))$  wählt. Zum Testen der Leistungsfähigkeit von `AllZeros()` sollte natürlich immer die Polynomdarstellung gewählt werden.



Abschließend noch einige Empfehlungen zur Anwendung der Funktion `AllZeros(...)`:

Bei einem einfachen Problem verwende man zunächst das Programm MPFR-15. Wenn dies scheitert, verwende man eines der Programme MPFR-21 oder MPFR-22, wobei  $k \geq 1$  schrittweise um eins zu vergrößern ist. Achten Sie darauf, nach Möglichkeit den günstigsten Funktionsterm zu wählen.

Scheitert auch dies, so ist das Startintervall so zu unterteilen, dass die Nullstellensuche in jedem Teilintervall erfolgreich ist, dabei sollte die relative Genauigkeit (`Tolerance`) mit etwa  $1e-20$  zunächst nicht zu groß gewählt werden. Nullstelleneinschließungen mit: *may contain a zero (not verified unique)*! sollten in einem einzelnen Intervall zusammengefasst, und nach Seite 169 unten sollte der Mindestwert der `CurrenPrecision` hinreichend vergrößert werden.

Falls die geforderte Genauigkeit einer Einschließung nach der ersten Anwendung von `AllZeros(...)` noch nicht ausreicht, muss mit dieser Einschließung als neues Startintervall das entsprechende Programm mit der gewünschten größeren Genauigkeit nochmals gestartet werden, vgl. Seite 176.

Mit diesen Empfehlungen können die gängigsten Nullstellenprobleme gelöst werden. Wenn ein Programm ohne Fehlermeldung beendet wird, so werden im Vektor `Zero` garantiert alle existierenden Nullstellen aus dem Startintervall erfasst. Eine Einschließung mit: *may contain a zero (not verified unique)*! kann entweder keine oder eine oder auch mehrere einfache oder mehrfache Nullstellen enthalten.

Bei einer Einschließung mit: *encloses a locally unique zero!* ist die Einschließung selbst mathematisch garantiert und sichergestellt, dass die Einschließung nur eine einzige Nullstelle enthält!



# A. Neue Funktionen vom Typ MpfrClass

Mit Hilfe der bereits in den beiden MPFR- und MPFI-Bibliotheken implementierten Funktionen aus den Tabellen auf Seite 37 und 44 sollen noch weitere (Hilfs-)Funktionen (insbesondere hilfreich im Zusammenhang mit der Realisierung weiterer komplexer Intervallfunktionen) für den Datentyp `MpfrClass` realisiert werden, wobei darauf zu achten ist, dass mit einem zusätzlichen Rundungsparameter `rnd` der Funktionswert, abweichend vom Current-Rundungsmodus, gerundet werden kann. Damit stellt sich die Aufgabe, einen gegebenen Funktionsterm so zu programmieren, dass dieser entweder zur nächsten Rasterzahl oder garantiert auf- bzw. abgerundet werden kann. Ist eine solche Funktion dann z.B. in einem Maschinenintervall  $[a, b]$  monoton wachsend, so ist eine Einschließung aller Funktionswerte  $f(x)$ ,  $x \in [a, b]$  durch das Intervall  $[f_d(a), f_u(b)]$  gegeben, wobei  $f_d(a)$  und  $f_u(b)$  die ab- bzw. aufgerundeten Funktionswerte bedeuten.

## A.1. Grundregeln für garantierte Rundungen

Soll ein neu implementierter Ausdruck  $A$  z.B. zur nächsten Rasterzahl gerundet werden, so ist der Rundungsmodus `rnd` auf `RoundNearest` zu setzen und  $A$  ist auszuwerten. Ist dann  $\tilde{A}$  das Ergebnis dieser Auswertung, so wird  $\tilde{A}$  i.a. nicht der Vorgänger oder Nachfolger von  $A$  sein. Entsprechende Aussagen gelten auch für die auf- oder abgerundeten Maschinenergebnisse. Wir fassen zusammen:

Wird eine neu implementierte Funktion mit `rnd = RoundNearest` mit der Maschinenzahl  $x_0$  ausgewertet, so muss der Maschinenwert  $\tilde{f}(x_0)$  nicht die zum exakten Funktionswert  $f(x_0)$  nächstgelegene Rasterzahl sein. Ist  $f_u(x_0)$  der mit `rnd = RoundUp` aufgerundete Funktionswert, so wird i.a. auch  $f_u(x_0)$  nicht der Nachfolger von  $f(x_0)$  sein;  $f(x_0) \leq f_u(x_0)$  ist aber stets garantiert. Für den abgerundeten Funktionswert  $f_d(x_0)$  gelten entsprechende Aussagen.

### Bezeichnungen:

$A, B$  sind exakte Ausdrücke, die im Zahlenformat nicht exakt darstellbar sein müssen.

- $A_d$  bezeichnet den abgerundeten Maschinenwert:  $A_d \leq A$ .
- $A_u$  bezeichnet den aufgerundeten Maschinenwert:  $A_u \geq A$ .
- $\oplus_d$  bezeichnet die abrundende Maschinen-Addition.
- $\oplus_u$  bezeichnet die aufrundende Maschinen-Addition.
- $\ominus_d$  bezeichnet die abrundende Maschinen-Subtraktion.
- $\ominus_u$  bezeichnet die aufrundende Maschinen-Subtraktion.
- $\odot_d$  bezeichnet die abrundende Maschinen-Multiplikation.
- $\odot_u$  bezeichnet die aufrundende Maschinen-Multiplikation.

- $\odot_d$  bezeichnet die abrundende Maschinen-Division.
- $\odot_u$  bezeichnet die aufrundende Maschinen-Division.

Auch jetzt wird nicht verlangt, dass  $A_d = \text{pred}(A)$  oder  $A_u = \text{succ}(A)$  erfüllt sind, aber die Beziehungen  $A_d \leq A$  bzw.  $A_u \geq A$  werden stets garantiert.

### A.1.1. Unitäre Operatoren

Unabhängig davon, ob  $A$  positiv, negativ oder gleich Null ist, gilt

$$\begin{aligned} \text{(A.1)} \quad & +(A_d) \leq A; \\ \text{(A.2)} \quad & +(A_u) \geq A; \\ \text{(A.3)} \quad & -(A_d) \geq -A; \\ \text{(A.4)} \quad & -(A_u) \leq -A; \end{aligned}$$

### A.1.2. Addition

Unabhängig davon, ob  $A$  oder  $B$  positiv, negativ oder gleich Null ist, gilt

$$\begin{aligned} \text{(A.5)} \quad & A_d \oplus_d B_d \leq A + B; \\ \text{(A.6)} \quad & A_u \oplus_u B_u \geq A + B; \end{aligned}$$

### A.1.3. Subtraktion

Unabhängig davon, ob  $A$  oder  $B$  positiv, negativ oder gleich Null ist, gilt

$$\begin{aligned} \text{(A.7)} \quad & A_d \ominus_d B_u \leq A - B; \\ \text{(A.8)} \quad & A_u \ominus_u B_d \geq A - B; \end{aligned}$$

Die Beweise für die obigen Sätze sind trivial und bleiben dem Leser überlassen.

### A.1.4. Multiplikation

Um bei Multiplikation und Division die Bedingungen an beide Operanden für ein Auf- bzw. Abrunden möglichst übersichtlich formulieren zu können, geben wir noch zwei einschränkende Eigenschaften<sup>1</sup> von  $A$  an, die aber nur in speziellen Fällen erfüllt sein müssen:

Für Ausdrücke  $A$ , die auf dem Rechner gezielt auf- bzw. abzurunden sind, muss zusätzlich erfüllt sein:

$$\begin{aligned} \text{(A.9)} \quad & A_u > 0 \implies A_u \geq A \geq 0; \\ \text{(A.10)} \quad & A_d < 0 \implies A_d \leq A \leq 0; \end{aligned}$$

Wir formulieren jetzt noch Bedingungen für einen arithmetischen Ausdruck  $A$ , mit denen die Eigenschaften (A.9) bzw. (A.10) abgesichert werden:

Wird ein Ausdruck  $A$  nur durch eine einzige Rechenoperation erzeugt, so gilt (A.9) bzw. (A.10). Diese Eigenschaften sind auch dann garantiert, wenn  $A$  nur durch eine einzige Funktion aus den Tabellen auf Seite 37 und 44 realisiert wird.

<sup>1</sup>Genauer gesagt sind es die Eigenschaften des Algorithmus, mit dem der arithmetische Ausdruck  $A$  ausgewertet wird.

Das folgende Beispiel soll zeigen, dass z.B die Eigenschaft (A.9) nicht immer erfüllt sein muss. Mit den Rasterzahlen  $a, b, c, d$  definieren wir dazu den exakten Ausdruck  $A$  durch:

$$A := a \cdot b + c \cdot d - c \cdot d = a \cdot b,$$

wobei  $-\text{minfloat}() < a \cdot b = A < 0$  gelten soll. Die Rasterzahlen  $c, d > 0$  seien so gewählt, dass  $c \cdot d \approx 1$  im vorliegenden Zahlenraster nicht exakt darstellbar ist, so dass bei der Produktberechnung eine Rundung notwendig wird. Der aufzurundende Ausdruck  $A$  ist wie folgt auszuwerten:

$$\begin{aligned} A_u &:= (a \odot_u b) \oplus_u (c \odot_u d) \ominus_u (c \odot_d d) \\ &= (c \odot_u d) \ominus_u (c \odot_d d) > 0; \end{aligned}$$

Zunächst gilt wegen der geforderten Aufrundung  $(a \odot_u b) = 0$ , und wegen der notwendigen Rundung bei den Produktberechnungen gilt:  $0 < (c \odot_u d) - (c \odot_d d) \leq (c \odot_u d) \ominus_u (c \odot_d d) = A_u$ , so dass jetzt im Gegensatz zu (A.9) aus  $A_u > 0$  nicht  $A \geq 0$  gefolgert werden kann.

Wir kommen jetzt zur Formulierung der Operandenbedingungen, mit denen bei der Multiplikation gezielt auf- bzw. abgerundet werden kann. Bei der gerundeten Multiplikation muss man dabei unterscheiden, ob die Operanden positiv, negativ oder gleich Null sind. Wir betrachten zunächst das **Abrunden**, dabei bedeutet  $*$ , dass nur einer der beiden gerundeten Operanden verschwinden darf:

$$(A.11) \quad A_d \geq 0; \quad B_d \geq 0 \quad \Longrightarrow \quad A_d \odot_d B_d \leq A \cdot B;$$

$$(A.12) \quad * \quad A_u \geq 0; \quad B_d \leq 0 \quad \stackrel{(A.9)(A.10)}{\Longrightarrow} \quad A_u \odot_d B_d \leq A \cdot B;$$

$$(A.13) \quad * \quad A_d \leq 0; \quad B_u \geq 0 \quad \stackrel{(A.9)(A.10)}{\Longrightarrow} \quad A_d \odot_d B_u \leq A \cdot B;$$

$$(A.14) \quad A_u \leq 0; \quad B_u \leq 0 \quad \Longrightarrow \quad A_u \odot_d B_u \leq A \cdot B;$$

Bei der Multiplikation betrachten wir jetzt das **Aufrunden**, und nachfolgend bedeutet  $*$ , dass nur einer der beiden gerundeten Operanden verschwinden darf:

$$(A.15) \quad * \quad A_u \geq 0; \quad B_u \geq 0 \quad \stackrel{(A.9)}{\Longrightarrow} \quad A_u \odot_u B_u \geq A \cdot B;$$

$$(A.16) \quad A_d \geq 0; \quad B_u \leq 0 \quad \Longrightarrow \quad A_d \odot_u B_u \geq A \cdot B;$$

$$(A.17) \quad A_u \leq 0; \quad B_d \geq 0 \quad \Longrightarrow \quad A_u \odot_u B_d \geq A \cdot B;$$

$$(A.18) \quad * \quad A_d \leq 0; \quad B_d \leq 0 \quad \stackrel{(A.10)}{\Longrightarrow} \quad A_d \odot_u B_d \geq A \cdot B;$$

### A.1.5. Division

Wir betrachten zunächst das **Abrunden**, wobei  $B \neq 0$  vorausgesetzt wird:

$$(A.19) \quad A_d \geq 0; \quad B_u > 0 \quad \stackrel{(A.9)}{\Longrightarrow} \quad A_d \odot_d B_u \leq A/B;$$

$$(A.20) \quad A_d \leq 0; \quad B_d > 0 \quad \stackrel{(A.10)}{\Longrightarrow} \quad A_d \odot_d B_d \leq A/B;$$

$$(A.21) \quad A_u \geq 0; \quad B_u < 0 \quad \stackrel{(A.9)}{\Longrightarrow} \quad A_u \odot_d B_u \leq A/B;$$

$$(A.22) \quad A_u \leq 0; \quad B_d < 0 \quad \stackrel{(A.10)}{\Longrightarrow} \quad A_u \odot_d B_d \leq A/B;$$

Bei der Division betrachten wir jetzt das **Aufrunden**, wobei wieder  $B \neq 0$  vorausgesetzt wird:

$$(A.23) \quad A_u \geq 0; \quad B_d > 0 \quad \stackrel{(A.9)}{\Longrightarrow} \quad A_u \odot_u B_d \geq A/B;$$

$$(A.24) \quad A_u \leq 0; \quad B_u > 0 \quad \stackrel{(A.9)}{\Longrightarrow} \quad A_u \odot_u B_u \geq A/B;$$

$$(A.25) \quad A_d \geq 0; \quad B_d < 0 \quad \stackrel{(A.10)}{\Longrightarrow} \quad A_d \odot_u B_d \geq A/B;$$

$$(A.26) \quad A_d \leq 0; \quad B_u < 0 \quad \stackrel{(A.10)}{\Longrightarrow} \quad A_d \odot_u B_u \geq A/B;$$

## Anmerkungen:

1. Zum Verständnis wird zunächst (A.12) bewiesen. Dabei zeigt sich, dass beide Bedingungen (A.9) und (A.10) von Seite 180 auch tatsächlich benötigt werden. Zu zeigen ist also:

$$(A.27) \quad A_u \geq 0, \quad B_d \leq 0 \implies A_u \odot_d B_d \leq A \cdot B;$$

Der Fall  $A_u = B_d = 0$  ist auszuschließen, da sonst folgt:  $A \leq 0$  und  $B \geq 0$ , d.h.  $A \cdot B \leq 0$ , und dies ist wegen  $A_u \odot_d B_d = 0$  ein Widerspruch zur Behauptung.

Im Fall  $A_u = 0$  und  $B_d < 0$  folgt zunächst  $A_u \odot_d B_d = 0$  und  $A \leq 0$ . Um damit  $A \cdot B \geq 0$  zu garantieren, benötigt man  $B \leq 0$ , und dies folgt nur, wenn für die gerundete Größe  $B_d$  die Forderung (A.10) erfüllt ist.

Im Fall  $B_d = 0$  und  $A_u > 0$  folgt zunächst  $A_u \odot_d B_d = 0$  und  $B \geq 0$ . Um damit  $A \cdot B \geq 0$  zu garantieren, benötigt man  $A \geq 0$ , und dies folgt nur, wenn für die gerundete Größe  $A_u$  die Forderung (A.9) erfüllt ist.

Jetzt bleibt noch:  $A_u > 0$  und  $B_d < 0$ . Nach Definition von  $\odot_d$  gilt ganz allgemein:  $A_u \odot_d B_d \leq A_u \cdot B_d$ . Weiter ergibt sich:

$$A_u \geq A \xrightarrow{B_d < 0} A_u \cdot B_d \leq A \cdot B_d, \quad \text{und} \quad B_d \leq B \xrightarrow{A \geq 0} B_d \cdot A \leq A \cdot B. \blacksquare$$

Zum Beweis benötigen wir damit  $A \geq 0$ , und diese Bedingung wird nur erfüllt, wenn für die gerundete Größe  $A_u > 0$  die Forderung (A.9) auch tatsächlich erfüllt ist. Man kann den Beweis auch etwas abändern, muss dann aber auf (A.10) zurückgreifen. Die restlichen Beweise bez. der gerundeten Multiplikation können ganz analog durchgeführt werden.

2. Als Beispiel für die Division wird jetzt noch (A.24) bewiesen. Dabei zeigt sich, dass eine der beiden Bedingungen (A.9) und (A.10) von Seite 180 auch tatsächlich benötigt wird. Zu zeigen ist:

$$(A.28) \quad A_u \leq 0; \quad B_u > 0 \implies A_u \oslash_u B_u \geq A/B;$$

Im Fall  $A_u = 0$  und  $B_u > 0$  folgt zunächst  $A_u \oslash_u B_u = 0$  und  $A \leq 0$ . Um  $A/B \leq 0$  zu garantieren, benötigt man  $B > 0$ , was wegen  $B_u > 0$  mithilfe von (A.9) gesichert ist.

Wir betrachten jetzt den Fall  $A_u < 0, B_u > 0$ . Zunächst gilt:  $A_u \oslash_u B_u \geq A_u/B_u$ ;

Wegen  $B_u > 0$  gilt nach (A.9):  $B_u \geq B > 0 \implies 1/B_u \leq 1/B \xrightarrow{A_u < 0} A_u/B_u \geq A_u/B$ ;

Es gilt nach Voraussetzung:  $0 > A_u \geq A \xrightarrow{B > 0} A_u/B \geq A/B. \blacksquare$

Die restlichen Beweise bez. der gerundeten Division können ganz analog durchgeführt werden.

3. Es sei noch einmal darauf hingewiesen, dass im Gegensatz zur Multiplikation und Division bei der gerichteten Addition und Subtraktion nicht untersucht werden muss, ob die Operanden positiv, negativ oder gleich Null sind.
4. Ein erstes einfaches Anwendungsbeispiel für gerichtete Rundungen findet man für die neu installierte Funktion  $f(x, y) = x^2 - y^2$  auf Seite 184.

Nachdem wir für die gerichteten Rundungen bei der Multiplikation und Division die entsprechenden Bedingungen in (A.11) bis (A.18) bzw. in (A.19) bis (A.26) beschrieben haben, stellt sich jetzt noch die Frage, wie z.B. in (A.12) die Bedingungen  $A_u \geq 0$  und  $B_d \leq 0$  garantiert werden können. Dazu formulieren wir zunächst:

Ist  $f(x)$  eine Funktion der MPFR-Bibliothek und bedeutet  $f_a(x)$  den von der Null weggerundeten Maschinenwert, so gilt:

$$(A.29) \quad f_a(x) = 0 \implies f(x) = 0;$$

$$(A.30) \quad f_a(x) > 0 \implies f(x) > 0;$$

$$(A.31) \quad f_a(x) < 0 \implies f(x) < 0;$$

Mithilfe von  $f_a(x)$  erhält man daher gesicherte Aussagen über den exakten Funktionswert  $f(x)$ , und mit den folgenden Sätzen erhält man schließlich Aussagen bez. der auf- bzw. abgerundeten Funktionswerte  $f_u(x)$ ,  $f_d(x)$

Ist  $f(x)$  eine Funktion der MPFR-Bibliothek und bedeuten  $f_u(x)$  und  $f_d(x)$  die auf- bzw. abgerundeten Funktionswerte, so gilt:

$$(A.32) \quad f(x) \geq 0 \implies f_u(x), f_d(x) \geq 0;$$

$$(A.33) \quad f(x) \leq 0 \implies f_u(x), f_d(x) \leq 0;$$

Wir fassen zusammen:

Ist  $f(x)$  eine Funktion der MPFR-Bibliothek und bedeuten  $f_u(x)$  und  $f_d(x)$  die auf- bzw. abgerundeten Funktionswerte und ist  $f_a(x)$  der von der Null weggerundete Funktionswert, so gilt:

$$(A.34) \quad f_a(x) \geq 0 \implies f_u(x), f_d(x) \geq 0;$$

$$(A.35) \quad f_a(x) \leq 0 \implies f_u(x), f_d(x) \leq 0;$$

Damit erhalten wir mithilfe von  $f_a(x)$  die gewünschten Aussagen bezüglich der auf- bzw. abgerundeten Funktionswerte  $f_u(x)$ ,  $f_d(x)$ . Wir formulieren noch zusätzlich:

Ist  $f(x)$  eine Funktion der MPFR-Bibliothek und bedeuten  $f_u(x)$  und  $f_d(x)$  die auf- bzw. abgerundeten Funktionswerte, so gilt:

$$(A.36) \quad f_u(x) > 0 \iff f(x) > 0, \quad \text{d.h. (A.9) ist erfüllt.}$$

$$(A.37) \quad f_u(x) = 0 \implies f(x) \leq 0,$$

$$(A.38) \quad f_u(x) < 0 \implies f(x) < 0,$$

$$(A.39) \quad f_d(x) > 0 \implies f(x) > 0,$$

$$(A.40) \quad f_d(x) = 0 \implies f(x) \geq 0,$$

$$(A.41) \quad f_d(x) < 0 \iff f(x) < 0, \quad \text{d.h. (A.10) ist erfüllt.}$$

## A.2. $x^2 - y^2$ , $x^2 + y^2$

Um einen vorzeitigen Überlauf zu vermeiden, benutzen wir  $x^2 - y^2 \equiv (|x| - |y|)(|x| + |y|)$ . Bei hinreichend großen Werten von  $|x|$  und  $|y|$  kann im Fall  $|x| \approx |y|$  auch jetzt noch die Summe  $(|x| + |y|)$  einen vorzeitigen Überlauf erzeugen, der jedoch mit der Skalierung  $|x| \cdot 2^{-2} + |y| \cdot 2^{-2}$  vermieden werden kann.

Es soll jetzt  $f(x, y) = x^2 - y^2$  abgerundet werden, wobei die Differenz  $A := |x| - |y| < 0$  als negativ vorausgesetzt wird. Die Summe  $B := |x| \cdot 2^{-2} + |y| \cdot 2^{-2} > 0$  ist positiv und wegen der Skalierung nur wenig kleiner als  $\text{MaxFloat}()$ . Nach (A.13) ist das Abrunden garantiert durch

$$(A.42) \quad * \quad A_d \leq 0; \quad B_u \geq 0 \quad \stackrel{(A.9)(A.10)}{\implies} \quad A_d \odot_d B_u \leq A \cdot B,$$

wobei  $*$  bedeutet, dass  $A_d = B_u = 0$  nicht eintreten darf. Um das Abrunden von  $f(x, y)$  zu gewährleisten, muss also  $A_d \leq 0$  und  $B_u \geq 0$  nachgewiesen werden.

Ganz allgemein gilt  $A_d \leq A$ , und wegen der Voraussetzung  $A < 0$  ist die Bedingung  $A_d \leq 0$  schon erfüllt, wobei  $A_d$  selbst mithilfe der MPFR-Funktion `mpfr_sub(..., ..., ..., RoundDown)` berechnet wird.

Da  $B > 0$  aufzurunden ist, müssen nach (A.6) die beiden Summanden  $|x| \cdot 2^{-2}$  und  $|y| \cdot 2^{-2}$  selbst aufgerundet werden durch:

$$\begin{aligned} \text{times2pown}(|x|, -2, \text{RoundUp}) &\longrightarrow |x|_u \geq |x| \cdot 2^{-2} > 0 \\ \text{times2pown}(|y|, -2, \text{RoundUp}) &\longrightarrow |y|_u \geq |y| \cdot 2^{-2} \geq 0; \end{aligned}$$

Da  $B > 0$  ist, gilt:  $B_u := |x|_u \oplus_u |y|_u \geq B > 0$ , d.h.  $B_u \geq B > 0$ , so dass damit auch die zweite Bedingung  $B_u \geq 0$  erfüllt ist und auch  $A_d = B_u = 0$  nicht eintreten kann. Zu beachten ist außerdem, dass bez.  $B_u$  die Bedingung (A.9) und bez.  $A_d$  die Bedingung (A.10) erfüllt ist, womit die gerichtete Abrundung des Funktionswertes  $f(x, y)$  gesichert ist. Die Berechnung von  $B_u$  erfolgt wieder mithilfe der MPFR-Funktion `mpfr_add(B_u, |x|_u, |y|_u, RoundUp)`;

Im Fall  $A := |x| - |y| \geq 0$  kann der Nachweis für eine gesicherte Abrundung analog geführt werden, und auch für die gesicherte Aufrundung von  $f(x, y)$  erfolgt der Nachweis ganz analog.

Wir betrachten noch die implementierte Funktion  $g(x, y) = x^2 + y^2$ . Da eine Skalierung einen Überlauf jetzt nicht verhindern kann, wird die Summe der Quadrate direkt ausgewertet. Um z.B. die Aufrundung von  $g(x, y)$  zu garantieren, muss nach (A.6)  $x \odot_u x \oplus_u y \odot_u y$  berechnet werden, wobei also beide Operanden von  $\oplus_u$  durch nur eine Rechenoperation realisiert werden, was nach den Bemerkungen von Seite 180 die garantierte Rundung beider Operanden gewährleistet. Das Aufrunden der Summe  $x \odot_u x \oplus_u y \odot_u y$  wird wieder realisiert mithilfe der MPFR-Funktion `mpfr_add(..., ..., ..., RoundUp)`;



### A.3. $x/(x^2 + y^2)$

Für  $x, y \in \mathbb{R} \setminus \{0\}$  sei  $f(x, y)$  definiert durch

$$(A.43) \quad f(x, y) := \frac{x}{x^2 + y^2}, \quad x \neq 0 \wedge y \neq 0.$$

Für Präzisionen  $\text{prec} \geq 2$  erfolgt die Auswertung von  $f(x, y)$  auf der Maschine durch

```
MpfrClass x_div_x2py2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

wobei bez. des Rundungsmodus  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  nahezu optimal gerundet wird. Ohne den obigen Parameter  $\text{rnd}$  erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus.

$f(x, y)$  wird mit  $x \neq 0 \wedge y \neq 0$  auf der Maschine berechnet durch:

$$(A.44) \quad f(x, y) = \begin{cases} 0, & x = 0 \\ 1/x, & y = 0 \\ \frac{x}{x^2 + y^2}, & \text{sonst.} \end{cases}$$

Die Auswertung von  $1/x$  erfolgt dabei mithilfe der Funktion

```
MpfrClass reci(const MpfrClass& x, RoundingMode rnd);
```

die in `mpfrclass.cpp` bereits definiert ist.

Um bei der Auswertung der Summe  $x^2 + y^2$  keinen vorzeitigen Überlauf zu erzeugen, suchen wir eine Konstante  $c$ , so dass für  $|x| \leq c$  und  $|y| \leq c$  auch im ungünstigsten Fall  $x^2 + y^2$  auf der Maschine ohne Überlauf berechnet werden kann. Der ungünstigste Fall ist dabei definiert durch  $\text{prec} = 2$  und durch den eingestellten Rundungsmodus `RoundUp`. Man erhält damit  $c := 2^{+536870910}$ , und um  $|x| = m \cdot 2^{\text{ex}_x} \leq 2^{+536870910}$  mit  $m \in [0.5, +1)$  garantieren zu können, verlangen wir daher mit  $\text{ex}_x = \text{expo}(x)$  die Beziehung  $\text{ex}_x \leq +536870910$ . Ganz entsprechend verlangen wir  $\text{ex}_y \leq c_1 := +536870910$ . Nur im Fall  $\text{ex}_x > c_1$  oder  $\text{ex}_y > c_1$  kann also ein vorzeitiger Überlauf eintreten. Mit  $\text{ex}_m := \text{Max}(\text{ex}_x, \text{ex}_y)$  wird dieser Überlauf verhindert, wenn im Nenner von  $x/(x^2 + y^2)$  die Beträge  $|x|, |y|$  mit  $2^{\text{ex}_s}$  und der Zähler  $x$  mit  $2^{2 \cdot \text{ex}_s}$  skaliert wird, wobei  $\text{ex}_s < 0$  definiert ist durch  $\text{ex}_m + \text{ex}_s = c_1$ .

Um bei der Auswertung von  $x^2 + y^2$  keinen vorzeitigen Unterlauf zu erhalten<sup>2</sup>, suchen wir z.B. für  $|x|$  eine möglichst kleine Obergrenze  $c$ , so dass im ungünstigsten Fall für  $|x| < c$  die Bedingung  $x^2 > m := \text{minfloat}()$  erfüllt ist. Der ungünstigste Fall ist dabei definiert durch  $\text{prec} = 2$  und durch den voreingestellten Rundungsmodus `RoundDown`. Man erhält damit  $c := 2^{-536870912}$ , so dass aus  $|x| \geq c \vee |y| \geq c$  stets folgt  $x^2 + y^2 > m$ , oder anders ausgedrückt:

Nur im Fall  $|x| < c \wedge |y| < c$  kann bei der Auswertung von  $x^2 + y^2$  ein Unterlauf eintreten.

Mit  $|x| = m_x \cdot 2^{\text{ex}_x}$ ,  $|y| = m_y \cdot 2^{\text{ex}_y}$ ,  $m_x, m_y \in [0.5, 1)$ ,  $\text{ex}_x = \text{expo}(x)$ ,  $\text{ex}_y = \text{expo}(y)$  und mit  $c_2 := -536870911$  erhält man die etwas gröbere Formulierung:

Nur für  $\text{ex}_x < c_2 \wedge \text{ex}_y < c_2$  kann bei der Auswertung von  $x^2 + y^2$  ein Unterlauf eintreten.

Mit  $\text{ex}_m := \text{Min}(\text{ex}_x, \text{ex}_y)$  wird dann ein solcher vorzeitiger Unterlauf vermieden, wenn man den Zähler  $x$  mit  $2^{2 \cdot \text{ex}_s}$  und im Nenner  $x, y$  jeweils mit  $2^{\text{ex}_s}$  skaliert, wobei  $\text{ex}_s$  definiert ist durch:  $\text{ex}_m + \text{ex}_s = c_2 + 1$ , um  $|x| > c$  und  $|y| > c$  garantieren zu können.

Als Beispiel wählen wir  $x = y = m := \text{minfloat}(\text{prec}) = 2^{-1073741824}$  und erhalten für den exakten Wert  $f(m, m) = 1/(2 \cdot m) = 2^{+1073741823} > \text{MaxFloat}(\text{prec})$ . `x_div_x2py2(..)` liefert für den aufgerundeten Funktionswert das korrekte Ergebnis `+infinity`, und für den optimal abgerundeten Funktionswert erhält man:  $\text{MaxFloat}(\text{prec}) = 2.0985787164673 \dots \cdot 10^{+323228496}$ .

<sup>2</sup>Ein vorzeitiger Unterlauf würde bei der Auswertung von  $x/(x^2 + y^2)$  eine Division durch Null erzeugen.

#### A.4. $(x^2 - y^2)/(x^2 + y^2)^2$

Für  $x, y \in \mathbb{R} \setminus \{0\}$  sei  $f(x, y)$  definiert durch

$$(A.45) \quad f(x, y) := \frac{x^2 - y^2}{(x^2 + y^2)^2}, \quad x \neq 0 \wedge y \neq 0.$$

Für Präzisionen  $\text{prec} \geq 2$  erfolgt die Auswertung von  $f(x, y)$  auf der Maschine durch

```
MpfrClass Re_rz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

wobei bez. des Rundungsmodus  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  nahezu optimal gerundet wird. Ohne den obigen Parameter  $\text{rnd}$  erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Die interne Auswertung von  $f(x, y)$  erfolgt beim Zähler und Nenner mittels geeigneter Skalierungen ohne einen vorzeitigen Über- oder Unterlauf. Im Fall  $x = y = 0$  erfolgt eine entsprechende Fehlermeldung.  $f(x, y)$  kommt u.a. zur Anwendung als Realteil der komplexen Funktion  $w(z) = 1/z^2$ ,  $z = x + i \cdot y \in \mathbb{C}$ .

Der Zähler  $Z(x, y) = x^2 - y^2$  wird realisiert mit Hilfe der internen Funktion

```
void Zae_x2my2(MpfrClass& d, long int& k, const MpfrClass& x,
               const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^k \cdot d = x^2 - y^2$ , wobei intern bei der Berechnung von  $d$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|d|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt.

Ganz entsprechend wird der Nenner  $N(x, y) = (x^2 + y^2)^2$  realisiert durch die interne Funktion

```
void Ne_x2py2p2(MpfrClass& r, long int& kn, const MpfrClass& x,
                const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^{kn} \cdot r = (x^2 + y^2)^2$ , wobei intern bei der Berechnung von  $r$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|r|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt. Anschließend kann die notwendige Division  $d/r$  ohne vorzeitigen Über- oder Unterlauf durchgeführt werden. Nur bei der letzten Multiplikation mit  $2^{k-kn}$  kann dann ein nicht vermeidbarer Über- oder Unterlauf eintreten.

##### A.4.1. Numerische Beispiele

Im **1. Beispiel** wähle wir  $x = \text{MaxFloat}()$ ,  $y = \text{MaxFloat}()/3$ , d.h.  $0 < f(x, y) < \text{minfloat}()$  ist sicher erfüllt. Mit dem Funktionsaufruf `res = Re_rz2(x, y, RoundDown);` erhalten wir daher `res = 0` und `res = Re_rz2(x, y, RoundUp);` liefert `res = minfloat()`.

Im **2. Beispiel** wähle wir  $x = \text{minfloat}()$ ,  $y = 3 \cdot \text{minfloat}()$ , d.h.  $f(x, y) < -\text{MaxFloat}()$  ist sicher erfüllt. Mit dem Funktionsaufruf `res = Re_rz2(x, y, RoundDown);` erhalten wir daher `res = -@Inf@`, und der Funktionsaufruf `res = Re_rz2(x, y, RoundUp);` liefert folglich: `res = -MaxFloat() = -2.098...10323228496`.

Im **3. Beispiel** wähle wir  $x = 1/3$ ,  $y = 2/3$  und erhalten  $f(1/3, 2/3) = -27/25 = -1.08000\dots$  Mit dem Funktionsaufruf `res = Re_rz2(x, y, RoundDown);` erhalten wir mit der Präzision  $\text{prec} = 300$  das nahezu optimal abgerundete Ergebnis

```
res = -1.0800000000000000...00000000000000388
```

wobei insgesamt 93 Dezimalstellen ausgegeben wurden. Beachten Sie, dass sowohl Zähler als auch Nenner mit den Werten  $-1/3$  bzw.  $25/81$  jeweils periodische Dezimalbrüche sind, die daher bei der internen Berechnung entsprechend gerundet werden müssen.

## A.5. $2xy/(x^2 + y^2)^2$

Für  $x, y \in \mathbb{R} \setminus \{0\}$  sei  $f(x, y)$  definiert durch

$$(A.46) \quad f(x, y) := \frac{2xy}{(x^2 + y^2)^2}, \quad x \neq 0 \wedge y \neq 0.$$

Für Präzisionen  $\text{prec} \geq 2$  erfolgt die Auswertung von  $f(x, y)$  auf der Maschine durch

```
MpfrClass mIm_rz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

wobei bez. des Rundungsmodus  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  nahezu optimal gerundet wird. Ohne den obigen Parameter  $\text{rnd}$  erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Die interne Auswertung von  $f(x, y)$  erfolgt beim Zähler und Nenner mittels geeigneter Skalierungen ohne einen vorzeitigen Über- oder Unterlauf. Im Fall  $x = y = 0$  erfolgt eine entsprechende Fehlermeldung.  $-f(x, y)$  kommt zur Anwendung als Imaginärteil der komplexen Funktion  $w(z) = 1/z^2$ ,  $z = x + i \cdot y \in \mathbb{C}$ .

Der Zähler  $Z(x, y) = 2xy$  wird realisiert mit Hilfe der internen Funktion

```
void Zae_2xy(MpfrClass& z, long int& k, const MpfrClass& x,  
            const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^k \cdot z = 2 \cdot (xy)$ , wobei intern bei der Berechnung von  $z$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|z|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt.

Ganz entsprechend wird der Nenner  $N(x, y) = (x^2 + y^2)^2$  realisiert durch die interne Funktion

```
void Ne_x2py2p2(MpfrClass& r, long int& kn, const MpfrClass& x,  
              const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^{kn} \cdot r = (x^2 + y^2)^2$ , wobei intern bei der Berechnung von  $r$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|r|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt. Anschließend kann die notwendige Division  $z/r$  ohne vorzeitigen Über- oder Unterlauf durchgeführt werden. Nur bei der letzten Multiplikation mit  $2^{k-kn}$  kann dann ein nicht vermeidbarer Über- oder Unterlauf eintreten.

Der Ausdruck  $-2xy/(x^2 + y^2)^2$  ist der Imaginärteil von  $1/z^2$ ,  $z = x + i \cdot y$ , und wird analog realisiert durch die Funktion  $\text{Im\_rz2}(x, y, \text{rnd})$ , die in `mpfrclass.cpp` definiert ist.

### A.5.1. Numerische Beispiele

**Im 1. Beispiel** wähle wir  $x = \text{MaxFloat}()$ ,  $y = \text{MaxFloat}()/3$ , d.h.  $0 < f(x, y) < +\text{minfloat}()$  ist sicher erfüllt. Mit dem Funktionsaufruf `res = mIm_rz2(x, y, RoundUp);` erhalten wir daher `res = +minfloat()` und `res = mIm_rz2(x, y, RoundDown);` liefert `res = 0`.

**Im 2. Beispiel** wähle wir  $x = -\text{minfloat}()$ ,  $y = 3 \cdot \text{minfloat}()$ , d.h.  $f(x, y) < -\text{MaxFloat}()$  ist sicher erfüllt. Mit dem ersten Funktionsaufruf `res = mIm_rz2(x, y, RoundUp);` erhalten wir folgerichtig `res = -MaxFloat() = -2.098... \cdot 10^{323228496}`, und der nächste Funktionsaufruf `res = mIm_rz2(x, y, RoundDown);` liefert wie erwartet: `res = -0Inf@`.

**Im 3. Beispiel** wähle wir  $x = 2$ ,  $y = 3$  und erhalten  $f(2, 3) = 12/169 = 7.10059116... \cdot 10^{-2}$ . Mit dem Funktionsaufruf `res = mIm_rz2(x, y, RoundDown);` erhalten wir mit der Präzision  $\text{prec} = 300$  das nahezu optimal abgerundete Ergebnis

$$\text{res} = 7.100591715976331360946745562130... 171597379 \cdot 10^{-2}$$

wobei insgesamt 93 Dezimalstellen ausgegeben wurden.



### A.7. $(1 + x^2 - y^2)/(4x^2y^2 + (1 + x^2 - y^2)^2)$

$f : \mathbb{R}^2 \setminus \{(0, \pm 1)\} \rightarrow \mathbb{R}$  sei definiert durch

$$f(x, y) := \frac{1 + x^2 - y^2}{4x^2y^2 + (1 + x^2 - y^2)^2}$$

wobei man sich bei der Auswertung auf den 1. Quadranten beschränken kann, d.h. wir betrachten nur die Funktion

$$(A.48) \quad f(x, y) := \frac{1 + x^2 - y^2}{4x^2y^2 + (1 + x^2 - y^2)^2} \quad \text{mit } x \geq 0, y \geq 0 \text{ und } (x, y) \neq (0, +1),$$

und mit  $\varepsilon > 0$  gilt zusätzlich:  $\lim_{\varepsilon \rightarrow 0} f(0, 1 \pm \varepsilon) = \mp\infty$ ,  $\lim_{\varepsilon \rightarrow 0} f(\varepsilon, +1) = 1/4$ .

Für die Präzisionen  $\text{prec} \geq 2$  erfolgt die Auswertung von  $f(x, y)$  auf der Maschine durch

```
MpfrClass Re_r1pz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

wobei bez. des Rundungsmodus  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  nahezu optimal gerundet wird. Ohne den obigen Parameter  $\text{rnd}$  erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Die interne Auswertung von  $f(x, y)$  erfolgt beim Zähler und Nenner mittels geeigneter Skalierungen ohne einen vorzeitigen Über- oder Unterlauf. Im Fall  $x = 0, y = +1$  erfolgt eine entsprechende Fehlermeldung.  $f(x, y)$  kommt u.a. zur Anwendung als Realteil der komplexen Funktion  $w(z) = 1/(1 + z^2)$ ,  $z = x + i \cdot y \in \mathbb{C}$ .

Der Zähler  $Z(x, y) := 1 + x^2 - y^2$  wird realisiert mit Hilfe der internen Funktion

```
void x2p1my2(MpfrClass& d, long int& k, const MpfrClass& x,
             const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^k \cdot d = 1 + x^2 - y^2$ , wobei intern bei der Berechnung von  $d$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|d|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt. Die interne Berechnung von  $d$  erfolgt in zwei Schritten.

Im Fall  $x \approx y$  wird  $x^2 - y^2$  über sein Produkt  $(x - y) \cdot (x + y)$  ausgewertet, da jetzt wegen der exakten Maschinenwerte  $x, y$  ihre Differenz  $(x - y)$  **ohne** Auslöschung, d.h. bis auf den normalen Rundungsfehler **ohne** Informationsverlust, berechnet wird. Anschließend erfolgt die Addition der Eins, wobei im Fall  $x^2 - y^2 \approx -1$  bzw.  $y \approx \sqrt{1 + x^2}$  eine jetzt wirkliche Auslöschung mit Informationsverlust durch die schrittweise Auswertung mit wachsender Präzision garantiert vermieden wird. Beachten Sie, dass selbst bei einer auftretenden Auslöschung bez. der mit  $\text{rnd}$  gewählten Rundungen **stets korrekt** auf- bzw. abgerundet wird!

Falls sich  $x, y$  deutlicher unterscheiden, wird  $1 - y^2$  über sein Produkt  $(1 - y) \cdot (1 + y)$  ausgewertet und danach erfolgt die Addition von  $x^2$ , wobei im Fall  $y \approx \sqrt{1 + x^2}$  eine Auslöschung mit Informationsverlust durch die schrittweise Auswertung mit wachsender Präzision garantiert vermieden wird.

Ganz analog wird der Nenner  $N(x, y) := 4x^2y^2 + (1 + x^2 - y^2)^2$  realisiert durch die interne Funktion

```
void Ne_r1px2(MpfrClass& r, long int& kn, const MpfrClass& x,
             const MpfrClass& y, RoundingMode rnd);
```

dabei gilt exakt, d.h. ohne Rundung:  $2^{kn} \cdot r = 4x^2y^2 + (1 + x^2 - y^2)^2$ , wobei intern bei der Berechnung von  $r$  bezüglich  $\text{rnd}$  gerundet wird. Mit geeigneter Skalierung erreicht man, dass  $|r|$  nur wenig kleiner als  $\text{MaxFloat}(\text{prec})$  ausfällt. Anschließend kann die notwendige Division  $d/r$  ohne vorzeitigen Über- oder Unterlauf durchgeführt werden. Nur bei der letzten Multiplikation mit  $2^{k-kn}$  kann dann ein nicht vermeidbarer Über- oder Unterlauf eintreten. Beachten Sie, dass in  $N(x, y)$  beide Summanden nicht negativ sind, so dass eine Auslöschung nicht auftreten kann.



## A.8. $x^2 + a \cdot x + b$

$f: \mathbb{R} \rightarrow \mathbb{R}$  sei definiert durch

$$f(x, a, b) := x^2 + a \cdot x + b; \quad x, a, b \in \mathbb{R}.$$

Für die Präzisionen  $\text{prec} \geq 2$  erfolgt die Auswertung von  $f(x, a, b)$  auf der Maschine zunächst durch die Funktion

```
void poly2(const MpfrClass& r, long int& k, const MpfrClass& x,
          const MpfrClass& a, const MpfrClass& b, RoundingMode rnd);
```

Dabei gilt ohne Berücksichtigung einer Rundung  $2^k \cdot r = x^2 + a \cdot x + b$ , wobei aber  $r$  bei der internen Berechnung bez. des Parameters  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  **optimal** gerundet wird. Ohne den obigen Parameter  $\text{rnd}$  erfolgt die Rundung nach dem voreingestellten Current-Rundungsmodus. Durch geeignete Skalierungen wird  $r \neq 0$  so berechnet, dass  $|r|$  nur etwas kleiner als  $\text{MaxFloat}()$  ausfällt, wobei ein vorzeitiger Überlauf oder Unterlauf vermieden wird. Im Fall  $r = 0$  gilt  $k = 0$ . Die optimale Rundung wird durch eine schrittweise Vergrößerung der internen Präzision erreicht. Im letzten Schritt wird dann in der Funktion

```
MpfrClass poly2(const MpfrClass& x, const MpfrClass& a,
               const MpfrClass& b, RoundingMode rnd);
```

das bez.  $\text{rnd}$  gerundete Produkt  $2^k \cdot r$  ausgewertet, wobei erst dabei ein dann unvermeidbarer Überlauf oder Unterlauf eintreten kann.

### A.8.1. Numerische Beispiele

Zunächst ein allgemeiner **Hinweis**: Bei der Funktion  $f(x, a, b)$ , wie auch bei allen Funktionen aus Tabelle 3.11.3, wird vorausgesetzt, dass  $x, a, b$  exakte Maschinenzahlen sind, die also nicht schon durch vorhergehende Rechnungen gerundet worden sind. Beispielsweise kann daher eine Maschinenzahl  $a$  nicht den Wert  $\sqrt{3}$  annehmen, wohl aber eine geeignete Maschinennäherung, die dann aber als **exakt** anzusehen ist.

Im **1. Beispiel** wählen wir  $x = 10$ ,  $a = 3$ ,  $b = -3$  und erhalten mit z.B.  $\text{prec} = 300$  für das auf- und abgerundete Ergebnis jeweils den gleichen, exakten Wert  $f(x, a, b) = 127$ .

Im **2. Beispiel** wählen wir  $a = 0$ ,  $b = -1$  und werten mit der Präzision  $\text{prec} = 300$  den Ausdruck  $f(x, a, b) = x^2 - 1$  in der unmittelbaren Nähe seiner Nullstelle  $x_0 = +1$  für  $x = \text{pred}(1) = 1 - 2^{-300}$  aus, so dass bei der internen Auswertung von

$$f(\text{pred}(1), 0, -1) = 2^{-300} \cdot (-2 + 2^{-300}) = -9.818186930595453106 \dots 2019576566474 \dots \cdot 10^{-91}$$

maximale Auslöschung auftreten muss. Mit der Funktion  $\text{poly2}(x, a, b, \text{rnd})$  erhält man für den **abgerundeten** und **aufgerundeten** Maschinenwert das Ergebnis

$$\text{poly2}(x, a, b, \text{rnd}) = -9.818186930595453106 \dots 2019576 \overset{325}{\underset{307}{566474}} \dots \cdot 10^{-91},$$

wobei die gerundeten Werte den exakten Wert einschließen und in den ersten 90 Dezimalziffern übereinstimmen. Wegen  $300 \cdot \log_{10}(2) = 90.308998 \dots$  erhalten wir daher trotz maximaler Auslöschung mit 90 die **Maximalzahl** übereinstimmender Dezimalziffern.

## A.9. $2 \sin(x) \cosh(y) / (\cosh(2y) - \cos(2x))$

Die reelle Funktion  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , mit

$$(A.49) \quad f(x, y) := \frac{2 \sin(x) \cdot \cosh(y)}{\cosh(2y) - \cos(2x)} = \frac{\sin(x) \cdot \cosh(y)}{\sin^2(x) + \sinh^2(y)}$$

ist der Realteil der komplexen Funktion  $1/\sin(z)$ ,  $z = x + i \cdot y \in \mathbb{C}$ . Die Funktionswerte  $f(x, y)$  werden in nahezu allen Fällen optimal approximiert durch

`MpfrClass Re_csc(const MpfrClass& x, const MpfrClass& y, const RoundingMode rnd),`

wobei für den Rundungsparameter `rnd` die folgenden Rundungsmodi zur Verfügung stehen:

$$\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}.$$

$f(x, y)$  besitzt Singularitäten an den Stellen  $S_k = (x_k, 0) := (k\pi, 0)$ ,  $k \in \mathbb{Z}$ . Parallel zur  $y$ -Achse im Abstand  $x_k$  verschwinden für  $y \neq 0$  alle Funktionswerte  $f(x_k, y) = 0$ , wobei die Funktion an den Stellen  $S_k = (x_k, 0)$  nicht definiert ist. Auf der  $x$ -Achse selbst besitzt  $f(x, 0)$  an den Stellen  $x_k$  Polstellen mit Vorzeichenwechsel, d.h. für  $x \rightarrow 0$  gilt beispielsweise

$$\lim_{x \rightarrow 0^-} f(x, 0) = \lim_{x \rightarrow 0^-} \frac{1}{\sin(x)} = -\infty, \quad \lim_{x \rightarrow 0^+} f(x, 0) = \lim_{x \rightarrow 0^+} \frac{1}{\sin(x)} = +\infty.$$

Bezüglich der Variablen  $x$  ist  $f(x, y)$   $2\pi$ -periodisch, d.h.  $f(x + k \cdot 2\pi, y) = f(x, y)$ ,  $k \in \mathbb{Z}$ . Wegen  $f(x, y) \equiv f(x, |y|)$  kann man sich auf die obere Halbebene, d.h.  $y \geq 0$ , beschränken.

Die interne Auswertung erfolgt intervallmäßig mithilfe von Punktintervallen  $x \in [x], y \in [y]$  in hinreichend hoher Präzision. Ist beispielsweise `Fxy` eine optimale Einschließung für  $f(x, y)$ , so ist z.B. der aufgerundete Funktionswert gegeben durch `Sup(Fxy) ≥ f(x, y)`, wobei  $x, y$  jetzt als Maschinenzahlen zu verstehen sind.

### A.9.1. Optimale Einschließung

Für eine optimale Einschließung von  $f(x, y)$  muss verhindert werden, dass bei intervallmäßiger Auswertung von (A.49) im Zähler und Nenner ein vorzeitiger Unter- oder Überlauf entsteht. Um bei  $\cosh(2y)$  oder  $\sinh^2(y)$  einen Überlauf zu vermeiden, betrachten wir zunächst den

**Fall 1.**  $0 \leq y \leq y_0 := 372130555$ , wodurch ein solcher Überlauf verhindert wird. Ausgewertet wird der zweite Ausdruck in (A.49), weil sein Nenner ohne Auslöschung berechnet werden kann. Bei der intervallmäßigen Auswertung von  $\sin([x])$ ,  $[x] = [\text{minfloat}()]$ , entsteht ein vorzeitiger Unterlauf, der vermieden wird, wenn zunächst mit 2 erweitert wird:

$$(A.50) \quad f(x, y) = \frac{2 \cdot \sin(x) \cdot \cosh(y)}{2 \cdot (\sin^2(x) + \sinh^2(y))}$$

und im Zähler im Fall  $|x| \leq 2 \cdot \text{minfloat}()$  der Term  $2 \cdot \sin(x)$  ersetzt wird durch

$$2 \cdot \sin(x) = \frac{\sin(2x)}{\cos(x)},$$

wobei  $\sin(2x)$  jetzt auch für  $[x] = [\text{minfloat}()]$  ohne Unterlauf berechnet werden kann.  $1/\cos(x)$  lässt sich mithilfe der geometrischen Reihe abschätzen durch

$$1 \leq \frac{1}{\cos(x)} = 1 + \frac{x^2}{2} + \frac{5x^4}{24} \dots \leq \frac{1}{1-x^2} \leq 1 + 2x^2, \quad |x| \leq 2 \cdot \text{minfloat}(),$$

und wir verlangen zusätzlich:  $1 + 2x^2 \leq \text{succ}(1) = 1 + 2^{1-\text{prec}} \iff x^2 \leq 2^{-\text{prec}}$ .  
Mit  $\text{minfloat}() = 2^{-1073741824}$  und  $x^2 \leq 4 \cdot \text{minfloat}()^2$  lautet die letzte Bedingung:



$2^{-2147483646} \leq 2^{-\text{prec}}$  bzw.  $\text{prec} \leq +2147483646$ , und diese Bedingung wird in der Praxis schon aus Laufzeitgründen stets erfüllt sein!  $1/\cos(x)$  kann damit sehr einfach eingeschlossen werden durch:

$$(A.51) \quad 1 \leq \frac{1}{\cos(x)} \leq \text{succ}(1), \quad \text{falls } |x| \leq 2 \cdot \text{minfloat}().$$

Der Zähler in (A.50) kann daher effektiv ohne Über- und Unterlauf eingeschlossen werden durch

$$2 \cdot \sin(x) \cdot \cosh(y) \in \sin(2 \diamond [x]) \diamond [1, \text{succ}(1)] \diamond \cosh([y]).$$

Um in (A.50) den Nenner  $2 \cdot (\sin^2(x) + \sinh^2(y)) \geq 2 \cdot \sinh^2(y) \geq 2y^2$  ohne vorzeitigen Unterlauf berechnen zu können, verlangen wir zunächst  $2y^2 \geq 32 \cdot \text{minfloat}() = 2^{-1073741819}$ , und diese Bedingung<sup>3</sup> ist für  $y \geq 2^{-536870910}$  erfüllt. Nur im Fall  $y < 2^{-536870910}$  bzw.  $\text{expo}(y) < -536870909$  kann daher bei intervallmäßiger Auswertung die Einschließung des Nenners die Null enthalten, wodurch bei der nachfolgenden Intervalldivision starke Überschätzungen auftreten<sup>4</sup>. Um dies zu vermeiden, werden Zähler und Nenner im Fall  $\text{expo}(y) < -536870909$  skaliert, d.h. Zähler und Nenner erhalten die Form  $m_{1,2} \cdot 2^{k_{1,2}}$ , wobei die schmalen Intervalle  $m_{1,2}$  in der Nähe von  $\pm 1$  liegen und die Division dieser Intervalle ohne Unter- oder Überlauf möglich ist. Erst bei der Multiplikation mit der Zweierpotenz  $2^{k_1 - k_2}$  kann dann ein unvermeidbarer Unter- oder Überlauf auftreten.

Beachten Sie, dass wegen der Taylorreihe  $\sinh(y) = y + y^3/3! + y^5/5! + \dots$  ein Unterlauf bei der Berechnung von  $\sinh([\text{minfloat}()])$  nicht auftreten kann, aber  $2 \cdot \sin^2(x)$  durch

$$2 \cdot \sin^2(x) = \frac{\sin^2(2x)}{2 \cdot \cos^2(x)} = \frac{1}{2} \cdot \left( \frac{\sin(2x)}{\cos(x)} \right)^2$$

zu ersetzen ist, um  $\sin(2[\text{minfloat}()])$  intervallmäßig ohne Unterlauf berechnen zu können. Das Quadrieren erfolgt mit skalierten Intervallen für  $\sin(2[x])$  und  $\sec([x]) = 1/\cos([x])$ , um einen möglichen Unterlauf zu vermeiden. Im Fall  $2^{-536870910} \leq y \leq 372130555$  kann der rechte Term in (A.49) direkt ohne Unter- oder Überlauf intervallmäßig ausgewertet werden. Weitere Einzelheiten findet man im Quellcode der Funktion `Re_csc(x, y, rnd)` in `mpfrclass.cpp`. Es bleibt jetzt noch der

**Fall 2.**  $y > y_0 = 372130555$ , in dem ein Überlauf bei der Auswertung von  $\cosh(2y)$  in (A.49) zu vermeiden ist. Wir schreiben (A.49) zunächst in der Form

$$(A.52) \quad f(x, y) = 2 \cdot \frac{\cosh(y)}{\cosh(2y)} \cdot \frac{\sin(x)}{1 - \frac{\cos(2x)}{\cosh(2y)}}, \quad \text{mit} \quad \frac{\cosh(y)}{\cosh(2y)} = e^{-y} \cdot \frac{1 + e^{-2y}}{1 + e^{-4y}}, \quad \text{und}$$

$$(A.53) \quad 1 < \frac{1 + e^{-2y}}{1 + e^{-4y}} < 1 + e^{-2y}, \quad \text{d.h.} \quad \frac{\cosh(y)}{\cosh(2y)} \in e^{-y} \cdot [1, 1 + e^{-2y}].$$

Mit  $1/\cosh(2y) = 2/(e^{2y} + e^{-2y})$  folgt weiter

$$(A.54) \quad 1 - \frac{\cos(2x)}{\cosh(2y)} \in 1 - \frac{[-1, +1]}{\cosh(2y)} = 1 - \frac{[-2, +2]}{e^{2y} + e^{-2y}} = \left[ 1 - \frac{2}{e^{2y} + e^{-2y}}, 1 + \frac{2}{e^{2y} + e^{-2y}} \right]$$

$$(A.55) \quad \subset [1 - 2e^{-2y}, 1 + 2e^{-2y}].$$

Mit diesen Abschätzungen erhält man schließlich nach einfachen Rechnungen:

$$\frac{\cosh(y)}{\cosh(2y) \cdot \left( 1 - \frac{\cos(2x)}{\cosh(2y)} \right)} \in \frac{e^{-y} \cdot [1, 1 + e^{-2y}]}{[1 - 2e^{-2y}, 1 + 2e^{-2y}]} = e^{-y} \cdot \left[ \frac{1}{1 + 2e^{-2y}}, \frac{1 + e^{-2y}}{1 - 2e^{-2y}} \right] \\ \subset e^{-y} \cdot [1 - 2e^{-2y}, 1 + 4e^{-2y}].$$

<sup>3</sup>Der Faktor 32 berücksichtigt mögliche Überschätzungen bei der intervallmäßigen Auswertung des Nenners.

<sup>4</sup>Im Fall  $x = y = 0$  ist die Division durch Null natürlich unvermeidbar, aber `Re_csc(0, 0, rnd)` erzeugt an dieser Singularität eine entsprechende Fehlermeldung.

Zur Vereinfachung verlangen wir noch zusätzlich

$$(A.56) \quad [1 - 2e^{-2y}, 1 + 4e^{-2y}] \subseteq [\text{pred}(1), \text{succ}(1)] = [1 - 2^{-\text{prec}}, 1 + 2^{1-\text{prec}}], \quad \text{d.h. wegen } y \geq y_0 \\ 1 - 2^{-\text{prec}} < 1 - 2 \cdot e^{-2y_0} \quad \text{und} \quad 1 + 4 \cdot e^{-2y_0} < 1 + 2^{1-\text{prec}}.$$

Die erste Ungleichung in (A.56) ist erfüllt für  $\text{prec} < 2y_0/\ln(2) - 1 = 1073741811,523\dots$ , und diese Bedingung wird in der Praxis schon aus Laufzeitgründen immer erfüllt sein. Es bleibt dem Leser überlassen, dass auch die zweite Ungleichung in (A.56) mit der gleichen Bedingung erfüllt wird.  $f(x, y)$  kann damit durch folgenden Intervallausdruck eingeschlossen werden:

$$f(x, y) \in ([\text{pred}(1), \text{succ}(1)] \diamond e^{-[y]+\ln(2)}) \diamond \sin([x]), \quad y \geq y_0 = 372130555.$$

Beachten Sie dabei, dass jetzt mit  $x = \text{minfloat}()$  der **Unterlauf** bei  $\sin(x) \in [0, \text{minfloat}())$  kein Unglück ist, da das Infimum des obigen Intervallausdrucks (...) sicher kleiner 1 ist, so dass das Infimum des ganzen Intervallausdrucks 0 sein muss. Ganz Entsprechendes gilt auch im Fall  $x = -\text{minfloat}()$ .

## A.9.2. Numerische Beispiele

Im **1. Beispiel** wählen wir mit  $x_1 = \text{minfloat}() = 2^{-1073741824}$  und  $y_1 = 0$  einen Punkt ganz in der Nähe der Singularität  $S_0 = (0, 0)$  und erhalten mit  $\text{prec} = 200$  für den optimal ab- und aufgerundeten Funktionswert

$$\begin{aligned} f(x_1, y_1, \text{RoundDown}) &= 2.09857871646738769240435811688\dots335262778664e323228496 \\ f(x_1, y_1, \text{RoundUp}) &= @\text{Inf}@ \end{aligned}$$

Im **2. Beispiel** wählen wir  $x_2 = 3 \cdot \text{minfloat}() = 3 \cdot 2^{-1073741824}$  und  $y_2 = 0$  und erhalten mit  $\text{prec} = 200$  für den optimal ab- und aufgerundeten Funktionswert

$$\begin{aligned} f(x_2, y_2, \text{RoundDown}) &= 1.399052477644925128269572077922\dots890175185776e323228496 \\ f(x_2, y_2, \text{RoundUp}) &= 1.399052477644925128269572077922\dots890175185777e323228496 \end{aligned}$$

Im **3. Beispiel** wählen wir  $x_3 = \text{minfloat}() = 2^{-1073741824}$  und  $y_3 = \text{Maxfloat}()$  und erhalten mit  $\text{prec} = 200$  für den optimal ab- und aufgerundeten Funktionswert

$$\begin{aligned} f(x_3, y_3, \text{RoundDown}) &= 0 \\ f(x_3, y_3, \text{RoundUp}) &= 2.382564904887951073216169781732\dots23978795503e-323228497 \end{aligned}$$

Im **4. Beispiel** wählen wir  $x_4 = 2$  und  $y_4 = 0.25$  und erhalten mit  $\text{prec} = 200$  für den optimal ab- und aufgerundeten Funktionswert

$$\begin{aligned} f(x_4, y_4, \text{RoundDown}) &= 1.0530256453589318310512458485994\dots69109247482859536177 \\ f(x_4, y_4, \text{RoundUp}) &= 1.0530256453589318310512458485994\dots69109247482859536178 \end{aligned}$$

Im **5. Beispiel** wählen wir  $x_5 = -2$  und  $y_5 = 372130600$  und erhalten mit  $\text{prec} = 200$  für den optimal ab- und aufgerundeten Funktionswert

$$\begin{aligned} f(x_5, y_5, \text{RoundDown}) &= -1.356406034350108628717756434978\dots813293104e-161614266 \\ f(x_5, y_5, \text{RoundUp}) &= -1.356406034350108628717756434978\dots813293103e-161614266 \end{aligned}$$

## A.10. $2 \cos(x) \sinh(y) / (\cos(2x) - \cosh(2y))$

Die reelle Funktion  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ , mit

$$(A.57) \quad f(x, y) := \frac{2 \cos(x) \cdot \sinh(y)}{\cos(2x) - \cosh(2y)} = \frac{-\cos(x) \cdot \sinh(y)}{\sin^2(x) + \sinh^2(y)}$$

ist der Imaginärteil der komplexen Funktion  $1/\sin(z)$ ,  $z = x + i \cdot y \in \mathbb{C}$ .  $f(x, y)$  wird in nahezu allen Fällen optimal approximiert durch

```
MpfrClass Im_csc(const MpfrClass& x, const MpfrClass& y, const RoundingMode rnd),
```

wobei für den Rundungsparameter `rnd` die folgenden Rundungsmodi zur Verfügung stehen:

`rnd`  $\in$  {RoundUp, RoundDown, RoundNearest}.

$f(x, y)$  besitzt Singularitäten an den Stellen  $S_k = (x_k, 0) := (k\pi, 0)$ ,  $k \in \mathbb{Z}$ . Parallel zur  $y$ -Achse im Abstand  $x_k$  gilt  $f(x_k, y) = (-1)^{k+1}/\sinh(y)$ , d.h. es liegen Polstellen mit Vorzeichenwechsel vor. Wegen  $f(x, 0) = 0/\sin^2(x)$  verschwinden die Funktionswerte auf der  $x$ -Achse für  $x \neq x_k$ . Bezüglich der Variablen  $x$  ist  $f(x, y)$   $2\pi$ -periodisch, d.h.  $f(x + k \cdot 2\pi, y) = f(x, y)$ ,  $k \in \mathbb{Z}$ . Außerdem gilt  $f(x, -y) = -f(x, y)$ . Die interne Auswertung erfolgt intervallmäßig mithilfe von Punktintervallen  $x \in [x], y \in [y]$  in hinreichend hoher Präzision. Ist beispielsweise `Fxy` eine optimale Einschließung für  $f(x, y)$ , so ist z.B. der aufgerundete Funktionswert gegeben durch  $\text{Sup}(\text{Fxy}) \geq f(x, y)$ , wobei  $x, y$  jetzt als Maschinenzahlen zu verstehen sind.

### A.10.1. Optimale Einschließung

Für eine optimale Einschließung von  $f(x, y)$  muss verhindert werden, dass bei intervallmäßiger Auswertung von (A.57) im Zähler und Nenner ein vorzeitiger Unter- oder Überlauf entsteht. Um bei  $\sinh^2(y)$  einen Überlauf zu vermeiden, betrachten wir zunächst den

**Fall 1.**  $0 \leq |y| \leq y_0 := 372130555$ , wodurch ein solcher Überlauf verhindert wird. Ausgewertet wird der zweite Ausdruck in (A.57), weil sein Nenner ohne Auslöschung berechnet werden kann.

Um bei der Auswertung von (A.57) einen vorzeitigen Über- oder Unterlauf zu vermeiden, werden Zähler und Nenner skaliert, d.h. Zähler und Nenner erhalten die Form  $m_{1,2} \cdot 2^{k_{1,2}}$ , wobei die schmalen Intervalle  $m_{1,2}$  in der Nähe von  $\pm 1$  liegen und die Division dieser Intervalle ohne Unter- oder Überlauf möglich ist. Erst bei der Multiplikation mit der Zweierpotenz  $2^{k_1 - k_2}$  kann dann ein unvermeidbarer Unter- oder Überlauf auftreten.

Beachten Sie, dass wegen der Taylorreihe  $\sinh(y) = y + y^3/3! + y^5/5! + \dots$  ein Unterlauf bei der Berechnung von  $\sinh(\text{[minfloat()]})$  nicht auftreten kann, aber  $\sin^2(x)$  durch

$$\sin^2(x) = \frac{1}{4} \cdot \left( \frac{\sin(2x)}{\cos(x)} \right)^2$$

zu ersetzen ist, um für  $x = \text{minfloat}()$  den rechten Term  $\sin(2[\text{minfloat()}])$  intervallmäßig ohne Unterlauf berechnen zu können. Das Quadrieren erfolgt mit skalierten Intervallen für  $\sin(2[x])$  und  $\cos([x])$ , um einen möglichen Unterlauf zu vermeiden. Weitere Einzelheiten findet man im Quellcode der Funktion `Im_csc(x, y, rnd)` in `mpfrclass.cpp`. Es bleibt jetzt noch der

**Fall 2.**  $|y| > y_0 = 372130555$ , in dem ein Überlauf bei der Auswertung von  $\sinh(2y)$  in (A.57) zu vermeiden ist. Wir schreiben mit  $\sinh(y) = \text{sign}(y) \cdot \sinh(|y|)$  (A.57) zunächst in der Form

$$(A.58) \quad f(x, y) = \frac{-\cos(x)}{\sinh(y) \cdot \left( 1 + \frac{\sin^2(x)}{\sinh^2(y)} \right)} = -\cos(x) \cdot \frac{2 \cdot \text{sign}(y)}{e^{|y|} - e^{-|y|}} \cdot \frac{1}{1 + \frac{\sin^2(x)}{\sinh^2(y)}}$$

und werden versuchen, die beiden letzten Brüche rechts weiter zu vereinfachen. Es gilt zunächst

$$(A.59) \quad \frac{2}{e^{|y|} - e^{-|y|}} = 2e^{-|y|} \cdot \frac{1}{1 - e^{-2|y|}}, \quad \text{mit: } 1 + e^{-2|y|} \leq \frac{1}{1 - e^{-2|y|}} \leq 1 + 2e^{-2|y|},$$

und daraus ergibt sich für  $f(x, y)$  die Einschließung

$$(A.60) \quad f(x, y) \in -2 \cdot \text{sign}(y) \cdot \cos([x]) \cdot e^{-|y|} \cdot \frac{[1 + e^{-2|y|}, 1 + 2e^{-2|y|}]}{1 + [0, 1]/\sinh^2(y)}$$

Für  $1/\sinh(y)$  und  $|y| > 372130555$  beweist man elementar die folgenden Ungleichungen

$$\begin{aligned} \frac{1}{\sinh(y)} &= \frac{4}{e^{2|y|} - 2 + e^{-2|y|}} < \frac{4}{e^{2|y|} - 2} < 5 \cdot e^{-2|y|} \\ \frac{1}{\sinh(y)} &= \frac{4}{e^{2|y|} - 2 + e^{-2|y|}} > \frac{4}{e^{2|y|} + e^{-2|y|}} > 3 \cdot e^{-2|y|} \quad \text{und erhält mit (A.60)} \end{aligned}$$

$$(A.61) \quad \begin{aligned} f(x, y) &\in -2 \cdot \text{sign}(y) \cdot \cos([x]) \cdot e^{-|y|} \cdot \frac{[1 + e^{-2|y|}, 1 + 2e^{-2|y|}]}{[1, 1 + 5 \cdot e^{-2|y|}]} \\ &= -2 \cdot \text{sign}(y) \cdot \cos([x]) \cdot e^{-|y|} \cdot \left[ \frac{1 + e^{-2|y|}}{1 + 5e^{-2|y|}}, 1 + 2e^{-2|y|} \right] \\ &\subset -2 \cdot \text{sign}(y) \cdot \cos([x]) \cdot e^{-|y|} \cdot [1 - 4e^{-2|y|}, 1 + 2e^{-2|y|}]. \end{aligned}$$

Zur weiteren Vereinfachung verlangen wir:  $[1 - 4e^{-2|y|}, 1 + 2e^{-2|y|}] \subseteq [\text{pred}(1), \text{succ}(1)]$ , d.h.

$$(A.62) \quad 1 - 4e^{-2y_0} > 1 - 2^{-\text{prec}}, \quad 1 + 2e^{-2y_0} < 1 + 2^{1-\text{prec}}, \quad \text{d.h. wir verlangen}$$

$$(A.63) \quad \text{prec} < \frac{2y_0}{\ln(2)} - 2 = 1\,073\,741\,810,523\dots; \quad \text{prec} < \frac{2y_0}{\ln(2)} = 1\,073\,741\,812,523\dots,$$

und beide Ungleichungen in (A.63) sind schon aus Laufzeitgründen sicher erfüllt. d.h es gilt

$$(A.64) \quad f(x, y) \in -\text{sign}(y) \cdot \cos([x]) \cdot e^{-|y|+\ln(2)} \cdot [\text{pred}(1), \text{succ}(1)].$$

Damit kann  $f(x, y)$  für  $|y| > 372130555$  ohne vorzeitigen Unter- oder Überlauf nahezu optimal und sehr effektiv eingeschlossen werden.

## A.11. $\sqrt{x^2 - 1}$

Da  $f(x) = \sqrt{x^2 - 1}$  für  $|x| \gg 1$  einen vorzeitigen Überlauf erzeugt, benutzen wir in diesem Fall für das gerichtete Auf- und Abrunden die Abschätzungen:

$$(A.65) \quad D(x) := |x| - \frac{0.5}{|x| - \frac{1}{|x|}} < \sqrt{x^2 - 1} < |x| - \frac{0.5}{|x|} =: U(x), \quad |x| \gg 1;$$

Für  $|x| \gg 1$  sind alle drei Terme in (A.65) positiv, so dass die Beweise der zwei Ungleichungen nach dem Quadrieren mit einfachen Umformungen leicht durchgeführt werden können.

Wenn also ein aufgerundeter Funktionswert  $f_u(x) \geq f(x)$  zu berechnen ist, so muss  $U(x)$  aufgerundet werden. Nach (A.8) muss dazu  $0.5/|x|$  abgerundet werden, was mithilfe der MPFR-Funktion `mpfr_div (... , ..., ..., RoundDown)` direkt realisiert werden kann, und die aufzurundende Differenz selbst wird berechnet mit `mpfr_sub (... , ..., ..., RoundUp)`.

Wenn jedoch ein abgerundeter Funktionswert  $f_d(x) \leq f(x)$  zu berechnen ist, so muss  $D(x)$  abgerundet und damit der Doppelbruch aufgerundet werden. Dazu muss  $N := |x| - 1/|x|$  abgerundet werden, wobei  $1/|x|$  mithilfe von `mpfr_div (... , ..., ..., RoundUp)` aufzurunden ist. Die aufzurundende Division  $0.5 \oslash_u N_d$  kann dann nach (A.23) problemlos durchgeführt werden, da der Zähler  $0.5$  rundungsfehlerfrei vorliegt und daher die geforderte Eigenschaft (A.9) automatisch erfüllt ist.

Man kann jetzt noch die Frage stellen, ob die Abschätzungen in (A.65) nicht zu grob sind, so dass bei großen Präzisionen die Fehler  $U(x) - \sqrt{x^2 - 1}$  bzw.  $\sqrt{x^2 - 1} - D(x)$  zu groß ausfallen. Da beide Fehler von der Größenordnung  $O(1/x)^3$  sind und weil (A.65) erst zur Anwendung kommt, wenn `expo(x) > 1073741813` erfüllt ist, was etwa  $323228493$  Dezimalstellen entspricht, so wird sich der Fehler erst nach insgesamt  $323228493 + 3 \cdot 323228493 = 1292913973$  Dezimalstellen bemerkbar machen, was in der numerischen Praxis absolut keine Rolle spielt.

Wenn ein vorzeitiger Überlauf nicht eintreten kann, wird  $x^2 - 1$  nach Bedarf auf- bzw. abgerundet, und weil die Wurzelfunktion eine monoton steigende Funktion ist, können  $f_u(x)$  oder  $f_d(x)$  problemlos berechnet werden.

### A.12. $x/\sqrt{1+x^2}$

Da die monoton wachsende Funktion  $f(x) := x/\sqrt{1+x^2}$  punktsymmetrisch ist, kann man sich bei der Auswertung auf  $x > 0$  beschränken. In diesem Bereich gelten mit  $m := \text{minfloat}()$  die Abschätzungen

$$x \cdot (1 - m) < x \cdot (1 - x^2/2) < \frac{x}{\sqrt{1+x^2}} \equiv \frac{1}{\sqrt{1+(1/x)^2}} < x, \quad x > 0.$$

Die erste Abschätzung  $x \cdot (1 - m) < x \cdot (1 - x^2/2)$  gilt für  $0 < x < \sqrt{2m}$  und wird beim Abrunden im Fall  $x \leq 4m$  benutzt. Falls der aufgerundete Funktionswert für  $x \ll 1$  größer ist als  $x$ , wird  $x$  selbst als Obergrenze gewählt. Die Auswertung von  $f(x)$  erfolgt mit Hilfe von  $f(x) = 1/\sqrt{1+(1/x)^2}$ , da die Wurzel jetzt für  $x \rightarrow \text{MaxFloat}()$  keinen vorzeitigen Überlauf verursacht.

Die nahezu optimale Auswertung von  $f(x)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfrClass xdsqrt1px2 (const MpfrClass& x, RoundingMode rnd)
```

wobei die Rundung mittels  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  gesteuert wird. Falls  $\text{rnd}$  nicht gesetzt wird, erfolgt die Rundung nach dem Current-Rundungsmodus.

### A.13. $x/\sqrt{1-x^2}$

Da die monoton wachsende Funktion  $f(x) := x/\sqrt{1-x^2}$  punktsymmetrisch ist, kann man sich bei der Auswertung auf  $0 < x < +1$  beschränken. Die nahezu optimale Auswertung von  $f(x)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfrClass xdsqrt1mx2 (const MpfrClass& x, RoundingMode rnd)
```

wobei die Rundung mittels  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  gesteuert wird. Falls  $\text{rnd}$  nicht gesetzt wird, erfolgt die Rundung nach dem Current-Rundungsmodus. Weitere Einzelheiten findet man in der Datei `mpfrclass.cpp`.

### A.14. $x/\sqrt{x^2-1}$

Da die monoton fallende Funktion  $f(x) := x/\sqrt{x^2-1}$  punktsymmetrisch ist, kann man sich bei der Auswertung auf  $x > +1$  beschränken. Die nahezu optimale Auswertung von  $f(x)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfrClass xdsqrtx2m1 (const MpfrClass& x, RoundingMode rnd)
```

wobei die Rundung mittels  $\text{rnd} \in \{\text{RoundUp}, \text{RoundDown}, \text{RoundNearest}\}$  gesteuert wird. Falls  $\text{rnd}$  nicht gesetzt wird, erfolgt die Rundung nach dem Current-Rundungsmodus. Weitere Einzelheiten findet man in der Datei `mpfrclass.cpp`.

## A.15. $\ln(\sin(x))$

Da  $f(x) := \ln(\sin(x))$  für  $x \in \mathbb{R}$  nur für  $\sin(x) > 0$  definiert ist, muss  $\sin(x) \leq 0$  bei der Auswertung von  $f(x)$  ausgeschlossen werden. Dies erreicht man durch

```
mpfr_sin(res.mpfr_rep, x.mpfr_rep, RoundFromZero);
// res: sin(x), gerundet weg von der Null;
if (res <= 0)
// ==> sin(x) <= 0:
{
    set_nan(res);
    return res;
}
// Fuer den exakten Funktionswert gilt jetzt: sin(x) > 0:
```

Da  $\sin(x)$  für  $x = \pi/2$  eine waagerechte Tangente besitzt, kann für  $x_0 \approx \pi/2$  der exakte Wert  $f(x_0)$  mithilfe des Funktionsterms  $\ln(\sin(x))$  nur grob approximiert werden. So erhält man z.B.  $f_n(x_0) = 0$ , obwohl  $f(x_0) \neq 0$  ist<sup>5</sup>. Benutzt man daher für  $\sin(x) \geq 0.5$  den Funktionsterm

$$(A.66) \quad f(x) = 0.5 \cdot \ln(1 - \cos^2(x)), \quad \sin(x) \geq 0.5,$$

so wird man nach Abb. A.1 eine sehr viel bessere Auswertung erwarten können, wenn man mit dem Argument  $-\cos^2(x)$  die rechte Seite von (A.66) mithilfe der Funktion  $\lnp1(\dots)$  auswertet.

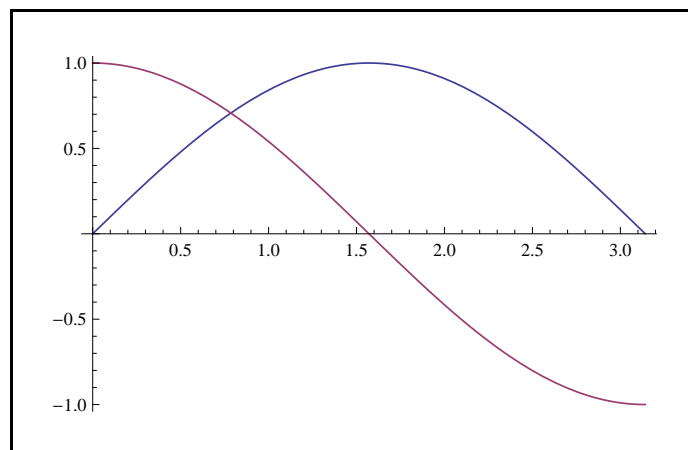


Abbildung A.1.:  $y = \sin(x), \cos(x)$

Für  $x_0 \approx \pi/2$  kann jetzt das Argument  $-\cos^2(x_0) \neq 0$ , im Gegensatz zu  $\sin(x_0)$ , sehr viel effektiver ausgewertet werden, da die Tangente der  $\cos$ -Funktion bei  $x = \pi/2$  die Steigung  $-1$  besitzt.

Jetzt wird gezeigt, wie man mit (A.66) einen garantiert **abgerundeten** Wert  $f_d(x)$  berechnet. Wegen der monoton wachsenden Funktion  $\ln(1+x)$  muss das Argument  $-\cos^2(x)$  abgerundet, d.h.  $\cos^2(x)$  muss aufgerundet werden. Für die folgende Fallunterscheidung muss  $\cos(x) < 0$  und  $\cos(x) > 0$  unterschieden werden, und dies gelingt, wenn man vorher wie oben

```
mpfr_cos(res.mpfr_rep, x.mpfr_rep, RoundFromZero);
```

aufruft und dann bez. des von Null weggerundeten Funktionswertes ( $\text{res} < 0$ ) bzw. ( $\text{res} > 0$ ) abfragt. Aus ( $\text{res} < 0$ ) folgt dann für den exakten Funktionswert die Aussage  $\cos(x) < 0$  und aus ( $\text{res} > 0$ ) folgt entsprechend  $\cos(x) > 0$ <sup>6</sup>.

<sup>5</sup> $f_n(x_0)$  bedeutet den zur nächsten Rasterzahl gerundeten Funktionswert.

<sup>6</sup>Diese Aussagen gelten auch dann, wenn die betrachtete Funktion, hier  $\cos(x)$ , mit `RoundFromZero` viel größer von der Null weggerundet wird. Wenn der exakte Funktionswert jedoch verschwindet, so muss dies auch für den auf- oder abgerundeten Funktionswert zutreffen, was bei allen implementierten Funktionen realisiert ist!

Sei also:  $\cos(x) < 0$ .

Dann gilt für den abgerundeten Funktionswert  $\cos_d(x)$

$$(A.67) \quad \cos_d(x) \leq \cos(x) < 0,$$

und mit (A.18) folgt dann

$$(A.68) \quad \cos_d(x) \odot_u \cos_d(x) \geq \cos^2(x),$$

wobei wegen (A.67) die Bedingung (A.18) automatisch erfüllt ist und auch  $\cos_d(x)$  nicht verschwinden kann.

Sei jetzt:  $\cos(x) > 0$ .

Dann gilt für den aufgerundeten Funktionswert  $\cos_u(x)$

$$(A.69) \quad \cos_u(x) \geq \cos(x) > 0,$$

und mit (A.15) folgt dann

$$(A.70) \quad \cos_u(x) \odot_u \cos_u(x) \geq \cos^2(x),$$

wobei wegen (A.69) die Bedingung (A.9) automatisch erfüllt ist und auch  $\cos_u(x)$  nicht verschwinden kann.

Jetzt wird gezeigt, wie man mit (A.66) einen garantiert **aufgerundeten** Wert  $f_u(x)$  berechnet, dabei werden die Bedingungen  $\cos(x) > 0$  bzw.  $\cos(x) < 0$  ganz analog zur Berechnung von  $f_d(x)$  abgefragt. Der Ausdruck  $\cos^2(x)$  muss jetzt abgerundet werden.

Sei also:  $\cos(x) < 0$ .

Dann gilt für den aufgerundeten Funktionswert  $\cos_u(x)$

$$\cos(x) \leq \cos_u(x) \leq 0,$$

und mit (A.14) folgt dann die gewünschte garantierte Abrundung:

$$\cos_u(x) \odot_d \cos_u(x) \leq \cos^2(x).$$

Sei jetzt:  $\cos(x) > 0$ .

Dann gilt für den abgerundeten Funktionswert  $\cos_d(x)$

$$0 \leq \cos_d(x) \leq \cos(x),$$

und mit (A.11) folgt dann wieder die gewünschte garantierte Abrundung:

$$\cos_d(x) \odot_d \cos_d(x) \leq \cos^2(x).$$

Im unkritischen Bereich  $0 < \sin(x) < 0.5$  wird  $f(x) := \ln(\sin(x))$  direkt ausgewertet, wobei wegen der Monotonie der  $\ln$ -Funktion zur Berechnung von  $f_u(x)$  bzw.  $f_d(x)$  auch  $\sin(x)$  entsprechend auf- bzw. abzurunden ist. Dabei ist jedoch der Bereich  $0 < x < 5 * \text{minfloat}()$  gesondert zu behandeln, um beim Abrunden zu große Überschätzungen zu vermeiden. So erhält man beispielsweise für  $x = \text{minfloat}()$  den abgerundeten Funktionswert  $f_d(x) = -\text{Inf}$ , der jedoch eine viel zu 'kleine' Unterschranke darstellt. Um eine größere und optimale Unterschranke zu berechnen, betrachten wir die Reihe:

$$\begin{aligned} \ln(\sin(x)) &= \ln(x) - \frac{x^2}{6} - \frac{x^4}{180} - \dots = - \sum_{k=1}^{\infty} \frac{2^{2k-1}}{k \cdot (2k)!} |B_{2k}| x^{2k}, \quad 0 < x < \pi, \\ &> \ln(x) - \frac{x^2}{6} (1 + x^2 + x^4 + x^6 + \dots) = \ln(x) - \frac{x^2}{6} \cdot \frac{1}{1-x^2}, \quad 0 < x < 1, \\ &> \ln(x) - \frac{x}{6} \cdot \frac{1}{1-x^2} > \ln(x) - \frac{x}{5} > \ln(x) - \text{minfloat}(), \text{ falls } x < 5 \cdot \text{minfloat}(). \end{aligned}$$

Für  $x < 5 \cdot \text{minfloat}()$  erhalten wir damit die Unterschranke:  $\ln(x) - \text{minfloat}() < \ln(\sin(x))$ .

Die Funktion  $f(x) = \ln(\cos(x))$  ist in `mpfrclass.cpp` ganz analog implementiert und muss deshalb hier nicht weiter diskutiert werden.



### A.16. $\ln(\sqrt{x^2 + y^2})$

Wir betrachten

$$(A.71) \quad f(x) := \ln(\sqrt{x^2 + y^2}), \quad \text{mit: } |x| \geq |y| \text{ und } |x|, |y| > 0.$$

Um im Fall  $|x| \gg 1$  einen vorzeitigen Überlauf zu vermeiden, wird folgende Darstellung benutzt:

$$(A.72) \quad f(x) := \ln(|x|) + \frac{1}{2} \ln(1 + (|y|/|x|)^2) > 0, \quad |y|/|x| > 0, |x| \gg 1.$$

Zur Berechnung von  $f_u(x)$  bzw.  $f_d(x)$  muss  $\ln(|x|)$  auf- bzw. abgerundet werden, was mithilfe der MPFR-Funktion `mpfr_log(...)` direkt realisiert werden kann. Da auch der zweite Summand in (A.72) auf- bzw. abzurunden ist und die `lnp1`-Funktion monoton wächst, muss das Quadrat  $(|y|/|x|)^2$  ebenfalls auf- bzw. abgerundet werden und damit auch der Quotient  $A := |y|/|x|$ . Da die Operanden von  $A$  rundungsfehlerfreie Rasterzahlen sind, können  $A_u$  und  $A_d$  direkt mithilfe der MPFR-Funktion `mpfr_div(...)` berechnet werden. Es gilt dann:

$$(A.73) \quad A_u \geq A > 0$$

$$(A.74) \quad 0 \leq A_d < A.$$

Es soll jetzt  $(|y|/|x|)^2$  aufgerundet werden. Mit (A.15) folgt dann

$$A_u \odot_u A_u \geq (|y|/|x|)^2,$$

wobei (A.9) wegen (A.73) automatisch erfüllt ist und  $A_u > 0$  garantiert ist.

Es soll jetzt  $(|y|/|x|)^2$  abgerundet werden. Mit (A.11) folgt dann

$$A_d \odot_d A_d \leq (|y|/|x|)^2,$$

wobei  $A_d \geq 0$  durch (A.74) garantiert wird.

Wenn  $f_u(x)$  und  $f_d(x)$  nach (A.71) zu berechnen sind, so ist dies kein Problem, da  $\sqrt{x^2 + y^2}$  mithilfe der Funktion `sqrtox2y2(..., rnd)` direkt auf- oder abgerundet werden kann und weil die  $\ln$ -Funktion monoton wachsend ist.

### A.17. $\ln(\sqrt{(1+x)^2 + y^2})$

Da  $u(x, y) = \ln(\sqrt{(1+x)^2 + y^2})$  auch als Realteil der komplexen Funktion  $\ln(1+z)$ , mit  $z \in \mathbb{C}$ , benötigt wird, findet man eine ausführliche Beschreibung des Algorithmus ab Seite 267.

## A.18. $\operatorname{arcosh}(1+x)$

Wir betrachten für  $x \geq 0$

$$(A.75) \quad \operatorname{arcosh}(1+x) = \ln(1+x + \sqrt{x \cdot (2+x)}), \quad x \geq 0,$$

$$(A.76) \quad = \ln(x) + \ln(1 + 1/x + \sqrt{1 + 2/x}), \quad x \gg 1;$$

In (A.76) ist die Rundung auch des zweiten Summanden unproblematisch, da bei den zwei auftretenden Divisionen beide Operanden jeweils rundungsfehlerfreie Rasterzahlen sind und neben der  $\ln$ -Funktion auch die Quadratwurzel eine monoton wachsende Funktion ist.

In (A.75) muss bei der Berechnung von  $f_u(x)$  und  $f_d(x)$  das Produkt  $P := x \cdot (2+x) = A \cdot B$  auf- bzw. abgerundet werden, wobei  $B := 2+x \geq 2$  ebenfalls auf- bzw. abzurunden ist. Es gilt:

$$(A.77) \quad 0 \leq A = A_u = A_d,$$

$$(A.78) \quad 2 \leq B_d \leq B,$$

$$(A.79) \quad 2 \leq B \leq B_u;$$

Nach (A.15) gilt:  $P_u := A_u \odot_u B_u \geq x \cdot (2+x)$ , wobei die Bedingung (A.9) wegen (A.77) automatisch erfüllt ist und auch  $A_u = A$  und  $B_u \geq 2$  nicht beide verschwinden können.

Nach (A.11) gilt wegen (A.77) und (A.78) unmittelbar:  $P_d := A_d \odot_d B_d \leq x \cdot (2+x)$ , womit für  $P_d$  und  $P_u$  die entsprechenden Rundungen garantiert sind. Den vollständigen Algorithmus findet man in `mpfrclass.cpp`.

## A.19. $\Gamma'(x)$

$\Gamma'(x)$  wird mithilfe der Digamma-Funktion  $\psi(x) = \Gamma'(x)/\Gamma(x)$  über das folgende Produkt berechnet:

$$(A.80) \quad \Gamma'(x) = \psi(x) \cdot \Gamma(x), \quad x \neq 0, -1, -2, \dots$$

Die Funktionen  $\psi(x)$  und  $\Gamma(x)$  stehen dabei als `digamma(x)` bzw. `gamma(x)` zur Verfügung. Um das Produkt  $\psi(x) \cdot \Gamma(x)$  gesichert auf- bzw. abzurunden, werden  $\psi(x)$  und  $\Gamma(x)$  jeweils durch die Intervalle `u` und `v` eingeschlossen. Mit dem Intervallprodukt  $u = u \diamond v$  ist dann das auf- bzw. abgerundete Produkt  $\psi(x) \cdot \Gamma(x)$  gegeben durch `Sup(u)` bzw. `Inf(u)`, und der nächstgelegene Wert ist gegeben durch `mid(u)`. Der Vorteil dieser Methode ist einmal die kurze Laufzeit, da zur Berechnung von `v` nur ein Funktionsaufruf `Inf(v) = gamma(x, RoundDown)` notwendig ist, wobei dann `Sup(v)` einfach durch den Nachfolger `succ(Inf(v)) = Sup(x)` zu berechnen ist. Die geschilderte Intervallmethode hat zusätzlich den Vorteil, dass die sonst notwendigen Fallunterscheidungen jetzt entfallen, wodurch ein übersichtlicher Programmcode gewährleistet wird. Aber jedes Verfahren besitzt auch einen Nachteil, denn wenn in seltenen Ausnahmefällen für `Inf(v) := gamma(x, RoundDown)` gilt: `Inf(v) = \Gamma(x)`, so wird mit `Sup(x) = succ(Inf(v))` das Intervall `v` leicht nach oben überschätzt, da bei diesem Beispiel das ideale Intervall `v` ein Punktintervall ist. Den vollständigen Algorithmus findet man in `mpfrclass.cpp`.

## A.20. $1/\Gamma(x)$

$f(x) := 1/\Gamma(x)$  ist in der ganzen komplexen Ebene analytisch und besitzt an den Polstellen  $x_k = 0, -1, -2, \dots$  der  $\Gamma$ -Funktion die Funktionswerte Null. Für alle anderen reellen  $x$ -Werte wird  $f(x)$  approximiert durch den Quotienten  $1/\Gamma(x)$  der auf der Maschine mittels `gamma(x)` ausgewertet wird.

Um  $f_u(x)$  zu berechnen, ist nach (A.23) und (A.25)  $\Gamma(x)$  unabhängig von seinem Wert stets abzurunden und die Division muss aufgerundet werden.

Um  $f_d(x)$  zu berechnen, ist nach (A.19) und (A.21)  $\Gamma(x)$  unabhängig von seinem Wert stets aufzurunden und die Division muss abgerundet werden.

## A.21. $(1/\Gamma(x))'$

Für die überall differenzierbare Funktion  $f(x) := (1/\Gamma(x))'$  gilt:

$$(A.81) \quad f(x) = \left(\frac{1}{\Gamma(x)}\right)' = \begin{cases} \frac{-\Gamma'(x)}{\Gamma^2(x)} = -\frac{\Gamma'(x)}{\Gamma(x)}/\Gamma(x), & x \neq 0, -1, -2, \dots \\ |x|!, & x = 0, -2, -4, \dots \\ -|x|!, & x = -1, -3, -5, \dots \end{cases}$$

$|x|!$  wird mithilfe der `faktorial(...)`-Funktion berechnet und die erste Zeile in (A.81) mit dem Quotienten `-digamma(x)/gamma(x)`, wobei zur Vermeidung von Fallunterscheidungen Zähler und Nenner und auch die Division selbst intervallmäßig ausgewertet werden. In der Umgebung von Null, d.h. für  $|x| \leq 2 \cdot \text{succ}(0)$  und  $x \neq 0$ , entsteht jedoch das Problem, dass `digamma(x)` und `gamma(x)` beide einen Überlauf erzeugen, obwohl ihr Quotient von der Größenordnung 1 ist. Um diesen Überlauf zu vermeiden, wird ausgenutzt, dass  $f'(x)$  für  $x \in [-1, 0.3]$  **positiv** ist, so dass  $f(x)$  in obiger Umgebung von Null monoton wächst. Für  $x = \text{succ}(0)$  und  $x = 2 \cdot \text{succ}(0)$  ist damit wegen  $f(0) = 1$  eine Einschließung von  $f(x)$  gegeben durch:

$$\begin{aligned} t &= 3 \cdot \text{succ}(0), & u &= \text{digamma}(t, \text{RoundDown}), & v &= \text{gamma}(t, \text{RoundDown}), \\ & & & u &= -u \diamond v; \\ f(x) &\in [1, \text{Sup}(u)], & x &= \text{succ}(0) \text{ oder } x = 2 \cdot \text{succ}(0) \end{aligned}$$

wobei wegen  $t = 3 \cdot \text{succ}(0)$  die Punktintervalle  $u$  und  $v$  jetzt keinen Überlauf mehr erzeugen. Natürlich ist das Einschließungsintervall  $[1, \text{Sup}(u)]$  eine gewisse Überschätzung von  $f(x)$ , aber dies wird in der numerischen Praxis kaum eine Rolle spielen.

Für die negativen Argumente  $x = \text{pred}(0)$  und  $x = 2 \cdot \text{pred}(0)$  gelten ganz entsprechende Aussagen. Den ausführlichen Algorithmus findet man in `mpfrclass.cpp`.

## B. Neue Funktionen vom Typ MpfiClass

### B.1. $\ln(\sqrt{(1+x)^2 + y^2})$

Für die vorgegebenen reellen Maschinenintervalle  $X \ni x$  und  $Y \ni y$  liefert die Intervallfunktion

```
MpfiClass ln_sqrtxp1_2y2(const MpfiClass& X, const MpfiClass& Y)
```

eine fast optimale Einschließung aller Funktionswerte

$$u(x, y) = \ln(\sqrt{(1+x)^2 + y^2}) = \frac{1}{2} \cdot \ln((1+x)^2 + y^2), \quad \text{mit: } x \in X, y \in Y.$$

Eine ausführliche Beschreibung von  $u(x, y)$  findet man ab Seite 267.  $u(x, y)$  ist als Realteilmfunktion der holomorphen Funktion  $\log(1+z)$ ,  $z \in \mathbb{C}$ , eine harmonische Funktion, wenn der singuläre Punkt  $P(-1, 0)$  nicht ein Element des zweidimensionalen reellen Intervalls  $Z$  ist.

$$(B.1) \quad Z := \{(x, y) \mid x \in X \wedge y \in Y\}, \quad (-1, 0) \notin Z.$$

$u(x, y)$  nimmt dann für alle  $(x, y) \in Z$  als harmonische Funktion ihre relativen Extrema auf dem Rand von  $Z$  an, vgl. dazu auch Seite 244. Die Berechnung dieser Extremstellen,  $m(x, y)$  für das relative Minimum und  $M(x, y)$  für das relative Maximum, kann wesentlich vereinfacht werden, wenn man  $u(x, y) \equiv u(x, |y|)$  beachtet und daher das Eingangsintervall  $Y$  durch das Intervall  $Y = \text{abs}(Y) = \{|y| \mid y \in Y\}$  seiner Absolutbeträge ersetzt. Die zu betrachtenden Intervalle  $Z$  aus (B.1) liegen damit alle in der oberen Halbebene, und für die interne Berechnung gilt mit  $X = X$ :

$$(B.2) \quad Z := \{(x, y) \mid x \in X = [x_1, x_2] \wedge y \in Y = [y_1, y_2]\}, \quad (-1, 0) \notin Z, \quad y_1 \geq 0.$$

Um die Koordinaten der Extrempunkte  $m, M$  bestimmen zu können, berechnen wir zunächst die partiellen Ableitungen von  $u(x, y)$  bezüglich  $x, y$ :

$$\begin{aligned} \frac{\partial u(x, y)}{\partial x} &= \frac{1+x}{(1+x)^2 + y^2} = 0, & \text{d.h. } x &= -1, \\ \frac{\partial u(x, y)}{\partial y} &= \frac{y}{(1+x)^2 + y^2} = 0, & \text{d.h. } y &= 0. \end{aligned}$$

Auf einer Parallelen zur  $x$ -Achse wächst bzw. fällt  $u(x, y)$  monoton für  $x > -1$  bzw. für  $x < -1$ , und auf einer Parallelen zur  $y$ -Achse ist  $u(x, y)$  (in der oberen Halbebene) monoton wachsend. Die Extremalkurven sind wegen  $x = -1$  einmal die Parallele zur  $y$ -Achse durch  $x = -1$  und zum anderen wegen  $y = 0$  die  $x$ -Achse, wobei letztere keine Rolle spielt, da wir nur Intervalle  $Z$  in der oberen Halbebene betrachten. Schneidet die Extremalkurve  $x = -1$  das Rechteck  $Z$ , so besitzt  $m$  die Koordinaten  $m(-1, y_1)$ . In Abb. B.1 geben die Pfeile auf dem Rand von  $Z$  die Richtung wachsender  $u(x, y)$ -Werte an, so dass die Extrempunkte  $m$  und  $M$  die folgenden Koordinaten besitzen:  $m(x_1, y_1)$ ,  $M(x_2, y_2)$ . Die breiten Pfeile links und rechts neben der **Extremalkurve**  $x = -1$  geben das Monotonieverhalten der Funktion  $u(x, y)$  an. Das zweidimensionale, **grüne** Intervall ist nicht erlaubt, da es den singulären Punkt  $P(-1, 0)$  enthält. Die Koordinaten der Extrempunkte  $m, M$  der erlaubten gelben, zweidimensionalen Intervalle  $Z$  sind stets Maschinenzahlen, so dass die Berechnung der Funktionswerte  $u(x, y)$  an diesen Stellen einfach realisiert werden kann.

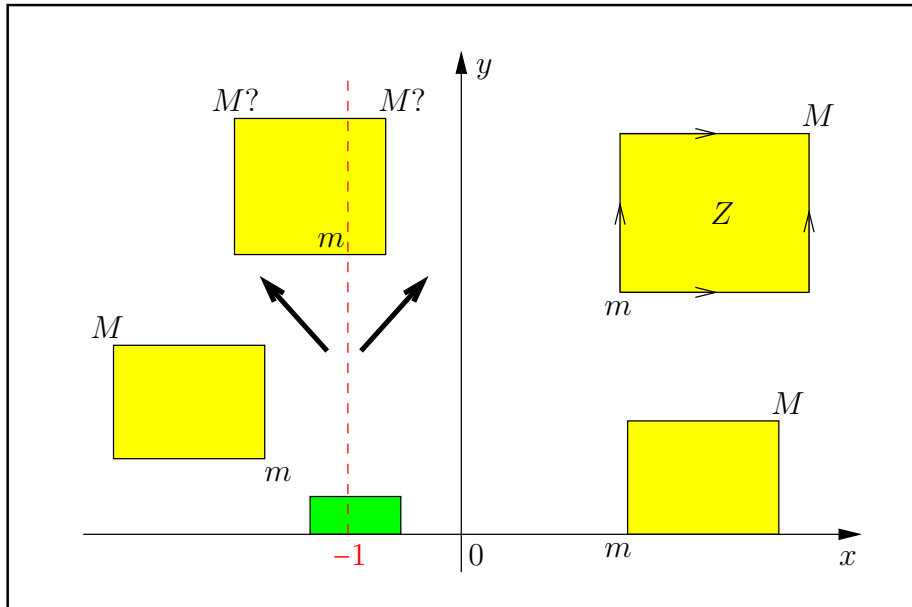


Abbildung B.1.: Zweidimensionale Intervalle  $Z$  mit Extrempunkten  $m, M$

Bei einem Intervall  $Z$ , das die **Extremalkurve**  $x = -1$  schneidet, muss noch entschieden werden, welcher der beiden möglichen Punkte  $M?$  das Maximum von  $u(x, y)$  wirklich liefert. Nach der obigen Abbildung ergibt sich

$$|1 + x_1| \geq |1 + x_2| \iff -1 - x_1 \geq x_2 - (-1) \iff x_1 + 2 \leq -x_2,$$

so dass nach obiger Abbildung für den Maximumpunkt gilt:  $M(x_1, y_2)$ . Die letzte Ungleichung muss auf der Maschine ausgewertet werden. Dazu berechnet man zunächst das reelle Intervall  $z = x_1 \diamond 2 = [z_1, z_2]$ , das wegen  $x_1 \leq -1$  keinen Überlauf erzeugt. Dabei ist  $z$  entweder ein Punktintervall oder in seinem Innern befindet sich keine Maschinenzahl, d.h. es gilt  $z_2 = z_1$  oder  $z_2 = \text{succ}(z_1)$ . Aus der Doppelungleichung  $x_1 + 2 \leq z_2 \leq -x_2$ , d.h. aus  $z_2 \leq -x_2$  folgt dann  $x_1 + 2 \leq -x_2$  und damit  $M(x_1, y_2)$ . Ganz entsprechend folgt aus  $z_1 \geq -x_2$  für den Maximumpunkt  $M(x_2, y_2)$ .

Nach diesen Überlegungen sind die **beiden** Abfragen  $z_2 \leq -x_2$  und  $z_1 \geq -x_2$  nötig, um die richtige  $x$ -Koordinate von  $M$  zu bestimmen. Wir zeigen jetzt, dass dazu **nur eine** Abfrage, z.B.  $z_2 \leq -x_2$ , notwendig ist und dabei nur der rechte Randpunkt  $z_2 \geq x_1 + 2$  bekannt sein muss. In der folgenden Abbildung findet man auf der oberen  $x$ -Achse die Situation, wenn  $z_2 \leq -x_2$  und  $z_2 = \text{succ}(z_1)$  erfüllt sind. Im grünen Bereich liegen dort alle möglichen  $-x_2$ -Werte, so dass die

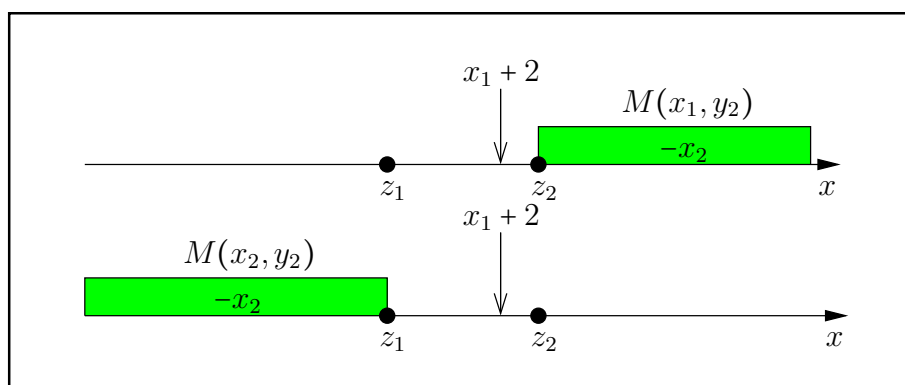


Abbildung B.2.:  $x$ -Koordinate von  $M(?, y_2)$  im Fall:  $-1 \in [x_1, x_2]$

Bedingung  $x_1 + 2 \leq -x_2$  stets erfüllt ist und damit  $M$  gegeben ist durch  $M(x_1, y_2)$ . Wenn jedoch  $z_2 \leq -x_2$  nicht erfüllt ist, so gilt  $z_2 > -x_2$  und da zwischen  $z_1$  und  $z_2$  keine Maschinenzahl liegt, können die möglichen  $-x_2$ -Werte jetzt nur in dem grünen Bereich auf der unteren  $x$ -Achse liegen, d.h. es gilt stets  $x_1 + 2 \geq -x_2$ , so dass jetzt  $M$  gegeben ist durch  $M(x_2, y_2)$ .

### Zusammenfassung:

- $x_1 \geq -1 \implies m = m(x_1, y_1), \quad M = M(x_2, y_2);$
- $x_2 \leq -1 \implies m = m(x_2, y_1), \quad M = M(x_1, y_2);$
- $-1 \in [x_1, x_2], \quad z_2 := x_1 \oplus_u 2 \implies$ 
  1.  $m = m(-1, y_1);$
  2.  $z_2 \leq -x_2 \implies M = M(x_1, y_2); \quad$  Nur diese Abfrage muss realisiert werden!
  3.  $z_2 > -x_2 \implies M = M(x_2, y_2);$

Beachten Sie, dass die Maschinenzahl  $z_2 := x_1 \oplus_u 2$  die optimal aufgerundete Summe  $(x_1 + 2)$  ist.

### Numerische Ergebnisse:

Bezeichnungen:  $t = \ln(\sqrt{(1+x)^2 + y^2}), \quad x \in X = [x_1, x_2], \quad y \in Y = [y_1, y_2];$

Vergleichen Sie für Punktargumente auch die auf- und abgerundeten Ergebnisse von Seite 269.

1. Mit  $X = Y = [\text{MaxFloat}(\text{prec}), \text{MaxFloat}(\text{prec})]$  und  $\text{prec} = 3000$  erhält man für  $t$  die nahezu optimale Einschließung mit 902 gemeinsamen Dezimalstellen:

$$t \in [7.4426111760831942759393\dots664834e8, 7.4426111760831942759393\dots664836e8].$$

2. Mit  $X = Y = [\text{minfloat}(\text{prec}), \text{minfloat}(\text{prec})]$  und  $\text{prec} = 3000000$  erhält man für  $t$  die notwendigerweise grobe Einschließung:

$$t \in [0, 2.3825649048879510732161697817326745\dots34539016124e - 323228497].$$

3. Mit  $X = Y = [2 \cdot \text{minfloat}(\text{prec}), 2 \cdot \text{minfloat}(\text{prec})]$  und  $\text{prec} = 30000$  erhält man für  $t$  die nahezu optimale Einschließung mit 9030 gemeinsamen Dezimalstellen:

$$t \in [4.765129809775\dots95148232e - 323228497, 4.765129809775\dots95148233e - 323228497].$$

4. Mit  $X = [-2, -2], \quad Y = [2^{-10000}, 2^{-10000}]$  und  $\text{prec} = 30000$  erhält man für  $t$  die nahezu optimale Einschließung mit 9030 gemeinsamen Dezimalstellen:

$$t \in [1.2561940288493722925\dots505035e - 6021, 1.2561940288493722925\dots505036e - 6021].$$

Die gleiche Einschließung erhält man natürlich auch mit  $X = [0, 0], \quad Y = [2^{-10000}, 2^{-10000}]$ .

5. Mit  $X = [\text{succ}(-1), \text{succ}(-1)], \quad \text{succ}(-1) = -1 + 2^{-\text{prec}}, \quad Y = [2^{-30000}, 2^{-30000}]$  und der Präzision  $\text{prec} = 30000$  erhält man für  $t$  die nahezu optimale Einschließung mit 9030 gemeinsamen Dezimalstellen:

$$t \in [-2.079406884320807930\dots15528832738e4, -2.079406884320807930\dots15528832737e4].$$

Beachten Sie, dass eine Präzision von 300000 Bits einer Präzision von  $300000/\log_2(10) = 30000/3.32192809\dots \approx 9031$  Dezimalstellen entspricht. Bezüglich  $\text{succ}(-1) = -1 + 2^{-\text{prec}}$  vergleiche Seite 17.

## B.2. $x/(x^2 + y^2)$

Für die vorgegebenen reellen Maschinenintervalle  $X \ni x$  und  $Y \ni y$  liefert die Intervallfunktion

`MpfiClass x_div_x2py2(const MpfiClass& X, const MpfiClass& Y)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x, y) = \frac{x}{x^2 + y^2} \quad \text{mit: } x \in X, y \in Y \text{ und } x \neq 0 \wedge y \neq 0.$$

Die Berechnung von  $f(x, y)$  wird ausführlich beschrieben auf Seite 185.  $f(x, y)$  ist als Realteilstfunktion der holomorphen Funktion  $1/z$ ,  $z \in \mathbb{C} \setminus \{0\}$ , eine harmonische Funktion, wenn der singuläre Punkt  $P(0, 0)$  nicht ein Element des zweidimensionalen reellen Intervalls  $Z$  ist.

$$(B.3) \quad Z := \{(x, y) \mid x \in X \wedge y \in Y\}, \quad 0 \notin X \wedge 0 \notin Y.$$

$f(x, y)$  nimmt daher für alle Punkte  $(x, y) \in Z$  als harmonische Funktion ihre Extrema auf dem Rand von  $Z$  an, vgl. dazu auch die Seite 244. Die Berechnung dieser Extremstellen,  $m(x, y)$  für das Minimum und  $M(x, y)$  für das Maximum, kann wesentlich vereinfacht werden, wenn man  $f(x, y) \equiv f(x, |y|)$  beachtet und daher das Eingangsintervall  $Y$  durch das Intervall  $Y = \text{abs}(Y) = \{|y| \mid y \in Y\}$  seiner Absolutbeträge ersetzt. Die zu betrachtenden Intervalle  $Z$  aus (B.3) liegen damit alle in der oberen Halbebene, und für die interne Berechnung gilt mit  $X = X$ :

$$(B.4) \quad Z := \{(x, y) \mid x \in X = [x_1, x_2] \wedge y \in Y = [y_1, y_2]\}, \quad (0, 0) \notin Z, y_1 \geq 0.$$

Wegen  $f(-x, y) \equiv -f(x, y)$  kann die Berechnung der Extrema  $m, M$  noch weiter auf den ersten Quadranten reduziert werden, wobei das Minimum  $m$  und das Maximum  $M$  berechnet werden mithilfe der Funktionen:

`MpfrClass x_div_m(const MpfiClass& x, const MpfiClass& y);`    `Inf(x)>=0, Inf(y)>=0,`  
`MpfrClass x_div_M(const MpfiClass& x, const MpfiClass& y);`    `Inf(x)>=0, Inf(y)>=0.`

Um die Extrema  $m, M$  mit den obigen Funktionen richtig berechnen zu können, sind drei Fälle zu unterscheiden:

1.  $x_1 \geq 0 \wedge y_1 \geq 0$ :    `m = x_div_m(X, Y); M = x_div_M(X, Y);`
2.  $x_2 < 0 \wedge y_1 \geq 0$ :    `res = -x; m = x_div_m(res, Y); M = x_div_M(res, Y);`  
`swap(m, M); m = -m; M = -M;`
3.  $x_1 < 0 \wedge x_2 \geq 0 \wedge y_1 \geq 0$ :    `res = [0, -x1]; m = -x_div_M(res, Y);`  
`res = [0, x2]; M = x_div_M(res, Y);`

Um jetzt im 1. Quadranten die beiden Funktionen `x_div_m(...)`, `x_div_M(...)` realisieren zu können, benötigen wir zunächst für  $x, y \geq 0$  die beiden partiellen Ableitungen von  $f(x, y)$ :

$$(B.5) \quad \frac{\partial f(x, y)}{\partial x} = \frac{y^2 - x^2}{(x^2 + y^2)^2},$$

$$(B.6) \quad \frac{\partial f(x, y)}{\partial y} = \frac{-2xy}{(x^2 + y^2)^2} \leq 0.$$

Oberhalb bzw. unterhalb der ersten Winkelhalbierenden gilt damit:  $\partial f/\partial x > 0$  bzw.  $\partial f/\partial x < 0$ , und auf einer Parallelen zur  $y$ -Achse nehmen die Funktionswerte von  $f(x, y)$  mit wachsenden  $y$ -Werten monoton ab, und wegen

$$(B.7) \quad f(x_1, x_1) = \frac{x_1}{x_1^2 + x_1^2} \geq \frac{x_2}{x_2^2 + x_2^2} \iff \frac{1}{2x_1} \geq \frac{1}{2x_2} \iff 0 < x_1 \leq x_2$$

nehmen die Funktionswerte von  $f(x, y)$  monoton zu, wenn man sich im 1. Quadranten auf der 1. Winkelhalbierenden dem Ursprung nähert. Mit diesen Ergebnissen erhält man in folgender Abbildung die Lage der Extremstellen  $m, M$  zu gegebenen zweidimensionalen Intervallen  $Z_i$ , wobei fünf grundsätzliche Lagen dieser  $Z_i$  zu unterscheiden sind.



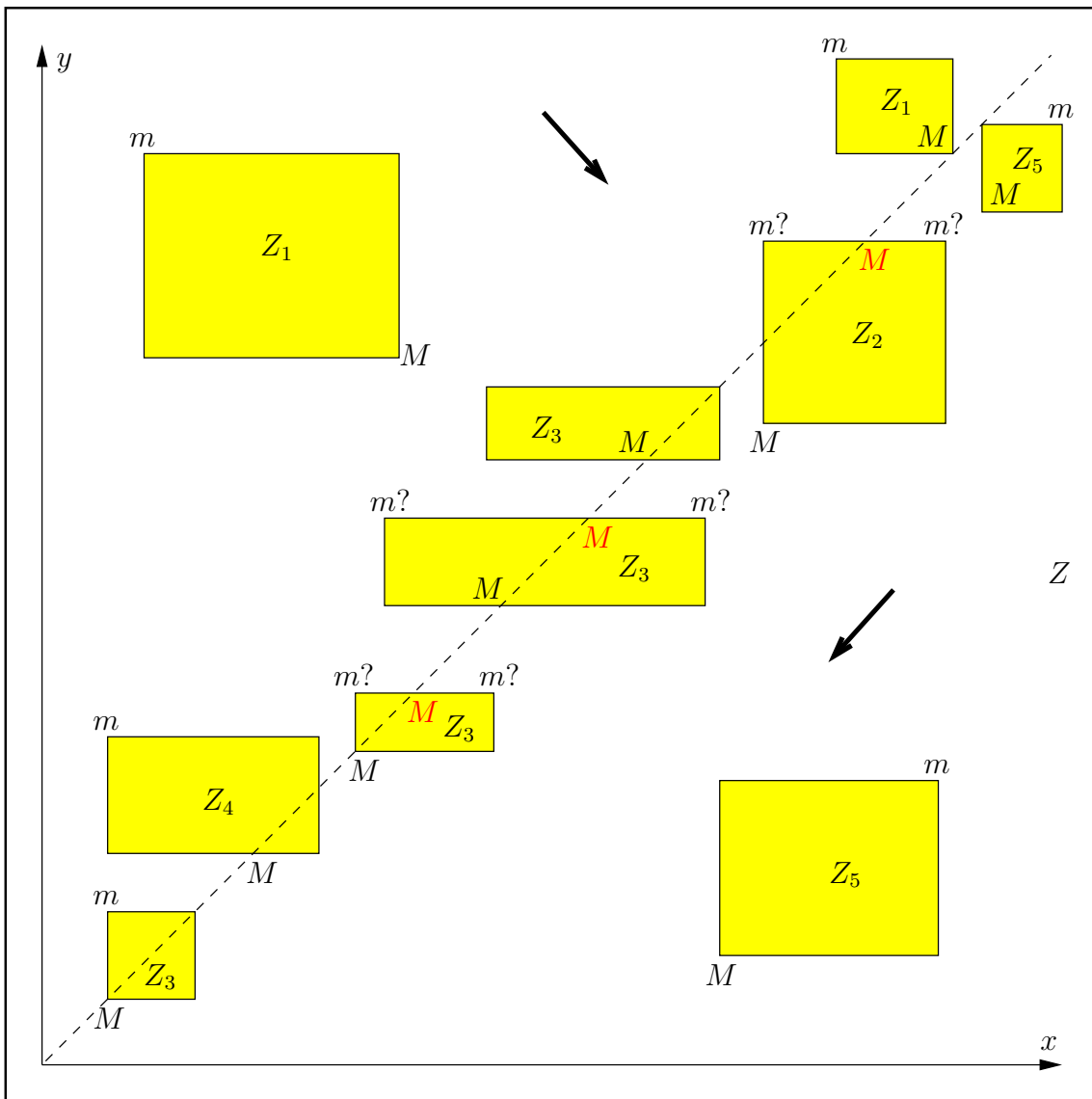


Abbildung B.3.: Zweidimensionale Intervalle  $Z_i$  mit Extrempunkten  $m, M$

Die fünf grundsätzlichen Lagen der  $Z_i$  sind definiert durch:

1.  $x_2 \leq y_1$ ,  $Z_1$ ;
2.  $y_1 < x_1 < y_2$ ,  $Z_2$ ;
3.  $x_1 \leq y_1 \leq x_2 \wedge x_1 \leq y_2 \leq x_2$ ,  $Z_3$ ;
4.  $y_1 < x_2 < y_2$ ,  $Z_4$ ;
5.  $y_2 \leq x_1$ ,  $Z_5$ ;

Die Punkte  $M$  auf der 1. Winkelhalbierenden sind Punkte mit relativen Maxima, wobei jedoch gilt  $M \geq \bar{M}$ . Bei den Punktepaaaren  $m?, m?$  ist jeweils der Punkt mit dem kleinsten Funktionswert auszuwählen. Die beiden Pfeile oberhalb und unterhalb der ersten Winkelhalbierenden geben die Richtung wachsender Funktionswerte von  $f(x, y)$  an.

### B.3. $(x^2 - y^2)/(x^2 + y^2)^2$

Für die vorgegebenen reellen Maschinenintervalle  $X \ni x$  und  $Y \ni y$  liefert die Intervallfunktion

`MpfiClass Re_rz2(const MpfiClass& X, const MpfiClass& Y)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x, y) = \frac{x^2 - y^2}{(x^2 + y^2)^2} \quad \text{mit: } x \in X, y \in Y \quad \text{und} \quad x \neq 0 \wedge y \neq 0.$$

Die Berechnung von  $f(x, y)$  wird ausführlich beschrieben auf Seite 186.  $f(x, y)$  ist als Realteilfunktion der holomorphen Funktion  $1/z^2$ ,  $z \in \mathbb{C} \setminus \{0\}$ , eine harmonische Funktion, wenn der singuläre Punkt  $P(0, 0)$  nicht ein Element des zweidimensionalen reellen Intervalls  $Z$  ist.

$$(B.8) \quad Z := \{(x, y) \mid x \in X \wedge y \in Y\}, \quad 0 \notin X \wedge 0 \notin Y.$$

$f(x, y)$  nimmt daher für alle Punkte  $(x, y) \in Z$  als harmonische Funktion ihre Extrema auf dem Rand von  $Z$  an, vgl. dazu auch die Seite 244. Die Berechnung dieser Extremstellen,  $m(x, y)$  für das Minimum und  $M(x, y)$  für das Maximum, kann wesentlich vereinfacht werden, wenn man  $f(x, y) \equiv f(|x|, |y|)$  beachtet und die Eingangsintervalle  $X, Y$  durch die jeweiligen Intervalle, z.B.  $X = \text{abs}(X) = \{|x| \mid x \in X\}$ , seiner Absolutbeträge ersetzt. Damit reduziert sich die Berechnung der Extremstellen  $m(x, y), M(x, y)$  auf ein Intervall  $Z := \{(x, y) \mid x \in X \wedge y \in Y\}$ ,  $x \geq 0, y \geq 0$ , das nur im 1. Quadranten liegt und den Ursprung nicht enthalten darf.  $X = [x_1, x_2]$ ,  $Y = [y_1, y_2]$ .

Zur Berechnung der Extremstellen  $m, M$  benötigt man zunächst die partiellen Ableitungen:

$$(B.9) \quad \frac{\partial f(x, y)}{\partial x} = \frac{2x \cdot (-x^4 + 2x^2y^2 + 3y^4)}{(x^2 + y^2)^4} = \frac{2x \cdot (4y^4 - (x^2 - y^2)^2)}{(x^2 + y^2)^4},$$

$$(B.10) \quad \frac{\partial f(x, y)}{\partial y} = \frac{-2y \cdot (-y^4 + 2x^2y^2 + 3x^4)}{(x^2 + y^2)^4} = \frac{-2y \cdot (4x^4 - (x^2 - y^2)^2)}{(x^2 + y^2)^4}.$$

Die Extremalkurven für  $\partial f/\partial x$  ergeben sich aus der Forderung  $\partial f/\partial x = 0$ , siehe Seite 247. Die erste Extremalkurve ist wegen  $x = 0$  die positive  $y$ -Achse, und die zweite ergibt sich aus der Forderung  $4y^4 = (x^2 - y^2)^2 \iff 2y^2 = |x^2 - y^2|$ , aus der  $y = x/\sqrt{3}$  folgt.

Die beiden Extremalkurven für  $\partial f/\partial y$  ergeben sich analog aus der Forderung  $\partial f/\partial y = 0$ . Die erste Extremalkurve ist dann wegen  $y = 0$  die positive  $x$ -Achse und die zweite Extremalkurve ist gegeben durch  $y = \sqrt{3} \cdot x$ .

#### Anmerkungen:

1. Da beide partiellen Ableitungen auf ihren Extremalkurven das Vorzeichen wechseln, wenn man parallel zu den beiden Koordinatenachsen diese Extremalkurven schneidet, besitzt  $f(x, y)$  an den Schnittpunkten der Intervallränder von  $Z$  mit den Extremalkurven längs dieser Intervallränder tatsächlich relative Extrema, d.h. in diese Fällen müssen die Extrempunkte  $m$  oder  $M$  nicht auf den Eckpunkten von  $Z$  liegen.
2. Es gilt folgendes:  $x > y \implies f(x, y) > 0$  und  $x < y \implies f(x, y) < 0$ , d.h. wird  $Z$  von der 1. Winkelhalbierenden geschnitten, so kann das Minimum nur auf dem Rand des oberen und das Maximum nur auf dem Rand des unteren  $Z$ -Abschnitts liegen.
3. In den folgenden Abbildungen sind alle typischen Lagen der Intervalle  $Z$  mit den entsprechenden Lagen der Extremstellen  $m(x, y), M(x, y)$  angegeben. Dabei bedeutet das Symbol  $(M)$ , dass hier zwar ein relatives Maximum vorliegt, dass aber das wirkliche Maximum an der anderen Stelle  $M$  angenommen wird. Die zwei Symbole  $m?$  bedeuten, dass das wirkliche Minimum nur an einer dieser beiden Stellen angenommen werden kann, so dass zwei Berechnungen und ein Vergleich erforderlich sind. Die 1. Winkelhalbierende und die beiden Extremalkurven  $y = \sqrt{3} \cdot x$  und  $y = x/\sqrt{3}$  sind leicht verdreht dargestellt. Die farbigen Pfeile geben die Richtung wachsender Funktionswerte von  $f(x, y)$  an.

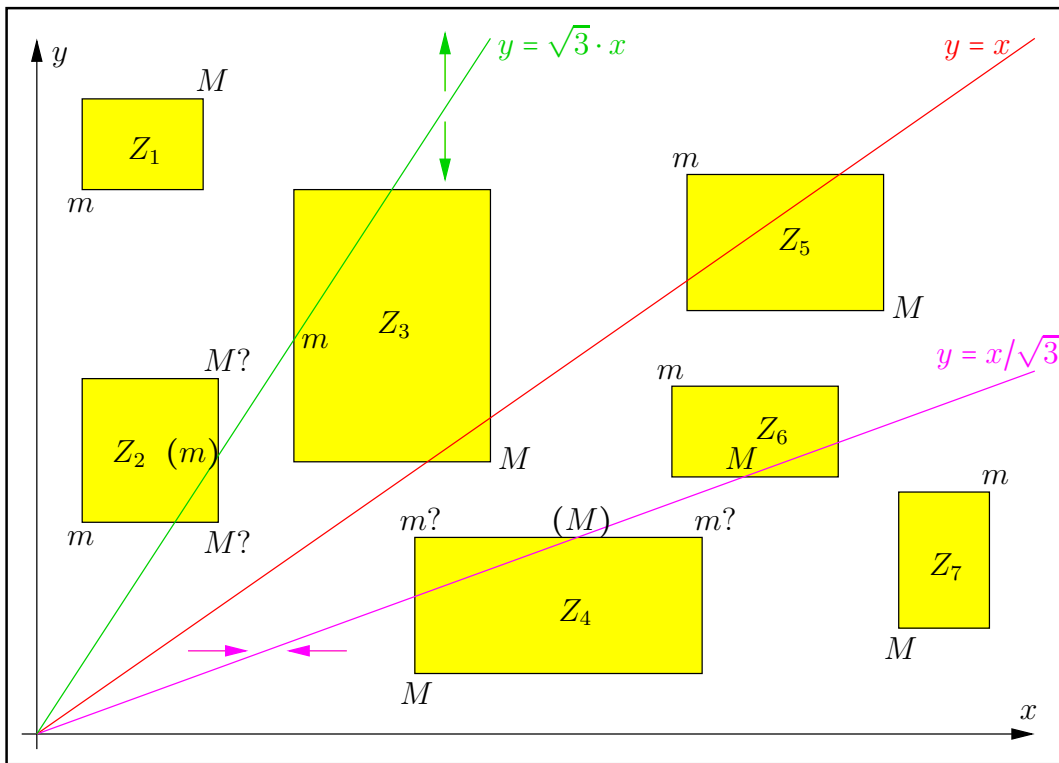


Abbildung B.4.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

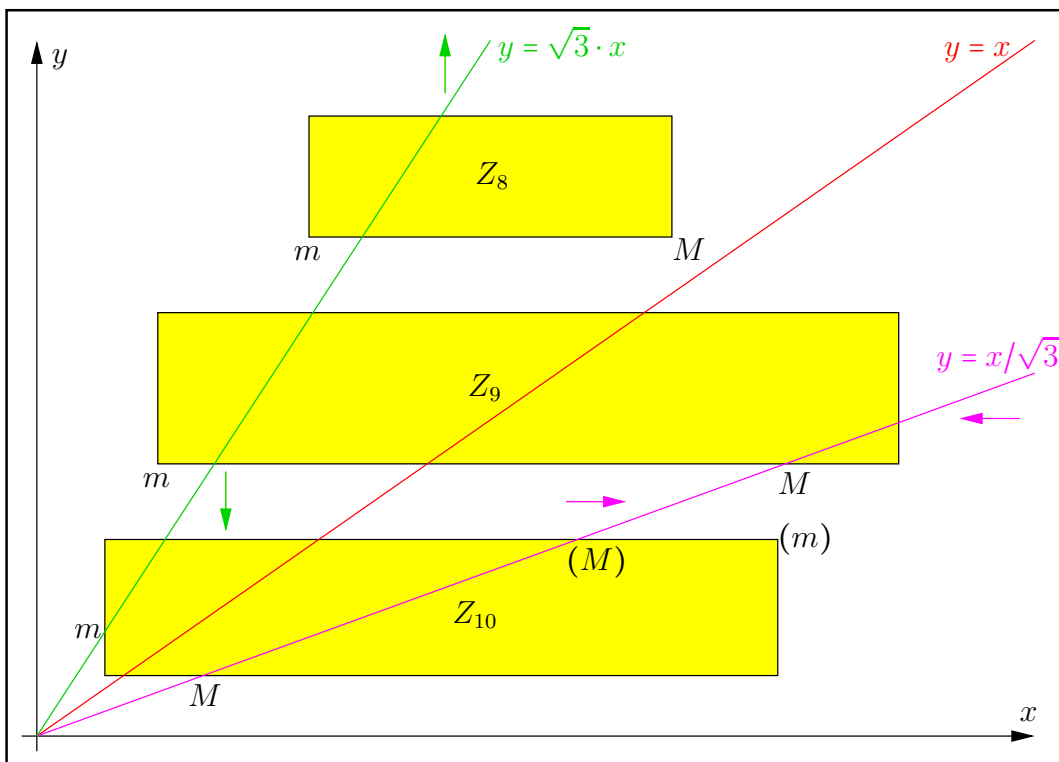


Abbildung B.5.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

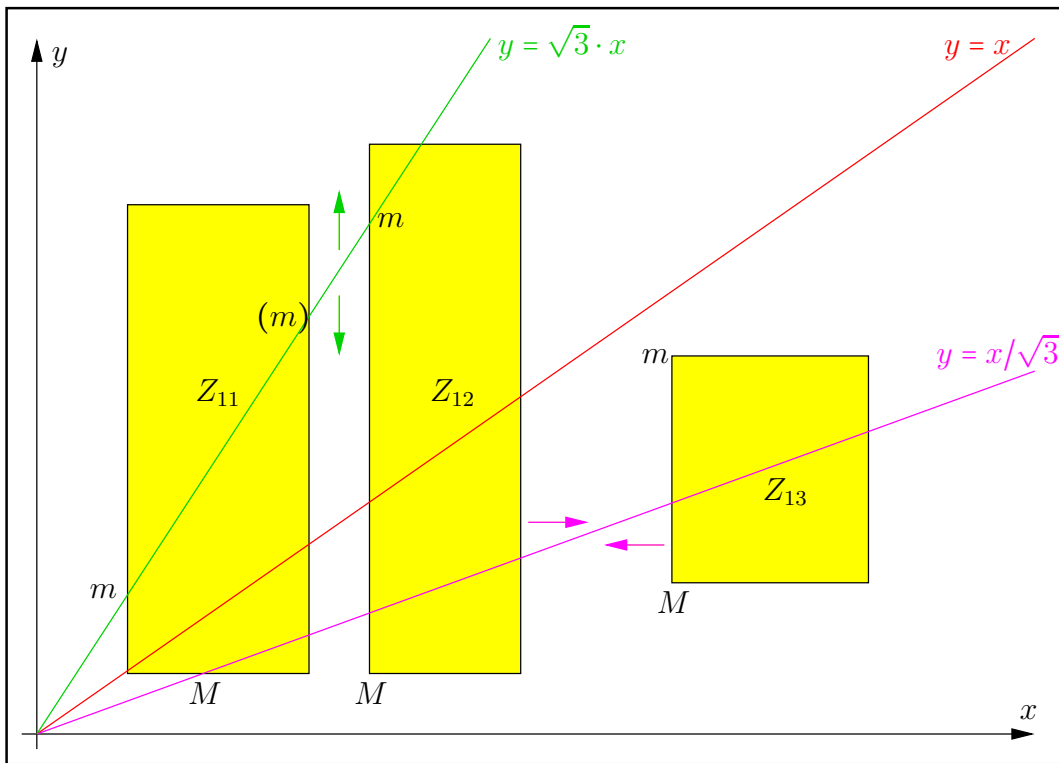


Abbildung B.6.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

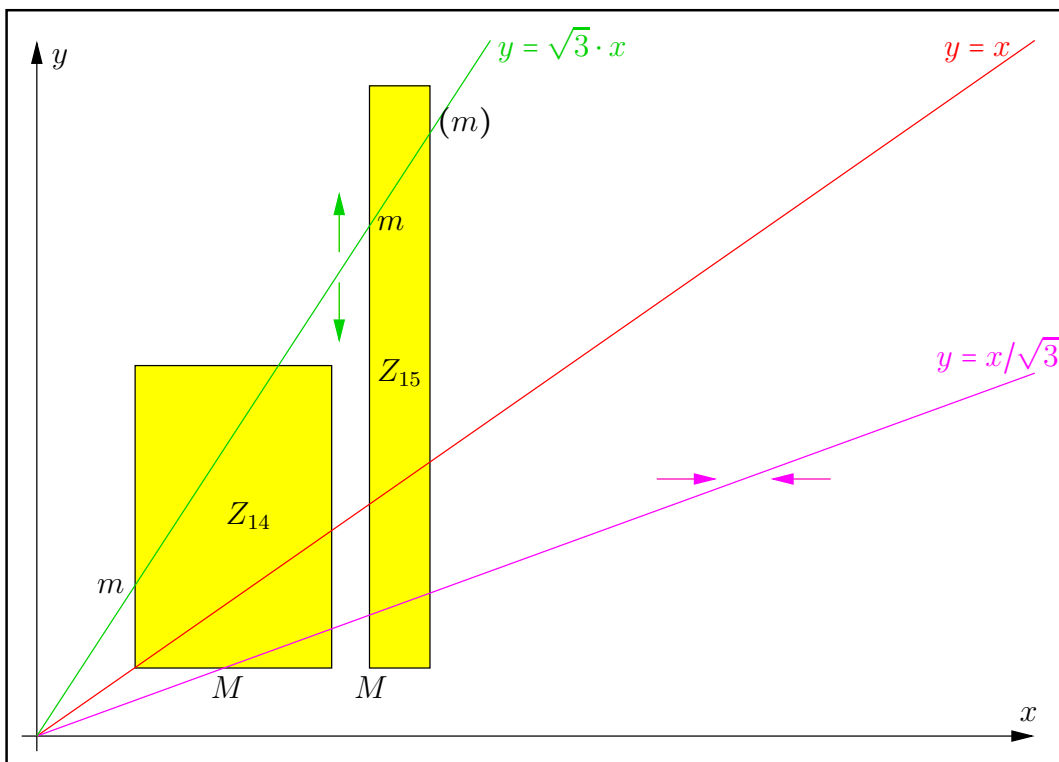


Abbildung B.7.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

## Anmerkungen:

1. Bezüglich des zweidimensionalen Intervalls  $Z_2$  muss noch geklärt werden, dass das Minimum nicht bei  $(m)$ , sondern bei  $m$  angenommen wird. Unter der Voraussetzung  $x_2 > x_1$  muss daher gezeigt werden:  $f(x_1, y_1) = (x_1^2 - y_1^2)/(x_1^2 + y_1^2)^2 < f(x_2, \sqrt{3} \cdot x_2) = -1/(8x_2^2)$ .

$$\text{Wegen } \sqrt{3} \cdot x_1 \leq y_1 \leq \sqrt{3} \cdot x_2 \text{ gilt zunächst } \frac{x_1^2 - y_1^2}{(x_1^2 + y_1^2)^2} \leq \frac{x_1^2 - 3 \cdot x_1^2}{(x_1^2 + 3 \cdot x_1^2)^2} = \frac{-1}{8 \cdot x_1^2} < \frac{-1}{8x_2^2} \blacksquare$$

2. Bezüglich der zweidimensionalen Intervalle  $Z_4, Z_{10}$  muss noch geklärt werden, dass das Maximum nicht bei  $(M)$ , sondern bei  $M$  angenommen wird. Der Beweis ist einfach, denn läuft man von  $(M)$  senkrecht nach unten und dann parallel zu  $x$ -Achse nach links bis  $M$ , so wachsen die Funktionswerte jeweils monoton.
3. Bezüglich des zweidimensionalen Intervalls  $Z_{10}$  muss noch geklärt werden, dass das Minimum nicht bei  $(m)$ , sondern bei  $m$  angenommen wird. Dies ist jedoch trivial, da die Funktionswerte von  $f(x, y)$  bei  $m$  negativ und bei  $(m)$  positiv sind.
4. Bezüglich der zweidimensionalen Intervalle  $Z_{11}, Z_{15}$  muss noch geklärt werden, dass das Minimum nicht bei  $(m)$ , sondern bei  $m$  angenommen wird. Im Fall  $x_2 > x_1$  muss daher gezeigt werden:  $f(x_1, \sqrt{3}x_1) < f(x_2, \sqrt{3}x_2)$ .

$$f(x_1, \sqrt{3}x_1) < f(x_2, \sqrt{3}x_2) \iff \frac{x_1^2 - 3x_1^2}{(x_1^2 + 3x_1^2)^2} < \frac{x_2^2 - 3x_2^2}{(x_2^2 + 3x_2^2)^2} \iff \frac{-1}{x_1^2} < \frac{-1}{x_2^2} \iff x_2 > x_1 \blacksquare$$

### B.3.1. Maximumbestimmung

Zur Bestimmung der Lage von  $M$  sind **nacheinander** folgende Fälle zu unterscheiden:

1. Wenn die untere Parallele zu  $x$ -Achse unterhalb der Extremalkurve  $y = x/\sqrt{3}$  liegt, so ist die Maximumstelle  $M$  die linke untere Ecke, vgl.  $Z_4, Z_7, Z_{12}, Z_{13}, Z_{15}$ .  
Sei  $Ux1$  die optimale Einschließung von  $x_1/\sqrt{3}$ , dann gilt: Die untere Parallele zu  $x$ -Achse liegt unterhalb der Extremalkurve  $y = x/\sqrt{3} \iff y_1 \leq \text{Inf}(Ux1)$ .
2. Wenn die untere Parallele zu  $x$ -Achse die Extremalkurve  $y = x/\sqrt{3}$  schneidet, so ist dieser Schnittpunkt gleich  $M$ , vgl.  $Z_6, Z_9, Z_{10}, Z_{11}, Z_{14}$ .  
Sei  $Ux1$  die optimale Einschließung von  $x_1/\sqrt{3}$  und  $Ux2$  die optimale Einschließung von  $x_2/\sqrt{3}$ , dann gilt: Die untere Parallele zu  $x$ -Achse schneidet die Extremalkurve  $y = x/\sqrt{3} \iff \text{Sup}(Ux1) \leq y_1 \leq \text{Inf}(Ux2)$ .  
Es gilt  $M(x_S, y_S) = M(\sqrt{3} \cdot y_1, y_1)$ , und für den maximalen Funktionswert erhält man:  
 $f(x_S, y_S) = (3y_1^2 - y_1^2)/(3y_1^2 + y_1^2)^2 = 1/(8y_1^2) = (1/(\sqrt{8} \cdot y_1))^2$ .
3. Wenn die rechte Parallele zu  $y$ -Achse die Extremalkurve  $y = \sqrt{3} \cdot x$  schneidet, so ist  $M$  der rechte untere oder obere Eckpunkt, vgl.  $Z_2$ .  
Sei  $Uy1$  die optimale Einschließung von  $y_1/\sqrt{3}$  und  $Uy2$  die optimale Einschließung von  $y_2/\sqrt{3}$ , dann gilt: Die rechte Parallele zu  $y$ -Achse schneidet die Extremalkurve  $y = \sqrt{3}x \iff \text{Sup}(Uy1) \leq x_2 \leq \text{Inf}(Uy2)$ .
4. Wenn der untere rechte Eckpunkt unterhalb der Extremalkurve  $y = \sqrt{3} \cdot x$  liegt, so ist dieser untere rechte Eckpunkt gleich  $M$ , vgl.  $Z_3, Z_5, Z_8$ .  
Sei  $Uy1$  die optimale Einschließung von  $y_1/\sqrt{3}$ , dann gilt: Der untere rechte Eckpunkt liegt unterhalb der Extremalkurve  $y = \sqrt{3} \cdot x \iff \text{Sup}(Uy1) \leq x_2$ .
5. Wenn der untere rechte Eckpunkt oberhalb der Extremalkurve  $y = \sqrt{3} \cdot x$  liegt, so ist  $M$  der obere rechte Eckpunkt, vgl.  $Z_1$ .  
Sei  $Uy1$  die optimale Einschließung von  $y_1/\sqrt{3}$ , dann gilt: Der untere rechte Eckpunkt liegt oberhalb der Extremalkurve  $y = \sqrt{3} \cdot x \iff \text{Inf}(Uy1) \geq x_2$ .

## Anmerkungen:

1. **nacheinander** bedeutet, dass z.B. die 3. Abfrage nur dann erfolgen darf, wenn vorher die 1. und 2. Abfrage jeweils negativ ausgefallen sind.
2. Die Schnittpunkte der zweidimensionalen Maschinenintervalle  $Z_\nu$  mit den Extremalkurven  $y = \sqrt{3} \cdot x$ ,  $y = x/\sqrt{3}$  können wegen der irrationalen Zahl  $\sqrt{3}$  keine Maschinenzahlen sein.

### B.3.2. Minimumbestimmung

Zur Bestimmung der Lage von  $m$  sind **nacheinander** folgende Fälle zu unterscheiden:

1. Wenn die linke Parallele zu  $y$ -Achse links der Extremalkurve  $y = \sqrt{3}x$  liegt, so ist die Minimumstelle  $m$  die linke untere Ecke, vgl.  $Z_1, Z_2, Z_8, Z_9$ .  
Sei  $Uy1$  die optimale Einschließung von  $y_1/\sqrt{3}$ , dann gilt: Die linke Parallele zu  $y$ -Achse liegt links der Extremalkurve  $y = \sqrt{3}x \iff x_1 \leq \text{Inf}(Uy1)$ .
2. Wenn die linke Parallele zu  $y$ -Achse die Extremalkurve  $y = \sqrt{3}x$  schneidet, so ist dieser Schnittpunkt gleich  $m$ , vgl.  $Z_3, Z_{10}, Z_{11}, Z_{12}, Z_{14}, Z_{15}$ .  
Sei  $Uy1$  die optimale Einschließung von  $y_1/\sqrt{3}$  und  $Uy2$  die optimale Einschließung von  $y_2/\sqrt{3}$ , dann gilt: Die linke Parallele zu  $y$ -Achse schneidet die Extremalkurve  $y = \sqrt{3}x \iff \text{Sup}(Uy1) \leq x_1 \leq \text{Inf}(Uy2)$ .  
Es gilt  $m(x_S, y_S) = m(x_1, \sqrt{3} \cdot x_1)$ , und für den minimalen Funktionswert erhält man:  
 $f(x_S, y_S) = (x_1^2 - 3x_1^2)/(x_1^2 + 3x_1^2)^2 = -1/(8x_1^2) = -(1/(\sqrt{8} \cdot x_1))^2$ .
3. Wenn die obere Parallele zu  $x$ -Achse die Extremalkurve  $y = x/\sqrt{3}$  schneidet, so ist  $m$  der obere rechte oder linke Eckpunkt, vgl.  $Z_4$ .  
Sei  $Ux1$  die optimale Einschließung von  $x_1/\sqrt{3}$  und  $Ux2$  die optimale Einschließung von  $x_2/\sqrt{3}$ , dann gilt: Die obere Parallele zu  $x$ -Achse schneidet die Extremalkurve  $y = x/\sqrt{3} \iff \text{Sup}(Ux1) \leq y_2 \leq \text{Inf}(Ux2)$ .
4. Wenn der obere linke Eckpunkt oberhalb der Extremalkurve  $y = x/\sqrt{3}$  liegt, so ist dieser obere linke Eckpunkt gleich  $m$ , vgl.  $Z_5, Z_6, Z_{13}$ .  
Sei  $Ux1$  die optimale Einschließung von  $x_1/\sqrt{3}$ , dann gilt: Der obere linke Eckpunkt liegt oberhalb der Extremalkurve  $y = x/\sqrt{3} \iff \text{Sup}(Ux1) \leq y_2$ .
5. Wenn der obere linke Eckpunkt unterhalb der Extremalkurve  $y = x/\sqrt{3}$  liegt, so ist  $m$  der obere rechte Eckpunkt, vgl.  $Z_7$ .  
Sei  $Ux1$  die optimale Einschließung von  $x_1/\sqrt{3}$ , dann gilt: Der obere linke Eckpunkt liegt unterhalb der Extremalkurve  $y = x/\sqrt{3} \iff y_2 \leq \text{Inf}(Ux1)$ .

Weitere Einzelheiten zum Algorithmus findet man in der Funktion `Re_rz2(X,Y)`, die in der Datei `mpficlass.cpp` definiert ist.

### B.3.3. Numerische Beispiele

Im **1. Beispiel** wählen wir  $X = [1, 4]$ ,  $Y = [1.5, 2]$ , womit die typische Lage des zweidimensionalen Intervalls  $Z_{10}$  realisiert wird. Mit dem Funktionsaufruf `T = Re_rz2(X,Y)`; erhalten wir mit `prec = 300` für den Wertebereich  $W_f := \{f(x,y) \mid x \in X, y \in Y\} = [-1/8, +1/18]$  die nahezu optimale Einschließung  $W_f \subseteq T$ .

$$T = [-0.12500000 \dots 00001, +5.55555 \dots 55556 \cdot 10^{-2}] \supset W_f$$

mit 90 ausgegeben Dezimalstellen. Rechnet man z.B. mit `prec = 3000000`, so erhält man das analoge Ergebnis mit  $3000000 \cdot \log_{10}(2) \approx 903090$  Dezimalstellen.

Im **2. Beispiel** wählen wir  $X = [0, 1]$ ,  $Y = [2, 3]$ , womit die typische Lage des zweidimensionalen Intervalls  $Z_1$  realisiert wird. Mit dem Funktionsaufruf  $T = \text{Re\_rz2}(X, Y)$ ; erhalten wir mit  $\text{prec} = 300$  für den Wertebereich  $W_f := \{f(x, y) \mid x \in X, y \in Y\} = [-1/4, -2/25]$  die nahezu optimale Einschließung  $W_f \subseteq T$ .

$$T = [-2.500000 \dots 000001 \cdot 10^{-1}, -7.999999 \dots 999999 \cdot 10^{-2}] \supset W_f$$

mit 90 ausgegeben Dezimalstellen.

In der folgenden Tabelle sind für alle typischen, zweidimensionalen Intervalle  $Z_\nu$  ihre möglichen reellen Intervalle  $X_\nu, Y_\nu$  zusammen mit den jeweiligen exakten Funktionswerten  $f_m, f_M$  an den Minimum- bzw. Maximumstellen  $m, M$  zusammengestellt:

$\nu$	$X_\nu$	$Y_\nu$	$f_m$	$f_M$	$\nu$	$X_\nu$	$Y_\nu$	$f_m$	$f_M$
1	[0, 1]	[2, 3]	-1/4	-2/25	2	[1, 2]	[3, 4]	-2/25	-5/169
3	[1, 2]	[1.5, 2]	-1/8	28/625	4	[2, 3]	[1, 1.5]	28/625	3/25
5	[4, 5]	[4.5]	-9/1681	+9/1681	6	[4, 5]	[2.6, 3]	7/625	25/1352
7	[1, 2]	[0.25, 0.5]	60/289	240/289	8	[3, 4]	[6, 6.5]	-27/2025	-5/676
9	[1, 3.5]	[2, 2.5]	-3/25	1/32	10	[1, 4]	[1.5, 2]	-1/8	1/18
11	[1, 2]	[1, 4]	-1/8	+1/8	12	[2, 3]	[1, 4]	-1/32	3/25
13	[4, 5]	[2, 3]	7/625	3/100	14	[1, 1.75]	[1, 2]	-1/8	+1/8
15	[2, 3]	[1, 6]	-1/32	3/25					

Der Algorithmus zur Berechnung der gesuchten Extrema  $f_m, f_M$  wurde u.a. mit allen obigen  $Z_\nu := \{(x, y) \mid x \in X_\nu \wedge y \in Y_\nu\}$  getestet.

## B.4. $2xy/(x^2 + y^2)^2$

Für die vorgegebenen reellen Maschinenintervalle  $X \ni x$  und  $Y \ni y$  liefert die Intervallfunktion

```
MpfiClass mIm_rz2(const MpfiClass& X, const MpfiClass& Y)
```

eine fast optimale Einschließung aller Funktionswerte

$$f(x, y) = \frac{2xy}{(x^2 + y^2)^2} \quad \text{mit: } x \in X = [x_1, x_2], y \in Y = [y_1, y_2] \quad \text{und } x \neq 0 \wedge y \neq 0.$$

Die Berechnung von  $f(x, y)$  wird ausführlich beschrieben auf Seite 187.  $f(x, y)$  ist als negative Imaginärteildfunktion der holomorphen Funktion  $1/z^2$ ,  $z \in \mathbb{C} \setminus \{0\}$ , eine harmonische Funktion, wenn der singuläre Punkt  $P(0, 0)$  nicht ein Element des zweidimensionalen reellen Intervalls  $Z$  ist.

$$(B.11) \quad Z := \{(x, y) \mid x \in X \wedge y \in Y\}, \quad 0 \notin X \wedge 0 \notin Y.$$

$f(x, y)$  nimmt daher für alle Punkte  $(x, y) \in Z$  als harmonische Funktion ihre Extrema auf dem Rand von  $Z$  an, vgl. dazu die Seite 244. Um die zugehörigen Extremstellen,  $m$  für das Minimum und  $M$  für das Maximum, möglichst einfach bestimmen zu können, betrachten wir zunächst die folgenden Symmetrieeigenschaften:

$$(B.12) \quad f(x, y) \equiv f(y, x), \quad \text{Spiegelung an der 1. Winkelhalbierenden,}$$

$$(B.13) \quad f(-x, -y) \equiv f(y, x), \quad \text{Spiegelung am Ursprung,}$$

$$(B.14) \quad f(-x, y) \equiv -f(x, y), \quad \text{Spiegelung an der } y\text{-Achse.}$$

Für Rechteckintervalle  $Z$  ergeben sich aus (B.12) und (B.13) die folgenden Eigenschaften:

1. Ist  $Z_W := \{(x, y) \mid x \in Y \wedge y \in X\}$  das an der 1. Winkelhalbierenden gespiegelte Rechteck, so nimmt  $f(x, y)$  über den Rechteckintervallen  $Z_W$  und  $Z$  genau die gleichen Funktionswerte an, d.h. die Wertebereiche über  $Z_W$  und  $Z$  sind identisch.
2. Bedeutet  $Z_U := \{(x, y) \mid x \in -X \wedge y \in -Y\}$  das am Ursprung gespiegelte Rechteck, so nimmt  $f(x, y)$  über den Rechteckintervallen  $Z_U$  und  $Z$  genau die gleichen Funktionswerte an.

Mit Hilfe dieser beiden Spiegelungen kann jetzt ein Rechteckintervall  $Z = \{(x, y) \mid x \in X \wedge y \in Y\}$ , das teilweise oder ganz in der unteren Halbebene liegt, ganz in die obere Halbebene so transformiert werden, dass die Wertebereiche von  $f(x, y)$  über  $Z$  und dem transformierten Rechteck  $Z_T$  wieder identisch sind.  $Z_T$  ergibt sich aus

```
if (y2<=0) { X=-X; Y=-Y; }
else
  if (y1<0)
    swap(X,Y);
    if (y1<0) {X=-X; Y=-Y; }
```

### Anmerkungen:

1. Nach dieser Transformation liegt  $Z_T = \{(x, y) \mid x \in X \wedge y \in Y\}$  ganz in der oberen Halbebene, und im Fall  $x_1 \geq 0$  liegt  $Z_T$  sogar ganz im ersten Quadranten, so dass dann die Bestimmung von  $m, M$  auf den 1. Quadranten beschränkt werden kann. Im Fall  $y_2 > x_2$  werden  $X$  und  $Y$  nochmals vertauscht, so dass dann im ersten Quadranten nur noch Rechtecke zu betrachten sind, deren obere rechte Eckpunkte stets unterhalb der 1. Winkelhalbierenden liegen. Hierdurch wird die Anzahl der notwendigen Fallunterscheidungen nochmals deutlich reduziert.



2. Im Fall  $x_2 \leq 0$  liegt  $Z_T$  ganz im 2. Quadranten, und mit der zusätzlichen Transformation  $X = -X$  liegt dann das neue  $Z_{T1}$  ganz im 1. Quadranten, nach (B.14) jedoch mit den Funktionswerten  $-f(x, y)$ . Ist dann  $[f_m, f_M]$  der berechnete Wertebereich von  $f(x, y)$  über  $Z_{T1}$ , so ist  $[-f_m, -f_M] = [-f_M, -f_m]$  der gesuchte Wertebereich von  $f(x, y)$  über  $Z_T$ .
3. Im Fall  $x_1 < 0 \wedge x_2 > 0$  wird mit  $X_2 := [x_1, 0]$  und  $X_1 := [0, x_2]$  das Rechteck  $Z_T$  aufgeteilt in  $Z_{T2} = \{(x, y) \mid x \in X_2 \wedge y \in Y\}$  und  $Z_{T1} = \{(x, y) \mid x \in X_1 \wedge y \in Y\}$ , so dass in  $Z_{T2}$  wegen  $f(x, y) \leq 0$  nur der Minimumpunkt  $m$  und in  $Z_{T1}$  wegen  $f(x, y) \geq 0$  nur der Maximumpunkt  $M$  liegen kann. Daher wird  $Z_{T1}$  wie unter Punkt 1. und  $Z_{T2}$  wie unter Punkt 2. behandelt.

### Zusammenfassung:

Um zu einem gegebenen Rechteckintervall  $Z \not\subset (0, 0)$  den Wertebereich  $[f_m, f_M]$  von  $f(x, y)$  zu bestimmen, wird  $Z$  zunächst nach Seite 216 als  $Z_T$  in die obere Halbebene transformiert. Nach den Punkten 1,2,3 von Seite 216 müssen dann die Extremwerte  $f_m$  und  $f_M$  nur noch für solche Rechteckintervalle aus dem 1. Quadranten bestimmt werden, deren obere rechte Ecken unterhalb der 1. Winkelhalbierenden liegen.

Die Berechnung der gerundeten Extrema  $f_m$  bzw.  $f_M$  für ein Rechteckintervall  $Z$ , nur aus dem 1. Quadranten mit seinem oberen rechten Eckpunkt unterhalb der 1. Winkelhalbierenden, erfolgt mit den beiden Funktionen

```
MpfrClass mIm_rz2_m (const MpfiClass& X, const MpfiClass& Y);
MpfrClass mIm_rz2_M (const MpfiClass& X, const MpfiClass& Y);
```

die in der Datei `mpficlass.cpp` definiert sind. Zur Berechnung der beiden Extremstellen  $m, M$  benötigen wir zunächst die partiellen Ableitungen:

$$(B.15) \quad \frac{\partial f(x, y)}{\partial x} = \frac{2y \cdot (-3x^2 + y^2)}{(x^2 + y^2)^3}$$

$$(B.16) \quad \frac{\partial f(x, y)}{\partial y} = \frac{2x \cdot (x^2 - 3y^2)}{(x^2 + y^2)^3}.$$

Die Extremalkurven für  $\partial f/\partial x$  ergeben sich aus der Forderung  $\partial f/\partial x = 0$ , siehe Seite 247. Die erste Extremalkurve ist wegen  $y = 0$  die positive  $x$ -Achse, und die zweite ergibt sich aus der Forderung  $y^2 = 3x^2$  aus der für den 1. Quadranten  $y = \sqrt{3} \cdot x$  folgt.

Die beiden Extremalkurven für  $\partial f/\partial y$  ergeben sich analog aus der Forderung  $\partial f/\partial y = 0$ . Die erste Extremalkurve ist dann wegen  $x = 0$  die positive  $y$ -Achse und die zweite Extremalkurve ist gegeben durch  $y = x/\sqrt{3}$ .

Da beide partiellen Ableitungen auf den Extremalkurven  $y = \sqrt{3} \cdot x$  und  $y = x/\sqrt{3}$  das Vorzeichen wechseln, wenn man parallel zu den beiden Koordinatenachsen diese Extremalkurven schneidet, besitzt  $f(x, y)$  an den Schnittpunkten der Intervallränder von  $Z$  mit den Extremalkurven längs dieser Intervallränder tatsächlich relative Extrema, d.h. in diese Fällen müssen die Extrempunkte  $m$  oder  $M$  nicht auf den Eckpunkten von  $Z$  liegen.

In den folgenden Abbildungen sind alle typischen Lagen der Intervalle  $Z_\nu$  aus dem 1. Quadranten mit den entsprechenden Lagen der Extremstellen  $m(x, y), M(x, y)$  angegeben. Dabei bedeutet das Symbol  $(M)$ , dass hier zwar ein relatives Maximum vorliegt, dass aber das wirkliche Maximum an der anderen Stelle  $M$  angenommen wird. Die zwei Symbole  $m?$  bedeuten, dass das wirkliche Minimum nur an einer dieser beiden Stellen angenommen werden kann, so dass zwei Berechnungen und ein Vergleich erforderlich sind. Die 1. Winkelhalbierende und die beiden Extremalkurven  $y = \sqrt{3} \cdot x$  und  $y = x/\sqrt{3}$  sind leicht gedreht dargestellt. Die farbigen Pfeile geben die Richtung wachsender Funktionswerte von  $f(x, y) = +2xy/(x^2 + y^2)^2$  an.

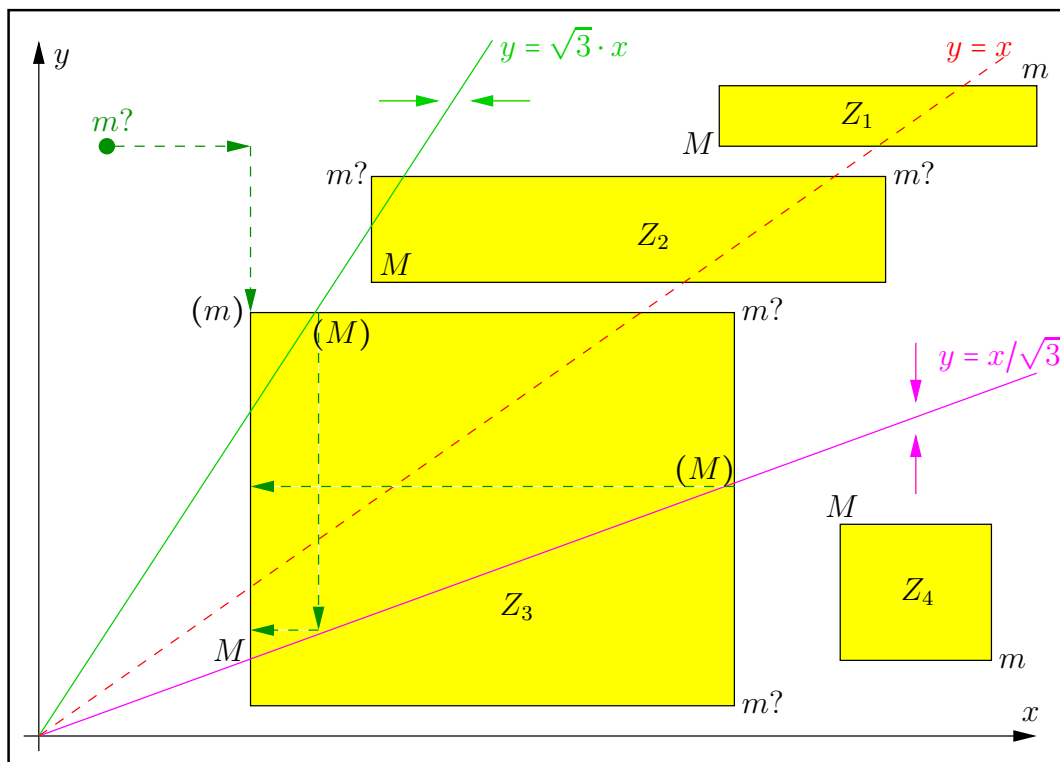


Abbildung B.8.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

**Anmerkungen:**

1. Die 1. Winkelhalbierende und die Extremalkurven  $y = \sqrt{3} \cdot x$  bzw.  $y = x/\sqrt{3}$  sind etwas gedreht dargestellt, so dass z.B. der an der 1. Winkelhalbierenden gespiegelte untere rechte Eckpunkt  $m?$  durch  $m?$  nicht ganz maßstabsgetreu dargestellt wird.
2. Bei den Rechtecken  $Z_1, Z_2, Z_4$  liegen die Extrempunkte  $m, M$  auf den Eckpunkten dieser Rechtecke, wobei bez.  $Z_2$  der Minimumpunkt  $m$  einer der beiden oberen Eckpunkte ist.
3. Wir zeigen zuerst, dass am oberen linken Eckpunkt  $(m)$  von  $Z_3$  die Funktion  $f(x, y)$  nicht ihr Minimum annehmen kann. Dazu spiegeln wir den unteren rechten Eckpunkt  $m?$  an der 1. Winkelhalbierenden und erhalten den Spiegelpunkt  $m?$ . Mit  $m?(x_2, y_1)$  gilt deshalb  $m?(y_1, x_2)$ , und nach (B.12) gilt  $f_{m?} := f(x_2, y_1) = f_{m?} := f(y_1, x_2)$ . Im nächsten Schritt zeigen wir, dass der Spiegelpunkt  $m?$  links oberhalb von  $(m)(x_1, y_2)$  liegt, d.h. zu zeigen ist  $y_1 < x_1$  und  $y_2 < x_2$ . Dazu benutzen wir die für die Lage von  $Z_3$  gültigen Bedingungen:

$$y_1 < x_1/\sqrt{3} < x_1 \quad \text{und} \quad y_2 < x_2,$$

aus denen die Behauptung unmittelbar folgt. In Abb. (B.8) wird mit den Pfeilen  $\text{-----}>$  gezeigt, wie man von  $m?$  nach  $(m)$  kommt, und weil  $\text{----}>$  die Richtung wachsender Funktionswerte angibt, gilt  $f_{m?} = f_{m?} < f_{(m)}$ , so dass  $(m)$  nicht der gesuchte Minimumpunkt  $m$  sein kann ■

4. Wir zeigen jetzt, dass der rechte Schnittpunkt  $(M)$  mit der Extremalkurve  $y = x/\sqrt{3}$  nicht der gesuchte Maximumpunkt  $M$  sein kann. Dazu bewegen wir uns zunächst von  $(M)$  aus in Richtung des Pfeils  $\text{<----}$  bis zum Schnittpunkt mit der linken Parallelen zur  $y$ -Achse von  $Z_3$  und von dort aus weiter in Richtung wachsender Funktionswerte nach unten bis  $M$ , so dass  $(M)$  nicht der Maximumpunkt sein kann ■ Ganz analog zeigt man, dass auch der obere Schnittpunkt  $(M)$  mit der Extremalkurve  $y = \sqrt{3} \cdot x$  kein Maximumpunkt sein kann.

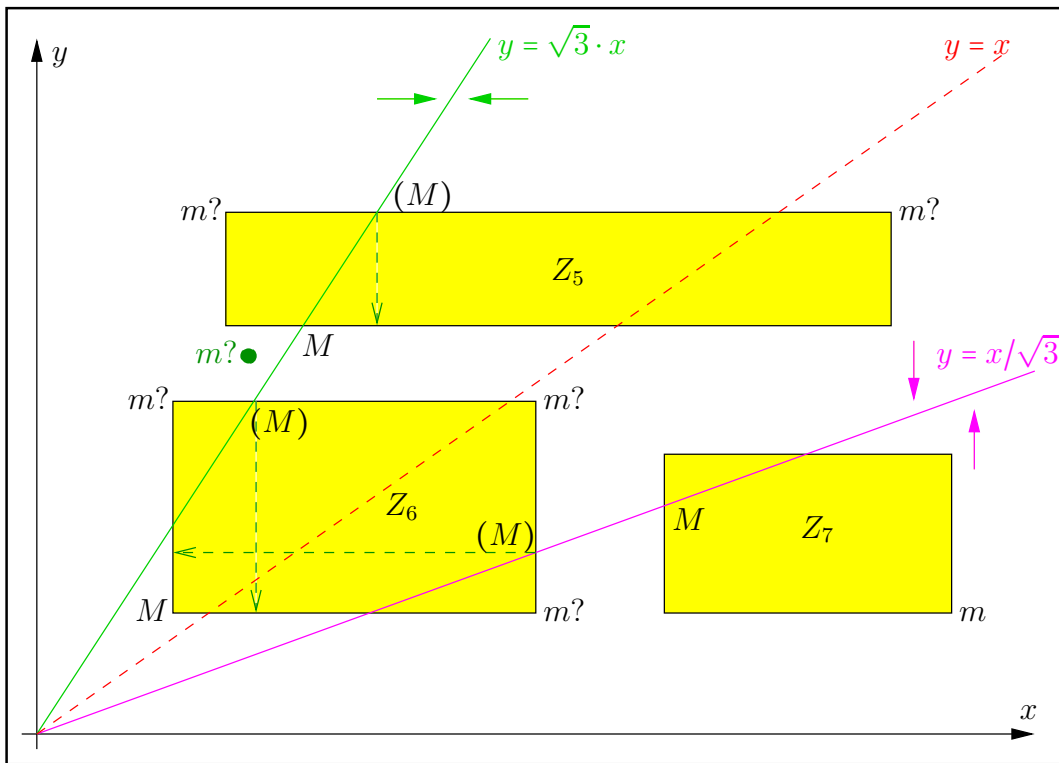


Abbildung B.9.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

Da in  $Z_6$  wegen  $y_1 > x_1$ , im Gegensatz zu  $Z_3$  aus Abb. B.8, der Spiegelpunkt  $m?(y_1, x_2)$  rechts von  $m?(x_1, y_2)$  liegt, kann  $m?(x_1, y_2)$  in  $Z_6$  nicht als Minimumpunkt ausgeschlossen werden.

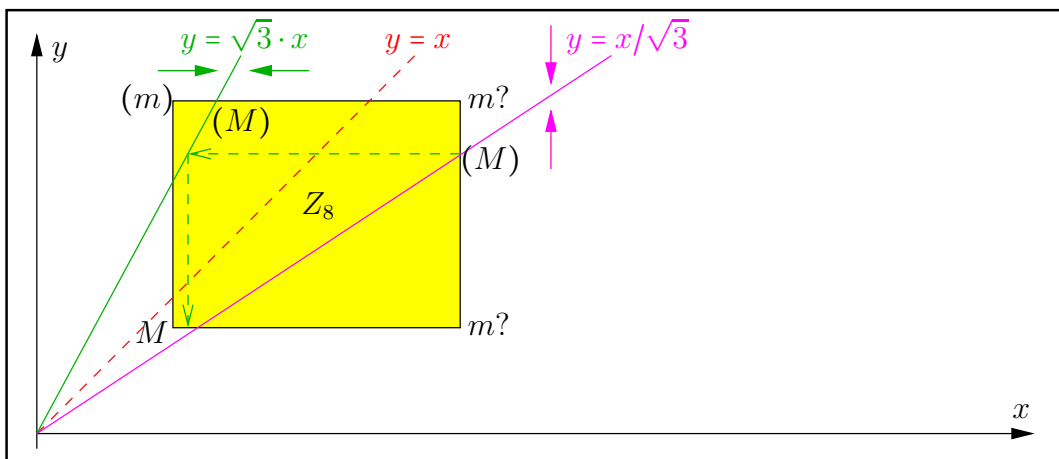


Abbildung B.10.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

In  $Z_8$  liegen jetzt bez. des Minimumpunktes  $m$  die gleichen Verhältnisse vor wie in  $Z_3$  aus Abb. (B.8), da wegen  $y_1 < x_1$  und  $y_2 < x_2$  der Spiegelpunkt  $m?(y_1, x_2)$  wieder links oberhalb von  $(m)(x_1, y_2)$  liegt, so dass  $(m)$  auch jetzt kein Minimumpunkt sein kann. Ganz analog zeigt man auch jetzt mit Hilfe der Pfeile  $\leftarrow$ , welche die Richtung wachsender Funktionswerte angeben, dass die beiden Punkte  $(M)$  keine Maximumpunkte sein können.

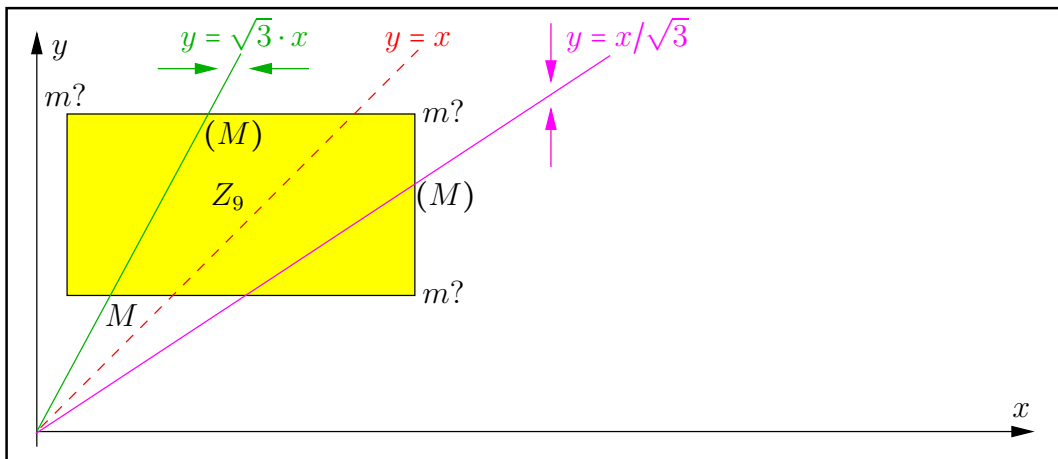


Abbildung B.11.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

Da in  $Z_9$  wegen  $y_1 > x_1$ , im Gegensatz zu  $Z_3$  aus Abb. B.8, der Spiegelpunkt  $m?(y_1, x_2)$  rechts von  $m?(x_1, y_2)$  liegt, kann  $m?(x_1, y_2)$  in  $Z_9$  nicht als Minimumpunkt  $m$  ausgeschlossen werden. Damit sind für  $m$  drei Berechnungen und zwei Vergleiche notwendig. Ganz analog zeigt man auch jetzt mit Hilfe der Pfeile  $\leftarrow$ , welche die Richtung wachsender Funktionswerte angeben, dass die beiden Punkte  $(M)$  keine Maximumpunkte sein können.

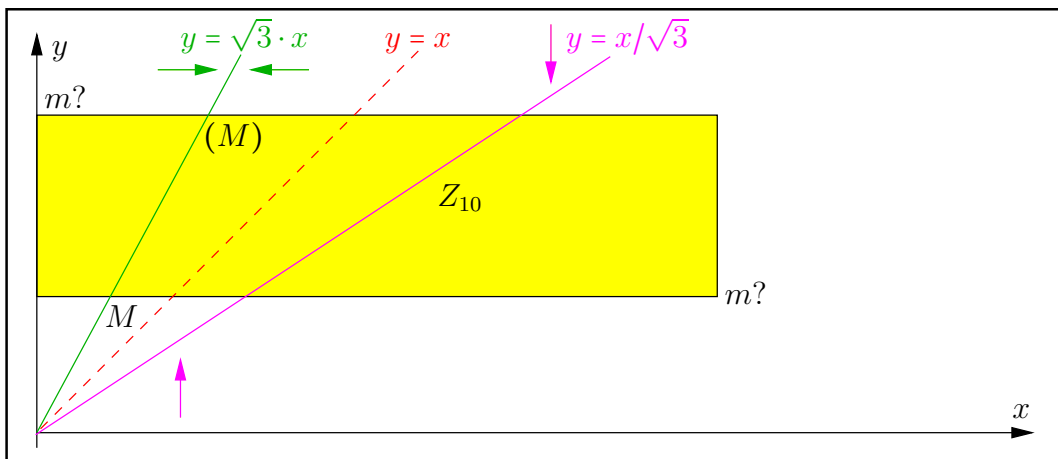


Abbildung B.12.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

Da in  $Z_{10}$  wegen  $y_1 > x_1$ , im Gegensatz zu  $Z_3$  aus Abb. B.8, der Spiegelpunkt  $m?(y_1, x_2)$  von  $m?(x_2, y_1)$  rechts von  $m?(x_1, y_2)$  liegt, kann  $m?(x_1, y_2)$  in  $Z_{10}$  nicht als Minimumpunkt  $m$  ausgeschlossen werden. Damit sind für  $m$  zwei Berechnungen und ein Vergleiche notwendig. Ganz analog zeigt man auch jetzt mit Hilfe des Pfeils  $\leftarrow$ , welcher die Richtung wachsender Funktionswerte angibt, dass der Punkte  $(M)$  mit seinem relativen Maximum nicht der Punkt  $M$  mit dem gesuchten absoluten Maximum auf dem Rand von  $Z_{10}$  sein kann.

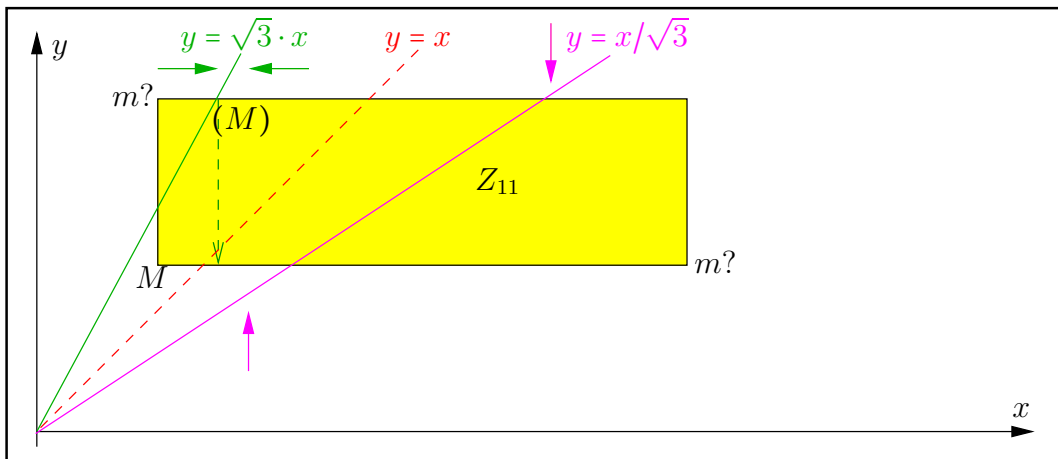


Abbildung B.13.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

Da in  $Z_{11}$  wegen  $y_1 > x_1$ , im Gegensatz zu  $Z_3$  aus Abb. B.8, der Spiegelpunkt  $m?(y_1, x_2)$  von  $m?(x_2, y_1)$  rechts von  $m?(x_1, y_2)$  liegt, kann  $m?(x_1, y_2)$  in  $Z_{11}$  nicht als Minimumpunkt  $m$  ausgeschlossen werden. Damit sind für  $m$  zwei Berechnungen und ein Vergleich notwendig. Ganz analog zeigt man auch jetzt mit Hilfe des Pfeils  $\leftarrow$ , welcher die Richtung wachsender Funktionswerte angibt, dass der Punkte  $(M)$  mit seinem relativen Maximum nicht der Punkt  $M$  mit dem gesuchten absoluten Maximum auf dem Rand von  $Z_{11}$  sein kann.

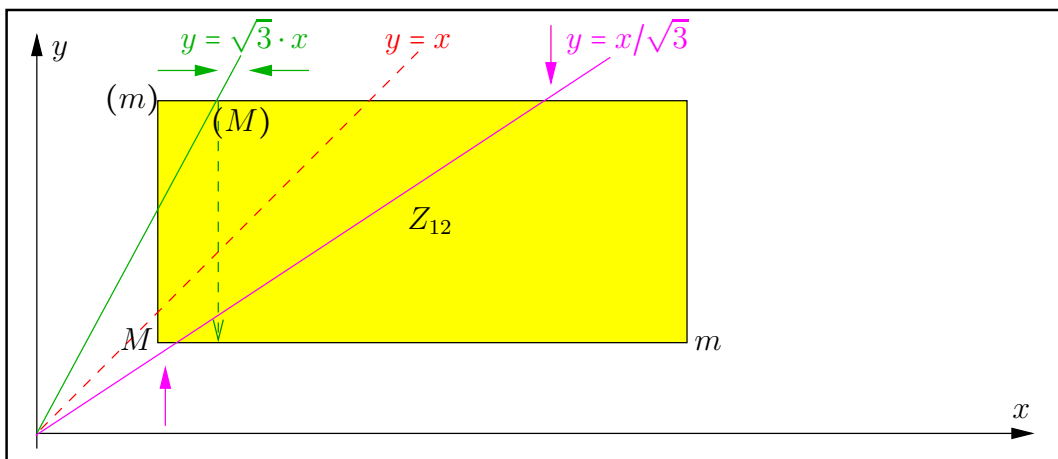


Abbildung B.14.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

In  $Z_{12}$  liegen jetzt bez. des Minimumpunktes  $m$  die gleichen Verhältnisse vor wie in  $Z_3$  aus Abb. (B.8), da wegen  $y_1 < x_1$  und  $y_2 < x_2$  der Spiegelpunkt  $m(y_1, x_2)$  von  $m(x_2, y_1)$  wieder links oberhalb von  $(m)(x_1, y_2)$  liegt, so dass  $(m)$  auch jetzt kein Minimumpunkt sein kann. Ganz analog zeigt man auch jetzt mit Hilfe des Pfeils  $\leftarrow$ , welcher die Richtung wachsender Funktionswerte angibt, dass der Punkte  $(M)$  mit seinem relativen Maximum nicht der Punkt  $M$  mit dem gesuchten absoluten Maximum auf dem Rand von  $Z_{12}$  sein kann.

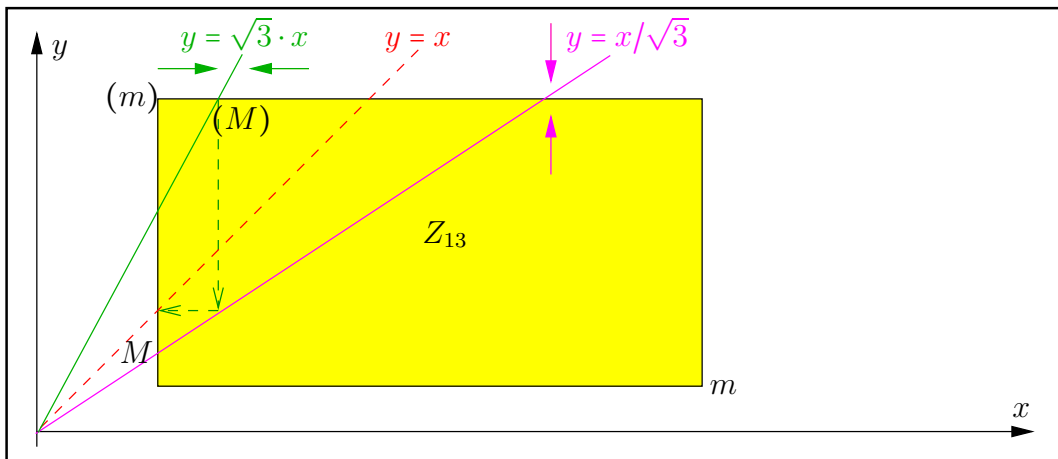


Abbildung B.15.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

In  $Z_{13}$  liegen jetzt bez. des Minimumpunktes  $m$  die gleichen Verhältnisse vor wie in  $Z_3$  aus Abb. (B.8), da wegen  $y_1 < x_1$  und  $y_2 < x_2$  der Spiegelpunkt  $m(y_1, x_2)$  von  $m(x_2, y_1)$  wieder links oberhalb von  $(m)(x_1, y_2)$  liegt, so dass  $(m)$  auch jetzt kein Minimumpunkt sein kann. Ganz analog zeigt man auch jetzt mit Hilfe der Pfeile  $\leftarrow$ , welche die Richtung wachsender Funktionswerte angeben, dass der Punkte  $(M)$  mit seinem relativen Maximum nicht der Punkt  $M$  mit dem gesuchten absoluten Maximum auf dem Rand von  $Z_{13}$  sein kann.

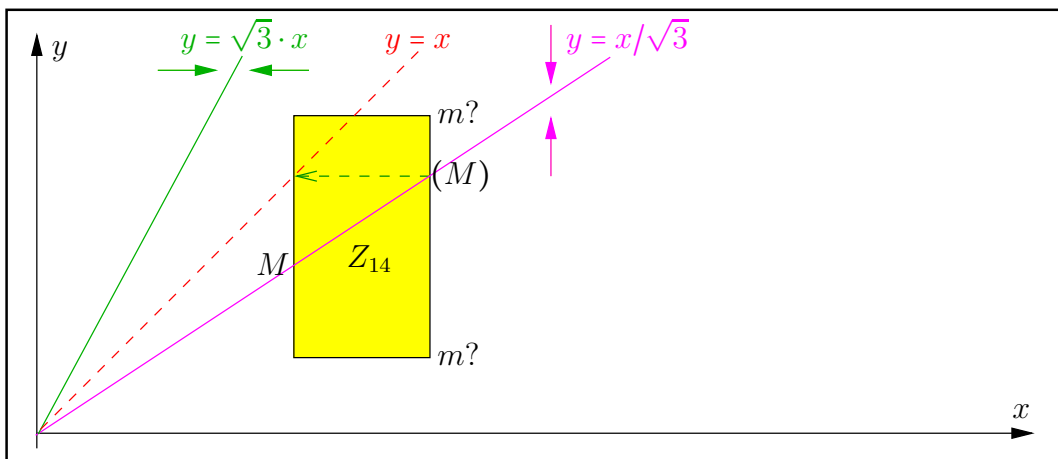


Abbildung B.16.: Typische Lagen der Intervalle  $Z_\nu$  mit den Extremstellen  $m, M$

In  $Z_{14}$  ist der Minimumpunkt der echte obere oder untere Eckpunkt, so dass zwei Berechnungen und ein Vergleich erforderlich sind.  $(M)$  kann nicht der Maximumpunkt sein, denn läuft man von  $(M)$  aus in Richtung des Pfeils  $\leftarrow$ , der die Richtung wachsender Funktionswerte angibt, so kommt man nur auf Wegen mit wachsenden Funktionswerten zum Punkt  $M$ .

In den Abbildungen B.8 bis B.16 sind alle typischen Lagen der Rechteckintervalle  $Z_\nu$ , mit  $\nu = 1, 2, \dots, 14$ , und  $y_2 \leq x_2$  zusammengestellt, wobei der rechte obere Eckpunkt eines  $Z_\nu$  nicht oberhalb der 1. Winkelhalbierenden liegt.

### B.4.1. Maximumbestimmung

Es gelten die folgenden äquivalenten Aussagen:

1. Der untere linke Eckpunkt  $C(x_1, y_1)$  liegt zwischen den Extremalkurven  $y = \sqrt{3} \cdot x$  und  $y = x/\sqrt{3} \iff x_1/\sqrt{3} < y_1 < x_1 \cdot \sqrt{3} \iff M = C(x_1, y_1)$ , vgl. die Rechtecke  $Z_1, Z_2, Z_6, Z_8, Z_{11}, Z_{12}$ .
2.  $S$  ist Schnittpunkt der linken Rechteck-Parallelen zur  $y$ -Achse mit der Extremalkurve  $y = x/\sqrt{3} \iff y_1 < x_1/\sqrt{3} < y_2 \iff M = S(x_1, x_1/\sqrt{3})$ , vgl.  $Z_3, Z_7, Z_{13}, Z_{14}$ .

$$f_M := f(x_1, x_1/\sqrt{3}) = \begin{cases} \left(\frac{1}{x_1}\right)^2 / \frac{8}{\sqrt{27}}, & x_1 > 1 \\ \left(\frac{1}{2x_1}\right)^2 / \frac{2}{\sqrt{27}}, & x_1 < 1. \end{cases}$$

3. Die untere Parallele zur  $x$ -Achse schneidet die Gerade  $y = \sqrt{3} \cdot x \iff x_1 < y_1/\sqrt{3} < x_2 \iff M = S(y_1/\sqrt{3}, y_1)$ , vgl. die Rechtecke  $Z_5, Z_9, Z_{10}$ .

**Anmerkung:** Die Schnittpunktbedingung  $x_1 < y_1/\sqrt{3} < x_2$  kann reduziert werden auf  $x_1 < y_1/\sqrt{3}$ , da nach Voraussetzung gilt:  $y_1/\sqrt{3} < y_1 \leq y_2 \leq x_2$ .

$f_M := f(y_1/\sqrt{3}, y_1)$  kann nach (B.12) wie unter Punkt 2. berechnet werden, wenn man  $x_1$  durch  $y_1$  ersetzt.

4. Der obere linke Eckpunkt  $C(x_1, y_2)$  liegt unterhalb der Extremalkurve  $y = x/\sqrt{3} \iff y_2 < x_1/\sqrt{3} \iff M = C(x_1, y_2)$ , vgl.  $Z_4$ .

Daraus ergibt sich der folgende Algorithmus zur Bestimmung des Maximumpunktes  $M$ :

```

if (y2 < x1/sqrt(3))
  M = C; // nach 4.
else // x1/sqrt(3) < y2
{
  if (y1 < x1/sqrt(3))
    M = S; // nach 2.
  else // y1 > x1/sqrt(3), x1/sqrt(3) < y2;
  {
    if (y1/sqrt(3) < x1)
      M = C; // nach 1.
    else // sqrt(3)*x1 < y1, y1 > x1/sqrt(3), x1/sqrt(3) < y2;
      M = S; // nach 3.
  }
}

```

Mit dem obigen recht einfachen Algorithmus, der alle typischen Lagen der vierzehn Rechtecke  $Z_\nu$  zusammenfasst, wird der maximale Funktionswert über dem Punkt  $M$  berechnet mit Hilfe der Funktion

```
MpfrClass mIm_rz2_M (const MpfiClass& x, const MpfiClass& y),
```

die in `mpficlass.cpp` definiert ist. Der Rückgabewert ist der nahezu optimal aufgerundete, maximale Funktionswert von  $f(x, y)$  über dem Rechteck  $Z_\nu$ , dessen oberer rechter Eckpunkt unterhalb der 1. Winkelhalbierenden liegen muss. Diese letzte Bedingung wird innerhalb der Funktion `mIm_rz2_M(...)` jedoch nicht überprüft und muss damit vorher abgesichert werden. Beachten Sie, dass bei der internen Auswertung von  $f(x, y)$  ein vorzeitiger Überlauf durch geeignete Skalierungen verhindert wird.

## B.4.2. Minimumbestimmung

Aus den Abbildungen B.8 bis B.16 ergibt sich, dass die Minimumpunkte  $m$  nur auf den Eckpunkten der  $Z_\nu$  liegen, da die Funktionswerte  $f(x, y)$  nur anwachsen, wenn man sich parallel zu den Koordinatenachsen auf die Extremalkurven  $y = \sqrt{3} \cdot x$  oder  $y = x/\sqrt{3}$  zubewegt. Um die Lage der Minimumpunkte  $m$  möglichst einfach, d.h. mit einem Minimum an Abfragen, bestimmen zu können, betrachten wir nacheinander die folgenden drei Fälle:

1. Der untere rechte Eckpunkt  $C(x_2, y_1)$  liegt oberhalb  $y = x/\sqrt{3}$ , d.h.  $x_2/\sqrt{3} < y_1$ .
2. Der obere rechte Eckpunkt  $C(x_2, y_2)$  liegt unterhalb  $y = x/\sqrt{3}$ , d.h.  $y_2 < x_2/\sqrt{3}$ .
3. Die rechte Parallele zur  $y$ -Achse schneidet die Extremalkurve  $y = x/\sqrt{3}$ . Dieser Fall tritt ein, wenn die Bedingungen aus 1. und 2. nicht erfüllt sind.

Die oberen drei Fälle werden jetzt genauer betrachtet:

1. Zu diesem Fall gehören die Rechtecke  $Z_1, Z_2, Z_5$ .  
Wenn zusätzlich der obere linke Eckpunkt  $C(x_1, y_2)$  oberhalb von  $y = \sqrt{3} \cdot x$  liegt, d.h. wenn gilt  $y_2/\sqrt{3} > x_1$ , dann können  $C(x_2, y_2)$  oder auch  $C(x_1, y_2)$  die Minimumpunkte sein, vgl.  $Z_2, Z_5$ , sonst ist nur  $C(x_2, y_2)$  der Minimumpunkt, vgl.  $Z_1$ .
2. Zu diesem Fall gehören die Rechtecke  $Z_4, Z_7, Z_{10}, Z_{11}, Z_{12}, Z_{13}$ .  
Für die Rechtecke  $Z_4, Z_7, Z_{12}, Z_{13}$  gilt:  $m = C(x_2, y_1)$ . Wenn jedoch der Spiegelpunkt  $C(y_1, x_2)$  von  $C(x_2, y_1)$  rechts oberhalb von  $C(x_1, y_2)$  liegt, d.h. wenn gilt  $y_1 > x_1 \wedge x_2 \geq y_2$ , so kann auch der obere linke Eckpunkt  $C(x_1, y_2)$  ein Minimum sein, vgl.  $Z_{10}, Z_{11}$ . Beachten Sie, dass die zweite Bedingung  $y_2 \leq x_2$  automatisch erfüllt ist, da vorausgesetzt wird, dass der obere rechte Eckpunkt  $C(x_2, y_2)$  unterhalb der 1. Winkelhalbierenden liegt. Außerdem ist zu beachten, dass die erste Bedingung  $y_1 > x_1$  für die Rechtecke  $Z_4, Z_7, Z_{12}, Z_{13}$  nicht erfüllt ist und damit nur für  $Z_{10}, Z_{11}$  zutrifft.
3. Zu diesem letzten Fall gehören die Rechtecke  $Z_3, Z_6, Z_8, Z_9, Z_{14}$ .  
Wenn also die rechte Parallele zur  $y$ -Achse die Extremalkurve  $y = x/\sqrt{3}$  schneidet, dann ist einer der beiden Eckpunkte  $C(x_2, y_1)$  oder  $C(x_2, y_2)$  gleich  $m$ . Wenn dann zusätzlich der Spiegelpunkt  $C(y_1, x_2)$  von  $C(x_2, y_1)$  rechts oberhalb von  $C(x_1, y_2)$  liegt, d.h. wenn gilt  $y_1 > x_1 \wedge y_2 \leq x_2$ , so kann auch  $C(x_1, y_2)$  der Minimumpunkt  $m$  sein. Beachten Sie bitte, dass die letzte Bedingung  $y_2 \leq x_2$  entfallen kann, da sie nach Voraussetzung erfüllt ist, denn der rechte obere Eckpunkt  $C(x_2, y_2)$  soll unterhalb der 1. Winkelhalbierenden liegen! Vergleichen Sie dazu die Rechtecke  $Z_3, Z_8, Z_{14}$  bzw.  $Z_6, Z_9$ . Beachten Sie außerdem, dass bei  $Z_{14}$  die Bedingung  $y_1 > x_1$  nicht erfüllt sein kann, denn der linke untere Eckpunkt  $C(x_1, y_1)$  soll nach Voraussetzung unterhalb von  $y = x/\sqrt{3}$  liegen, d.h. es muss gelten  $y_1 < x_1/\sqrt{3} < x_1$ .

Mit Hilfe der obigen drei Hauptfälle werden alle typischen Lagen der vierzehn Rechtecke  $Z_\nu$  beschrieben, und der minimale Funktionswert über dem Punkt  $m$  lässt sich berechnet mit der Funktion

```
MpfrClass mIm_rz2_m (const MpfiClass& x, const MpfiClass& y),
```

die in `mpficlass.cpp` definiert ist. Der Rückgabewert ist der nahezu optimal abgerundete, minimale Funktionswert von  $f(x, y)$  über dem Rechteck  $Z_\nu$ , dessen oberer rechter Eckpunkt unterhalb der 1. Winkelhalbierenden liegen muss. Diese letzte Bedingung wird innerhalb der Funktion `mIm_rz2_m(...)` jedoch nicht überprüft und muss damit vorher abgesichert werden. Beachten Sie, dass bei der internen Auswertung von  $f(x, y)$  ein vorzeitiger Überlauf durch geeignete Skalierungen verhindert wird.



### B.4.3. Numerische Beispiele

Im **1. Beispiel** wählen wir  $X_1 = [5, 8]$  und  $Y_1 = [5, 6]$ , womit die typische Lage des zweidimensionalen Intervalls  $Z_1$  realisiert wird. Mit dem Funktionsaufruf  $T = \text{mIm\_rz2\_m}(X, Y)$ ; erhalten wir mit  $\text{prec} = 300$  für den Wertebereich  $W_f := \{f(x, y) \mid x \in X, y \in Y\} = [+6/625, +1/50]$  die nahezu optimale Einschließung  $W_f \subseteq T$ .

$$T = [9.59999999 \dots 99999 \cdot 10^{-3}, 2.00000 \dots 00001 \cdot 10^{-2}] \supset W_f$$

mit 90 ausgegeben Dezimalstellen. Rechnet man z.B. mit  $\text{prec} = 3\,000\,000$ , so erhält man das analoge Ergebnis mit  $3\,000\,000 \cdot \log_{10}(2) \approx 903\,090$  Dezimalstellen.

In der folgenden Tabelle sind für alle typischen, zweidimensionalen Intervalle  $Z_\nu$  ihre möglichen reellen Intervalle  $X_\nu, Y_\nu$  zusammen mit den jeweiligen exakten Funktionswerten  $f_m, f_M$  an den Minimum- bzw. Maximumstellen  $m, M$  zusammengestellt:

$\nu$	$X_\nu$	$Y_\nu$	$f_m$	$f_M$	$\nu$	$X_\nu$	$Y_\nu$	$f_m$	$f_M$
1	[5, 8]	[5, 6]	6/625	1/50	2	[4, 6]	[3.5, 4.5]	32/1875	448/12769
3	[2, 5]	[1, 4]	5/338	9/(32√3)	4	[6, 9]	[1, 2]	9/3362	3/200
5	[2, 6]	[4, 5]	60/3721	9/(128√3)	6	[2, 6]	[3, 5]	60/3721	12/169
7	[4, 5]	[2, 2.5]	20/841	9/(128√3)	8	[2, 4]	[1.5, 3]	192/5329	96/625
9	[0.5, 2]	[1, 2]	32/289	9√3/24	10	[0, 4]	[1, 2]	0	9√3/24
11	[1, 6]	[1.5, 3]	32/2601	48/169	12	[2, 7]	[1.5, 4]	336/42025	96/625
13	[2, 7]	[1, 4]	14/2500	9√3/96	14	[2, 5]	[1, 3]	5/338	9√3/96
15	[1, 4]	[2, 7]	14/2500	9√3/96	16	[-7, -2]	[1, 4]]	-9√3/96	-14/2500
17	[-4, 4]	[1, 2]	-9√3/96	+9√3/96	18	[-2, -0.5]	[-2, -1]	32/289	9√3/24

#### Anmerkungen:

- Die Rechtecke  $Z_\nu$ ,  $\nu = 1, 2, \dots, 14$ , liegen alle mit dem oberen rechten Eckpunkt nicht oberhalb der 1. Winkelhalbierenden, d.h. es gilt  $y_2 \leq x_2$ .
- $Z_{15}$  entsteht aus  $Z_{13}$  durch Spiegelung an der 1. Winkelhalbierenden, so dass  $f(x, y)$  über  $Z_{15}$  und  $Z_{13}$  den gleichen Wertebereich besitzt.
- $Z_{16}$  entsteht aus  $Z_{13}$  durch Multiplikation von  $X_{13}$  mit  $-1$ . Nach (B.14) ist dann der Wertebereich von  $f(x, y)$  über  $Z_{16}$  der negative Wertebereich von  $f(x, y)$  über  $Z_{13}$ .
- $Z_{17}$  entsteht aus  $Z_{10}$  durch Vereinigung von  $X_{10}$  mit  $-X_{10}$ , so dass das Minimum von  $f(x, y)$  über  $Z_{17}$  das negative Maximum von  $X_{10}$  sein muss.
- $Z_{18}$  entsteht aus  $Z_9$  durch Multiplikation von  $X_9$  und  $Y_9$  durch  $-1$ , so dass nach (B.13)  $f(x, y)$  über  $Z_9$  und  $Z_{18}$  den gleichen Wertebereich besitzt.

Der Algorithmus zur Berechnung der gesuchten Extrema  $f_m, f_M$  wurde u.a. mit allen obigen  $Z_\nu := \{(x, y) \mid x \in X_\nu \wedge y \in Y_\nu\}$ ,  $\nu = 1, 2, \dots, 18$ , getestet.

## B.5. $x^2 + a \cdot x + b$

Für die gegebenen reellen Maschinenintervalle  $X = [x_1, x_2] \ni x$ ,  $A = [a_1, a_2] \ni a$  und  $B = [b_1, b_2] \ni b$  liefert die Funktion

```
MpfiClass poly2(const MpfiClass& X, const MpfiClass& A, const MpfiClass& B)
```

für

$$f(x, a, b) = x^2 + a \cdot x + b; \quad x, a, b \in \mathbb{R}$$

eine optimale Maschineneinschließung  $f(X, A, B)$  der Wertemenge

$$W_{f,X,A,B} := \{f(x, a, b) \mid x \in X, a \in A, b \in B\} \subset f(X, A, B),$$

wobei die **optimale** Einschließung nur garantiert ist, wenn  $A = [a_1, a_2]$  ein Punktintervall<sup>1</sup> ist. Die optimale Einschließung ist gewährleistet, wenn  $f(x, a, b)$  mit folgendem Ausdruck

$$(B.17) \quad f(x, a, b) = (x + a/2)^2 - (a/2)^2 + b$$

intervallmäßig ausgewertet wird, denn jetzt kommen beide Variablen  $x, b$  explizit nur einmal vor und weil  $a$  zweimal auftritt, muss  $A = [a_1, a_2]$  ein Punktintervall sein. Um bei der Auswertung von B.17 einen vorzeitigen Über- oder Unterlauf zu vermeiden, wird in der Funktion

```
void poly2(MpfiClass& R, long int& k, const MpfiClass& X,  
          const MpfiClass& A, const MpfiClass& B);
```

durch geeignete Skalierungen erreicht, dass das Intervall  $2^k \cdot R$  die gewünschte optimale Einschließung  $f(X, A, B)$  liefert, wobei  $\text{Sup}(|R|)$  nur wenig kleiner als  $\text{MaxFloat}()$  ausfällt. Bei der Berechnung von  $R$  werden mögliche Auslöschungen durch schrittweise Vergrößerungen der Präzision **prec** vermieden. Dabei wird **prec** so lange erhöht, bis sich die jeweiligen  $R$ -Werte nicht mehr ändern. Beim Vergleich dieser  $R$ -Werte müssen diese natürlich jeweils **vorher** nach außen in die Ausgangspräzision gerundet werden. Beachten Sie, dass die berechneten Einschließungen natürlich nicht optimal sein können, wenn mögliche Auslöschungen nicht vermieden werden.

### B.5.1. Numerische Beispiele

Im **1. Beispiel** wählen wir die Intervalle  $X = [-3, 2]$ ,  $A = [-6, -6]$ ,  $B = [-1, 2]$  und erhalten mit den Zwischenergebnissen

$$X + A/2 = [-6, -1], \quad (X + A/2)^2 = [1, 36], \quad (X + A/2)^2 - (A/2)^2 = [-8, 27]$$

das exakte Ergebnis  $f(X, A, B) = W_{f,X,A,B} = [-9, 29]$ , und genau diese optimale Einschließung liefert für **prec**  $\geq 5$  der Funktionsaufruf  $F = \text{poly2}(X, A, B)$ ; mit  $F, X, A, B$  vom Typ `MpfiClass`. Für **prec** = 4 erhält man die gröbere Einschließung  $F = [-9, 30]$ , da bei dieser Präzision das exakte Intervall  $[-9, 29]$  nicht mehr darstellbar ist.

Im **2. Beispiel** wählen wir mit der Präzision **prec** = 300 die nachfolgenden Maschinenintervalle  $X = [\text{pred}(1), \text{pred}(1)] = [1 - 2^{-300}, 1 - 2^{-300}]$ ,  $A = [0, 0]$ ,  $B = [-1, -1]$  und werten damit die Funktion  $f(x, 0, -1) = x^2 - 1$  in der unmittelbaren Umgebung ihrer Nullstelle  $x_0 = +1$  aus, so dass erhebliche Auslöschungen durch eine hinreichend große interne Präzision zu vermeiden ist. Der Funktionsaufruf  $F = \text{poly2}(X, A, B)$ ; liefert die optimale Einschließung

$$T = [-9.81818693059545310 \dots 19576808 \cdot 10^{-91}, -9.81818693059545310 \dots 19576325 \cdot 10^{-91}]$$

mit 93 ausgegebenen Dezimalziffern, vergleichen Sie dazu das entsprechende Ergebnis auf Seite 191, das mit einem ganz anderen Algorithmus ebenfalls optimal berechnet wurde.

<sup>1</sup>Gilt  $a_2 = \text{succ}(a_1)$ , so ist die Einschließung nahezu optimal.

## B.6. $(1 + x^2 - y^2)/(4x^2y^2 + (1 + x^2 - y^2)^2)$

Für die vorgegebenen reellen Maschinenintervalle  $\mathbf{X} \ni x$  und  $\mathbf{Y} \ni y$  liefert die Intervallfunktion

`MpfiClass Re_r1pz2(const MpfiClass& X, const MpfiClass& Y)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x, y) = \frac{1 + x^2 - y^2}{(4x^2y^2 + (1 + x^2 - y^2)^2)} \quad \text{mit: } x \in \mathbf{X}, y \in \mathbf{Y} \text{ und } (x, y) \neq (0, \pm 1).$$

Die Berechnung von  $f(x, y)$  wird ausführlich beschrieben auf Seite 189.  $f(x, y)$  ist als Realteilfunktion der holomorphen Funktion  $1/(1 + z^2)$ ,  $z \in \mathbb{C} \setminus \{0 \pm i\}$ , eine harmonische Funktion, wenn die singulären Punkte  $P(0, \pm 1)$  nicht Element des zweidimensionalen reellen Intervalls  $Z$  sind.

$$(B.18) \quad Z := \{(x, y) \mid x \in \mathbf{X} \wedge y \in \mathbf{Y}\}, \quad (x, y) \neq (0, \pm 1).$$

$f(x, y)$  nimmt daher für alle Punkte  $(x, y) \in Z$  als harmonische Funktion ihre Extrema auf dem Rand von  $Z$  an, vgl. dazu auch die Seite 244. Die Berechnung dieser Extremstellen,  $m(x, y)$  für das Minimum und  $M(x, y)$  für das Maximum, kann wesentlich vereinfacht werden, wenn man  $f(x, y) \equiv f(|x|, |y|)$  beachtet und die Eingangsintervalle  $\mathbf{X}, \mathbf{Y}$  durch die jeweiligen Intervalle, z.B.  $X = \text{abs}(\mathbf{X}) = \{|x| \mid x \in \mathbf{X}\}$ , seiner Absolutbeträge ersetzt. Damit reduziert sich die Berechnung der Extremstellen  $m(x, y), M(x, y)$  auf ein Intervall  $Z := \{(x, y) \mid x \in X \wedge y \in Y\}$ ,  $x \geq 0, y \geq 0$ , das nur im 1. Quadranten liegt und den Punkt  $(0, 1)$  nicht enthalten darf.  $X = [x_1, x_2]$ ,  $Y = [y_1, y_2]$ .

Zur Berechnung der Extremstellen  $m, M$  benötigt man zunächst die partiellen Ableitungen:

$$(B.19) \quad \frac{\partial f(x, y)}{\partial x} = - \frac{2x \cdot (1 + 2x^2 + x^4 + 2y^2 - 2x^2y^2 - 3y^4)}{(1 + x^2 - 2y + y^2)^2 \cdot (1 + x^2 + 2y + y^2)^2},$$

$$(B.20) \quad \frac{\partial f(x, y)}{\partial y} = - \frac{2y \cdot (-1 + 2x^2 + 3x^4 + 2y^2 + 2x^2y^2 - y^4)}{(1 + x^2 - 2y + y^2)^2 \cdot (1 + x^2 + 2y + y^2)^2}.$$

Die Extremalkurven für  $\partial f/\partial x$  ergeben sich aus der Forderung  $\partial f/\partial x = 0$ , siehe Seite 247. Die erste ist wegen des Faktors  $x = 0$  die positive  $y$ -Achse, und die zweite Extremalkurve ergibt sich aus der Forderung  $1 + 2x^2 + x^4 + 2y^2 - 2x^2y^2 - 3y^4 = 0$ , aus der für den 1. Quadranten folgt<sup>2</sup>

$$y = h_1(x) = \sqrt{(1 - x^2 + 2 \cdot \sqrt{1 + x^2 + x^4})/3}, \quad 0 \leq x \leq 16,$$

$$h_1(x) = x \cdot \sqrt{(u - 1 + 2 \cdot \sqrt{1 + u + u^2})/3}, \quad u := (1/x)^2, \quad x > 16.$$

Die Extremalkurven für  $\partial f/\partial y$  ergeben sich analog aus der Forderung  $\partial f/\partial y = 0$ . Die erste ist wegen des Faktors  $y = 0$  die positive  $x$ -Achse und die zweite und dritte Extremalkurve ergibt sich aus  $-1 + 2x^2 + 3x^4 + 2y^2 + 2x^2y^2 - y^4 = 0$ , und für den 1. Quadranten folgt wie oben

$$(B.21) \quad y = h_2(x) := \sqrt{1 + x^2 + 2x \cdot \sqrt{1 + x^2}}, \quad 0 \leq x \leq 16,$$

$$(B.22) \quad := x \cdot \sqrt{1 + u + 2 \cdot \sqrt{1 + u}}, \quad u := (1/x)^2, \quad x > 16.$$

$$(B.23) \quad y = h_3(x) := \sqrt{1 + x^2 - 2x \cdot \sqrt{1 + x^2}}, \quad 0 \leq x \leq 1/\sqrt{3}.$$

Für die Bestimmung der Extremalpunkte  $m, M$  ist es noch wichtig zu wissen, in welchen Bereichen des 1. Quadranten  $f(x, y)$  positiv oder negativ ist. Für  $y > h_4(x) := \sqrt{1 + x^2}$  gilt  $f(x, y) < 0$  und für  $y < h_4(x)$  gilt  $f(x, y) > 0$ .

Die Funktionen  $h_i(x), i = 1, 2, 3, 4$  besitzen den gemeinsamen Punkt  $P(0, 1)$  und erfüllen die Ungleichungen

$$(B.24) \quad h_3(x) \leq h_1(x), \text{ falls } 0 \leq x \leq 1/\sqrt{3}; \quad h_1(x) \leq h_4(x) \leq h_2(x), \text{ falls } x \geq 0.$$

Der Beweis der Ungleichungen ist trivial und bleibt dem Leser überlassen.

<sup>2</sup> $h_1(x)$  wurde berechnet mit *Mathematica* und der `Solve[...]`-Funktion.

In der folgenden Abbildung sind die Funktionen  $h_1(x)$ ,  $h_2(x)$ ,  $h_3(x)$ ,  $h_4(x)$  im 1. Quadranten entsprechend farbig und leicht verzerrt dargestellt.

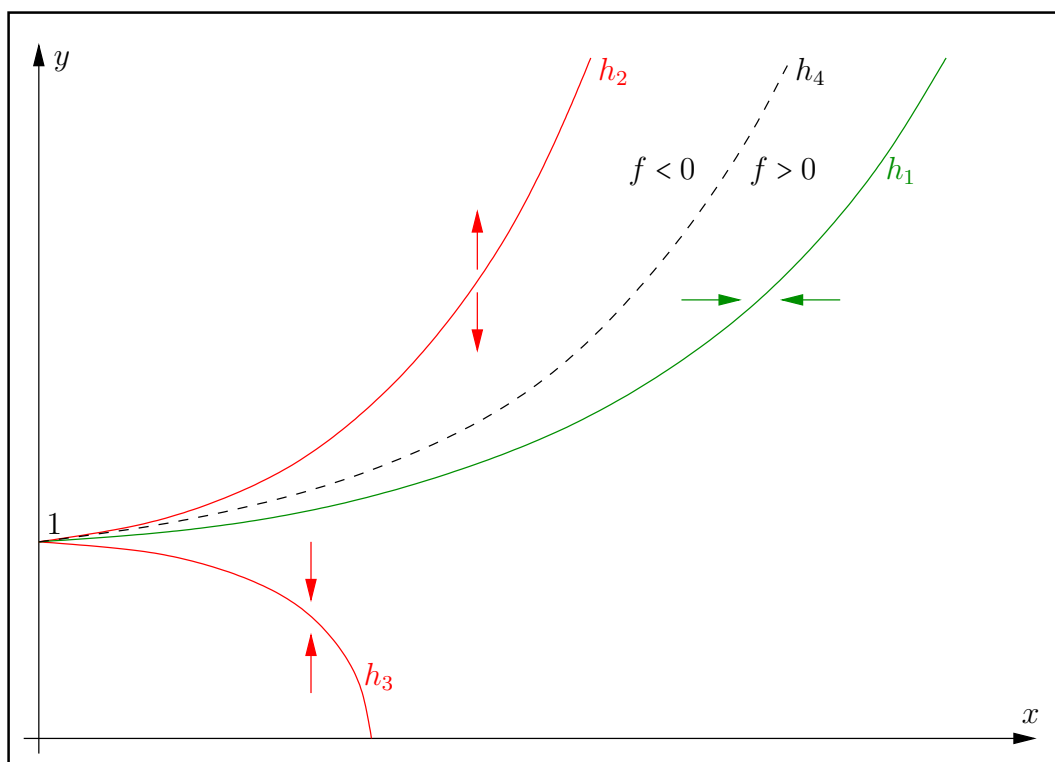


Abbildung B.17.: Extremalkurven und Richtungspfeile für wachsende  $f(x, y)$ -Werte

Die roten Pfeile geben parallel zur  $y$ -Achse und die grünen Pfeile parallel zur  $x$ -Achse die Richtung wachsender  $f(x, y)$ -Werte an.

### B.6.1. Maximumbestimmung

Vergleicht man die Extremalkurven der obigen Abbildung z.B. mit denen aus Abbildung B.4, so erkennt man, dass sich die Extremalkurven von  $y = h_2(x)$  und  $y = \sqrt{3} \cdot x$  bzw. von  $y = h_1(x)$  und  $y = x/\sqrt{3}$  genau entsprechen, da sie jeweils im gemeinsamen Punkt  $(0, 0)$  bzw.  $(0, +1)$  starten, monoton sind und die gleiche gegenseitige Lage zueinander besitzen. Außerdem stimmen auch die jeweiligen Pfeile, welche die Richtungen wachsender Funktionswerte von  $f(x, y)$  angeben, genau überein. Zusätzlich liegen auch die monoton verlaufenden Nulllinien, d.h. die Graphen der Funktionen  $y = x$  bzw.  $y = h_4(x) = \sqrt{1 + x^2}$  stets zwischen diesen Extremalkurven.

Prinzipiell **neu** ist in Abbildung B.17 jedoch die zusätzliche Extremalkurve  $y = h_3(x)$ , die im Definitionsbereich  $0 \leq x \leq 1/\sqrt{3}$  streng monoton fällt. Wenn wir daher zur Bestimmung von  $M$  den Algorithmus von Seite 213 sinngemäß übernehmen wollen, so muss in einem vorgeschalteten Punkt 0. untersucht werden, wie der Extrempunkt  $M$  zu bestimmen ist, wenn der untere linke Eckpunkt des Rechtecks  $Z$  unterhalb von  $h_3(x)$  liegt. In der folgenden Abbildung wird dieser Fall daher genauer untersucht.

Wir müssen noch zeigen, dass die Nachweise bei den Anmerkungen auf Seite 213 auch mit unserer Funktion  $f(x, y)$  und den entsprechenden Extremalkurven geführt werden können. Nach **Punkt 1.** ist unter den Voraussetzungen  $x_2 > x_1$ ,  $y_1 > \sqrt{1 + x_1^2}$  zu zeigen:

$$f(x_1, y_1) < f(x_2, h_2(x_2)) = \frac{-1}{4x_2 \cdot (x_2 + \sqrt{1 + x_2^2})}.$$

Wegen  $1 \leq h_2(x_1) \leq y_1 \leq h_2(x_2)$  gilt die Doppelungleichung

$$1 + x_1^2 + 2x_1 \cdot \sqrt{1 + x_1^2} \leq y_1^2 \leq 1 + x_2^2 + 2x_2 \cdot \sqrt{1 + x_2^2},$$

und damit erhält man nach einfachen Umrechnungen die Abschätzung

$$f(x_1, y_1) = \frac{1 + x_1^2 - y_1^2}{4x_1^2y_1^2 + (1 + x_1^2 - y_1^2)^2} \leq \frac{-1}{4x_1(x_1 + \sqrt{1 + x_1^2})}.$$

Wegen  $x_2 > x_1$  folgt daraus unmittelbar die Behauptung ■ Die Punkte 2. und 3. auf Seite 213 sind auch hier trivial und der Nachweis unter Punkt 4. bleibt dem Leser überlassen.

In der folgenden Abbildungen zeigen wir unter der Voraussetzung *Der linke untere Eckpunkt von  $Z$  liegt unterhalb oder auf dem Graphen von  $h_3(x)$*  drei typische Lagen eines Rechteckintervalls  $Z$ .

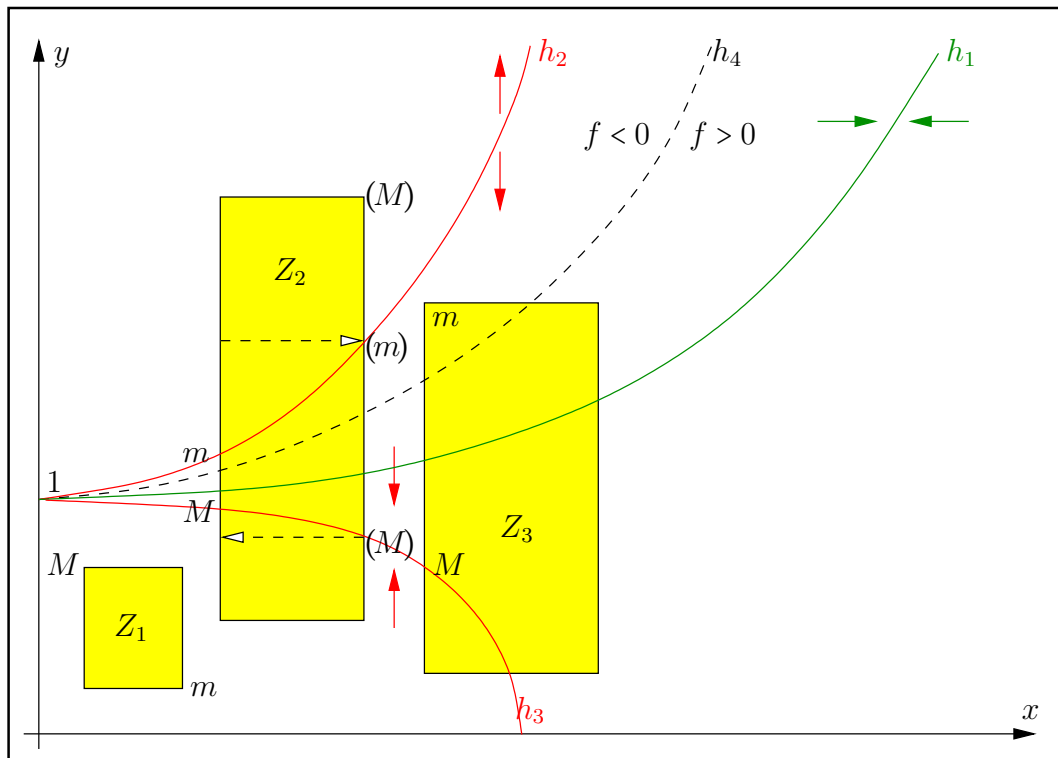


Abbildung B.18.: Linke untere Ecke von  $Z$  unterhalb von  $h_3(x)$

### Anmerkungen:

1. In  $Z_2$  kann der Eckpunkt  $(M)$  oben rechts kein Maximumpunkt sein, da er im Bereich  $f(x, y) < 0$  liegt. Der Schnittpunkt  $(M)$  mit  $h_3(x)$  kann ebenfalls kein Maximumpunkt sein, da man längs des punktierten Pfeils in Richtung wachsender Funktionswerte nach  $M$  kommt, und ganz entsprechend kann  $(m)$  kein Minimumpunkt sein.
2. Bewegt man bei  $Z_2$  und  $Z_3$  die obere Parallele zu  $x$ -Achse so, dass der obere rechte Eckpunkt immer oberhalb oder auf dem Graphen von  $h_3$  liegt, so befindet sich  $M$  stets auf der linken Parallelen zur  $y$ -Achse. Schneidet diese Parallele den Graphen von  $h_3(x)$ , so ist  $M$  dieser Schnittpunkt, sonst ist  $M$  die linke obere Ecke von  $Z$ . Der Leser möge dies anhand aller möglichen, prinzipiellen Lagen von  $Z$  überprüfen!
3. Wenn also der untere linke Eckpunkt von  $Z$  unterhalb oder auf dem Graphen von  $h_3(x)$  liegt, so liegt  $M$  stets auf der linken Parallelen zur  $y$ -Achse. Schneidet diese Parallele den Graphen von  $h_3(x)$ , so ist  $M$  dieser Schnittpunkt, sonst ist  $M$  der obere linke Eckpunkt von  $Z$ .

Wir formulieren jetzt noch die Bedingungen aus dem letzten Punkt 3. in der Sprache der Intervallarithmetik. Dazu benutzen wir die nur für interne Anwendungen definierte Funktion

```
MpfiClass Re_r1pz2_H3(const MpfrClass& x)
```

welche für eine Maschinenzahl  $0 \leq x < 1/\sqrt{3}$  den Funktionswert  $h_3(x)$  optimal einschließt. Eine ebenfalls optimale Einschließung von  $1/\sqrt{3}$  liefert mit  $x = 1$  und  $n = 3$  die Funktion

```
MpfiClass xdivsqrtn (const MPFR::MpfrClass& x, unsigned long int n)
```

welche in `mpfiClass.cpp` definiert ist.

Die Bedingung *Die linke untere Ecke von  $Z$  liegt unterhalb oder auf dem Graphen von  $h_3(x)$*  ist definiert durch:

1.  $x_1 < 1/\sqrt{3}$  und wird realisiert durch:  $x_1 \leq \text{Inf}(\text{xdivsqrtn}(1,3))$ ,
2.  $y_1 \leq h_3(x_1)$  und wird realisiert durch:  $y_1 \leq \text{Inf}(\text{Re\_r1pz2\_H3}(x_1))$ .

Wenn oben  $y_1 = \text{Inf}(\text{Re\_r1pz2\_H3}(x_1))$  erfüllt ist, so kann gelten  $y_1 = h_3(x_1)$ , aber auch dann ist  $M$  der Schnittpunkt der linken Parallelen von  $Z$  mit dem Graphen von  $h_3(x)$ .

Die Bedingung *Die linke Parallele von  $Z$  schneidet den Graphen von  $h_3(x)$*  ist definiert durch:

$$y_1 \leq h_3(x) \leq y_2 \quad \text{und wird realisiert durch: } h_3([x_1]) \leq [y_1, y_2],$$

wobei  $h_3([x_1])$  die mit `MpfiClass Re_r1pz2_H3(const MpfrClass& x)` berechnete **optimale** Einschließung von  $h_3(x_1)$  bedeutet.

Zur Bestimmung von  $M$  benötigen wir zusätzlich die internen Funktionen

```
MpfiClass Re_r1pz2_H1(const MpfrClass& x)
MpfiClass Re_r1pz2_H2(const MpfrClass& x)
MpfiClass Re_r1pz2_H4(const MpfrClass& x)
```

die für eine Maschinenzahl  $x$  **optimale** Einschließungen der Funktionswerte  $h_1(x), h_2(x), h_4(x)$  liefern.

### Anmerkung:

Um sicherzugehen, dass z.B. für Maschinenzahlen  $x_1, y_1$  die Ungleichung  $y_1 \leq h_3(x_1)$  erfüllt ist, wenn  $y_1 \leq \text{Inf}(\text{Re\_r1pz2\_H3}(x_1))$  gilt, muss vorausgesetzt werden, dass mit `Re_r1pz2_H3(x_1)` eine **optimale** Einschließung von  $h_3(x_1)$  berechnet wird. Der Leser möge mit einem einfachen Beispiel zeigen, dass die **optimale** Einschließung wirklich vorausgesetzt werden muss. Die obigen fünf Funktionen liefern stets die notwendigen optimalen Einschließungen, wenn für Maschinenzahlen  $x$  gilt:

$$x > 0; \quad h_i(x) \text{ ist keine Maschinenzahl;} \quad \text{Eine Fehlermeldung tritt nicht auf.}$$

Falls jedoch die Funktionswerte  $h_i(x), i = 1, 2, 3$  für Maschinenzahlen  $x$  selbst wieder Maschinenzahlen sind, so kann der Fall eintreten, dass  $h_i(x)$  nicht durch ein Punktintervall sondern durch ein echtes Intervall eingeschlossen wird, das jedoch in seinem Innern keine Maschinenzahl enthält. In diesem Fall kann z.B. die Ungleichung  $y_1 \leq h_3(x_1)$  durch die Maschinenabfrage

$$y_1 \leq \text{Inf}(\text{Re\_r1pz2\_H3}(x_1))$$

nicht mehr garantiert werden. Mir ist jedoch kein Fall bekannt, in dem  $h_i(x), i = 1, 2, 3$  für  $x > 0$  wirklich eine Maschinenzahl ist. Für  $h_4(x) = \sqrt{1+x^2}$  gibt es natürlich Beispiele, bei denen für Maschinenzahlen  $x > 0$  die Funktionswerte  $h_4(x)$  selbst wieder Maschinenzahlen sind, z.B. gilt  $\sqrt{1+(3/4)^2} = 5/4$ . Die Funktion `Re_r1pz2_H4(x)` liefert jedoch für diese ersten Spezialfälle einschließende **Punktintervalle**. Diese Spezialfälle berechnen sich aus der speziellen Pythagoräischen Gleichung

$$(2^{n+1})^2 + (2^{2(n-p)} - 2^{2p})^2 = (2^{2(n-p)} + 2^{2p})^2, \quad n \in \mathbb{N}, \quad p \in \mathbb{N}_0, \quad p = 0, 1, \dots, n/2 + 1,$$

wobei  $x$  gegeben ist durch:  $x = (2^{2(n-p)} - 2^{2p})/2^{n+1}$ , mit z.B.  $x = 3/4$  für  $n = 1$  und  $p = 0$ .

Um die maximalen oder minimalen Funktionswerte  $f(x, y)$  an den Schnittpunkten der Rechteckseiten von  $Z$  mit den Extremalkurven  $h_1(x), h_2(x), h_3(x)$  einschließen zu können, benötigen wir noch die folgenden Funktionen

$$f_{h_1}(x) := f(x, h_1(x)) = \begin{cases} \frac{3}{4 \cdot (2 + x^2 + \sqrt{1 + x^2 + x^4})}, & 0 \leq x \leq 8, \\ u \cdot \frac{3}{4 \cdot (1 + 2u + \sqrt{1 + u + u^2})}, & u = (1/x)^2, \quad x > 8. \end{cases}$$

$$f_{h_2}(x) := f(x, h_2(x)) = \begin{cases} \frac{-1}{4x \cdot (\sqrt{1 + x^2} + x)}, & 0 \leq x \leq 8, \\ u \cdot \frac{-1}{4 \cdot (1 + \sqrt{1 + u})}, & u = (1/x)^2, \quad x > 8. \end{cases}$$

$$f_{h_3}(x) := f(x, h_3(x)) = \frac{\sqrt{1 + x^2} + x}{4x}, \quad 0 < x \leq 1\sqrt{3}.$$

deren Funktionswerte mithilfe der internen Funktionen

```
MpfiClass Re_r1pz2_fH1(const MpfrClass& x)
MpfiClass Re_r1pz2_fH2(const MpfrClass& x)
MpfiClass Re_r1pz2_fH3(const MpfrClass& x)
```

eingeschlossen werden. Beachten Sie dabei, dass diese Einschließungen nicht unbedingt optimal sondern nur nahezu optimal sein müssen. Die obigen Funktionsterme  $f_{hi}(x)$  wurden mit dem Algebrasystem *Mathematica* und `FullSimplify[...]` so berechnet, dass bei der Auswertung im 1. Quadranten wegen  $x \geq 0$  keine Auslöschungen auftreten können.

Um analog zu Punkt 2. auf Seite 213 den Funktionswert über dem Schnittpunkt der unteren Parallelen von  $Z$  mit dem Graphen der Extremalkurve  $h_1(x)$  zu berechnen, gilt zunächst für den Schnittpunkt  $M(x_S, y_S) = M(x_S, y_1)$ , wobei  $y_1$  und  $x_S$  gegeben sind durch<sup>3</sup>

$$y_1 = h_1(x_S) = \sqrt{(1 - x_S^2 + 2 \cdot \sqrt{1 + x_S^2 + x_S^4})/3} \quad \rightsquigarrow \quad x_S = \sqrt{y_1^2 - 1 + 2y_1 \cdot \sqrt{y_1^2 - 1}}.$$

Der Funktionswert über  $M(x_S, y_1)$  ergibt sich daraus mit  $y = y_1$  zu

$$f(x_S, y_1) = h_5(y) := \frac{1}{4y \cdot (y + \sqrt{y^2 - 1})} = u \cdot \frac{1}{4 \cdot (1 + \sqrt{1 - u})}, \quad u := (1/y)^2 \leq 1,$$

wobei  $y^2 \geq 1$  automatisch erfüllt ist. Die Funktionswerte von  $h_5(x)$  werden für Maschinenzahlen  $x \geq 1$  nahezu optimal eingeschlossen mit Hilfe der internen Funktion

```
MpfiClass Re_r1pz2_H5(const MpfrClass& x);
```

## B.6.2. Minimumbestimmung

Analog zur Maximumbestimmung können wir auch jetzt die Lage der Minimumpunkte  $m$  analog zum Fall  $1/z^2$  bestimmen, wenn alle Rechtecke  $Z$  stets oberhalb des Graphen der Extremalkurve  $y = h_3(x)$  liegen, und dies gilt auch dann noch, wenn die untere Parallele von  $Z$  den Graphen der Extremalkurve  $y = h_3(x)$  schneidet. In der folgenden Abbildung wird dies am Beispiel  $Z_0$  verdeutlicht, das dem Rechteck  $Z_{13}$  in Abbildung B.6 entspricht, wobei der Minimumpunkt  $m$  der linke obere Eckpunkt des jeweiligen Rechtecks ist. Der Leser möge diesen Sachverhalt für alle restlichen typischen Lagen von  $Z$  bestätigen, wenn die untere Parallele von  $Z$  den Graphen der Extremalkurve  $y = h_3(x)$  schneidet.

Zur Bestimmung der Lage von  $m$  muss daher nur noch der Fall gesondert behandelt werden, dass der rechte untere Eckpunkt von  $Z$  unterhalb des Graphen von  $y = h_3(x)$  liegt. In diesem

<sup>3</sup>Die folgenden Umformungen wurden mit *Mathematica* und der Funktion `Solve[...]` durchgeführt.

Fall genügt es, den Minimumpunkt  $m$  auch wieder analog zum Fall  $1/z^2$  zu bestimmen und den Funktionswert  $f(x, y)$  über diesem Punkt mit dem Funktionswert über dem rechten unteren Eckpunkt von  $Z$  zu vergleichen. In der nachfolgenden Abbildung wird dies mit dem Beispiel  $Z_1$  gezeigt. Der Leser möge anhand weiterer typischer Beispiele für die Lage von  $Z$  zeigen, dass so der minimale Funktionswert  $f(x, y)$  über  $Z$  stets korrekt berechnet wird.

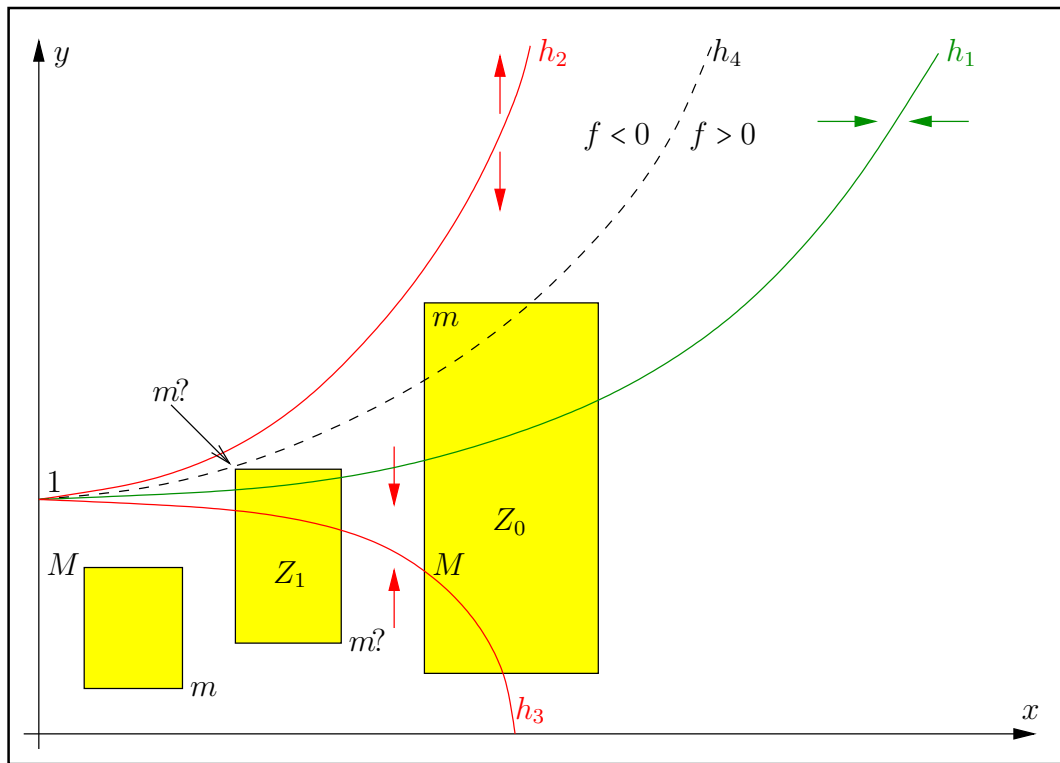


Abbildung B.19.: Zur Lage der Minimumpunkte  $m$

### B.6.3. Numerische Beispiele

Im **1. Beispiel** wählen wir  $Z_1$  aus Abb. B.19 mit  $X = [1/4, 1/2]$ ,  $Y = [1/4, 5/4]$ . Mit einfachen Rechnungen erhält man  $f_m = f(1/4, 5/4) = -32/41 = -0.78048\dots$  und  $f_M = f(1/4, h_3(1/4)) = (1 + \sqrt{17})/4 = 1.2807764\dots$ . Mit dem Funktionsaufruf  $T = \text{Re\_r1pz2}(X, Y)$ ; erhalten wir mit  $\text{prec} = 300$  für den Wertebereich  $W_f := \{f(x, y) \mid x \in X, y \in Y\} = [f_m, f_M]$  die nahezu optimale Einschließung  $W_f \subseteq T$ .

$$T = [-7.80487804878048\dots78049 \cdot 10^{-1}, +1.2807764064044\dots96592] \supset W_f$$

mit 90 ausgegeben Dezimalstellen.

Im **2. Beispiel** wählen wir  $X = [2, 2]$ ,  $Y = [3, 3]$  und erhalten mit  $T = \text{Re\_r1pz2}(X, Y)$  für den Wertebereich  $W_f := \{f(x, y) \mid x \in X, y \in Y\} = [f_m, f_M]$  die nahezu optimale Einschließung  $W_f \subseteq T$ .

$$T = [-2.5000000000000\dots00001 \cdot 10^{-2}, -2.4999999999999\dots99999 \cdot 10^{-2}] \supset W_f$$

mit 90 ausgegeben Dezimalstellen.



### B.7. $2xy/(4x^2y^2 + (1 + x^2 - y^2)^2)$

Für die vorgegebenen reellen Maschinenintervalle  $X \ni x$  und  $Y \ni y$  liefert die Intervallfunktion

`MpfiClass Im_r1pz2(const MpfiClass& X, const MpfiClass& Y)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x, y) = \frac{2xy}{4x^2y^2 + (1 + x^2 - y^2)^2} \quad \text{mit: } x \in X, y \in Y \text{ und } (x, y) \neq (0, \pm 1).$$

Die Berechnung von  $f(x, y)$  wird ausführlich beschrieben auf Seite 188.  $-f(x, y)$  ist als Imaginärteil der holomorphen Funktion  $1/(1 + z^2)$ , mit  $z \in \mathbb{C} \setminus \{0 \pm i\}$ , eine harmonische Funktion, wenn die singulären Punkte  $P(0, \pm 1)$  nicht Element des zweidimensionalen reellen Intervalls  $Z$  sind.

$$(B.25) \quad Z := \{(x, y) \mid x \in X \wedge y \in Y\}, \quad (x, y) \neq (0, \pm 1).$$

$f(x, y)$  nimmt daher für alle Punkte  $(x, y) \in Z$  als harmonische Funktion ihre Extrema auf dem Rand von  $Z$  an, vgl. dazu auch die Seite 244. Um die zugehörigen Extremstellen,  $m$  für das Minimum und  $M$  für das Maximum, möglichst einfach bestimmen zu können, betrachten wir zunächst die folgenden Symmetrieeigenschaften:

$$(B.26) \quad f(-x, -y) \equiv f(y, x), \quad \text{Spiegelung am Ursprung,}$$

$$(B.27) \quad f(-x, y) \equiv -f(x, y), \quad \text{Spiegelung an der } y\text{-Achse,}$$

$$(B.28) \quad f(x, -y) \equiv -f(x, y), \quad \text{Spiegelung an der } x\text{-Achse.}$$

Mit den obigen drei Transformationen kann die Bestimmung der Extrempunkte  $m, M$  auf den 1. Quadranten reduziert werden. Wir betrachten dazu ein typisches Beispiel, wobei das gegebene Rechteck  $Z = \{(x, y) \mid x \in X = [-1, 3] \wedge y \in Y = [0.321, 1.10]\} = Z1 \cup Z2$  in der oberen Halbebene liegt.

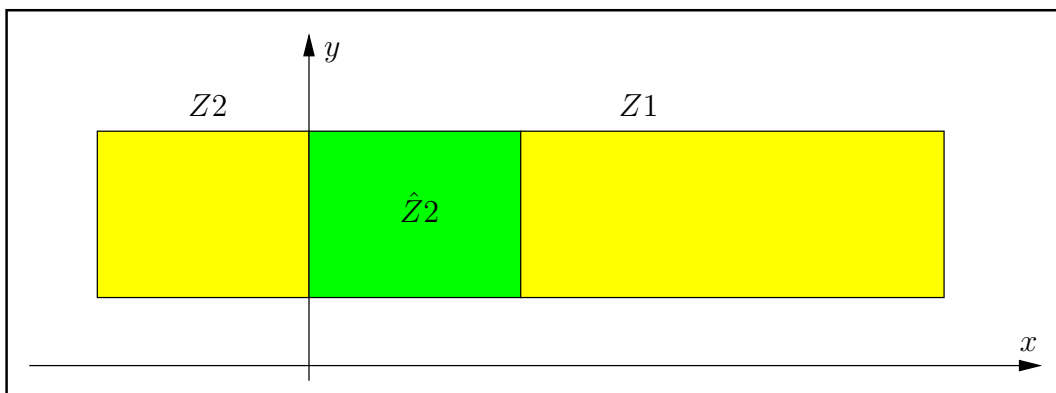


Abbildung B.20.: Transformation in den 1. Quadranten

$Z = Z1 \cup Z2$  wird zuerst in die beiden Intervalle  $Z1, Z2$  aus dem 1. und 2. Quadranten zerlegt. Da  $f(x, y)$  in  $Z2$  nur negative Funktionswerte besitzt, kann  $f(x, y)$  über ganz  $Z$  sein Minimum nur über  $Z2$  annehmen. Ist dann  $\hat{Z}2$  das an der  $y$ -Achse gespiegelte Rechteck  $Z2$  und ist  $f_M$  das Maximum über diesem Rechteck, so ist  $f_m := -f_M$  das Minimum über ganz  $Z$ . Da  $f(x, y)$  in  $Z1$  nur positive Funktionswerte besitzt, ist das Maximum  $f_M$  über  $Z1$  auch das Maximum über dem ganzen Rechteck  $Z$ . Zur Berechnung von  $f_m, f_M$  wurden damit Maximum- und Minimumberechnungen nur über Rechtecken im 1. Quadranten vorgenommen. Ganz analog können auch für alle anderen Lagen von  $Z$  die Berechnungen stets nur im 1. Quadranten durchgeführt werden. Weitere Einzelheiten findet man in der Funktion

`MpfiClass Im_r1pz2(const MpfiClass& X, const MpfiClass& Y)`

die in `mpficlass.cpp` definiert ist.

Zur Berechnung der Extremstellen  $m, M$  von Rechtecken  $Z$  im 1. Quadranten benötigt man zunächst die partiellen Ableitungen:

$$(B.29) \quad \frac{\partial f(x, y)}{\partial x} = \frac{2y \cdot (1 - 2x^2 + y^4 - 2y^2 - 2x^2y^2 - 3x^4)}{(1 + x^2 - 2y + y^2)^2 \cdot (1 + x^2 + 2y + y^2)^2},$$

$$(B.30) \quad \frac{\partial f(x, y)}{\partial y} = \frac{2x \cdot (1 + 2x^2 + x^4 + 2y^2 - 2x^2y^2 - 3y^4)}{(1 + x^2 - 2y + y^2)^2 \cdot (1 + x^2 + 2y + y^2)^2}.$$

Die Extremalkurven für  $\partial f/\partial x$  ergeben sich aus der Forderung  $\partial f/\partial x = 0$ , siehe Seite 247. Die erste ist wegen des Faktors  $y = 0$  die positive  $x$ -Achse, und die zweite und dritte Extremalkurve ergibt sich aus  $1 - 2x^2 + y^4 - 2y^2 - 2x^2y^2 - 3x^4 = 0$ . Für den 1. Quadranten folgt daraus<sup>4</sup>

$$(B.31) \quad y = h_2(x) := \sqrt{1 + x^2 + 2x \cdot \sqrt{1 + x^2}}, \quad 0 \leq x \leq 16,$$

$$(B.32) \quad := x \cdot \sqrt{1 + u + 2 \cdot \sqrt{1 + u}}, \quad u := (1/x)^2, \quad x > 16.$$

$$(B.33) \quad y = h_3(x) := \sqrt{1 + x^2 - 2x \cdot \sqrt{1 + x^2}}, \quad 0 \leq x \leq 1/\sqrt{3}.$$

Die Extremalkurven für  $\partial f/\partial y$  ergeben sich analog aus der Forderung  $\partial f/\partial y = 0$ . Die erste ist wegen des Faktors  $x = 0$  die positive  $y$ -Achse und die zweite ergibt sich aus der Forderung  $1 + 2x^2 + x^4 + 2y^2 - 2x^2y^2 - 3y^4 = 0$ , aus der für den 1. Quadranten wie oben folgt

$$(B.34) \quad y = h_1(x) := \sqrt{(1 - x^2 + 2 \cdot \sqrt{1 + x^2 + x^4})/3}, \quad 0 \leq x \leq 16,$$

$$(B.35) \quad := x \cdot \sqrt{(u - 1 + 2 \cdot \sqrt{1 + u + u^2})/3}, \quad u := (1/x)^2, \quad x > 16.$$

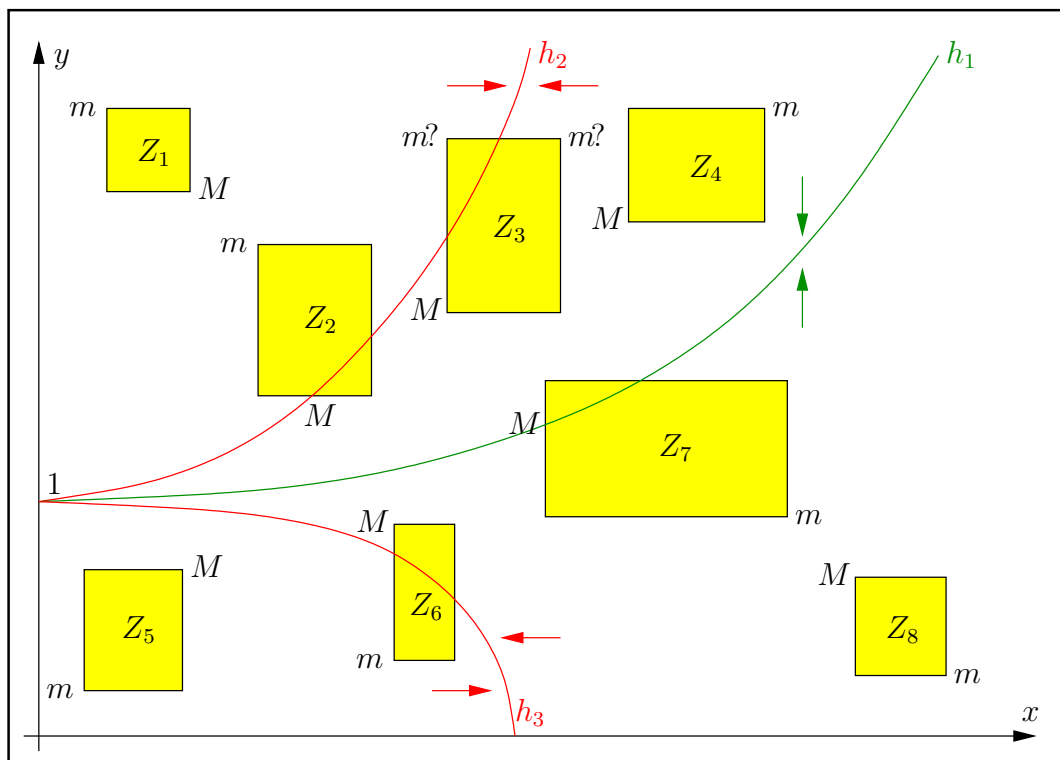


Abbildung B.21.: Rechtecke  $Z_i$  und Extremalkurven  $h_1(x), h_2(x), h_3(x)$  im 1. Quadranten

In den folgenden Abbildungen sind weitere typische Lagen der Rechtecke  $Z$  angegeben. Dabei bedeuten z.B. zwei oder auch mehrere Symbole  $m?$ , dass jeweils hier ein Minimum vorliegen kann, weshalb entsprechend mehrere Rechnungen und Vergleiche notwendig werden.

<sup>4</sup> $h_1(x)$  wurde berechnet mit *Mathematica* und der `Solve[...]`-Funktion.

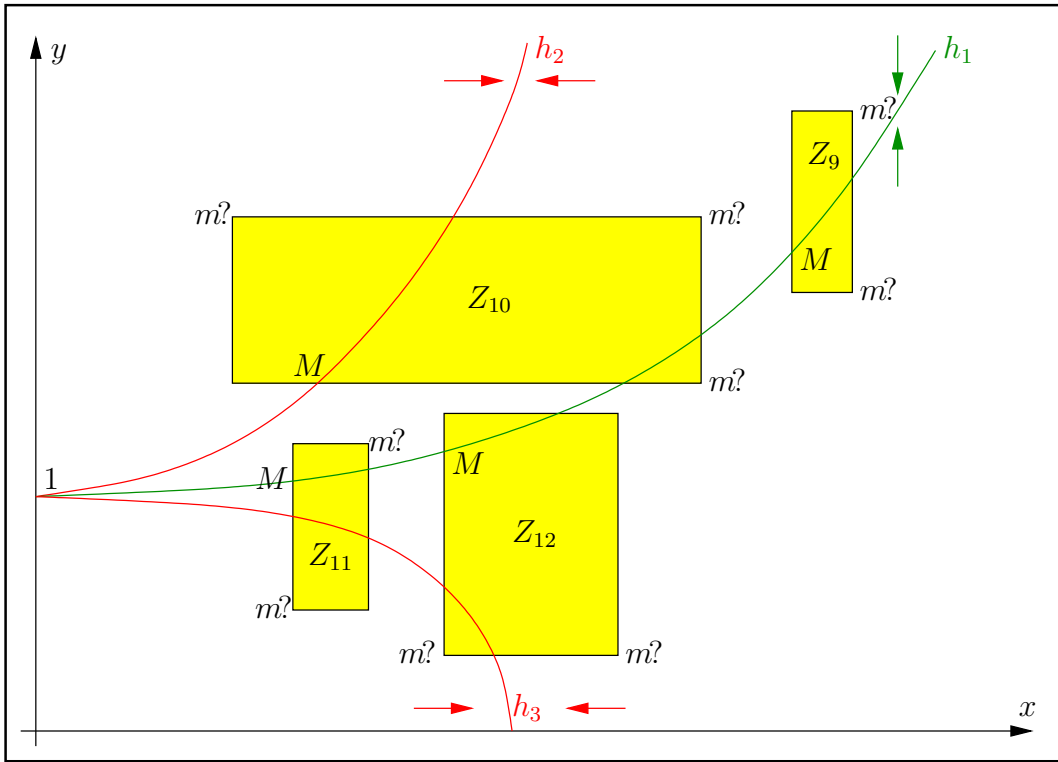


Abbildung B.22.: Rechtecke  $Z_i$  und Extremalkurven  $h_1(x), h_2(x), h_3(x)$  im 1. Quadranten

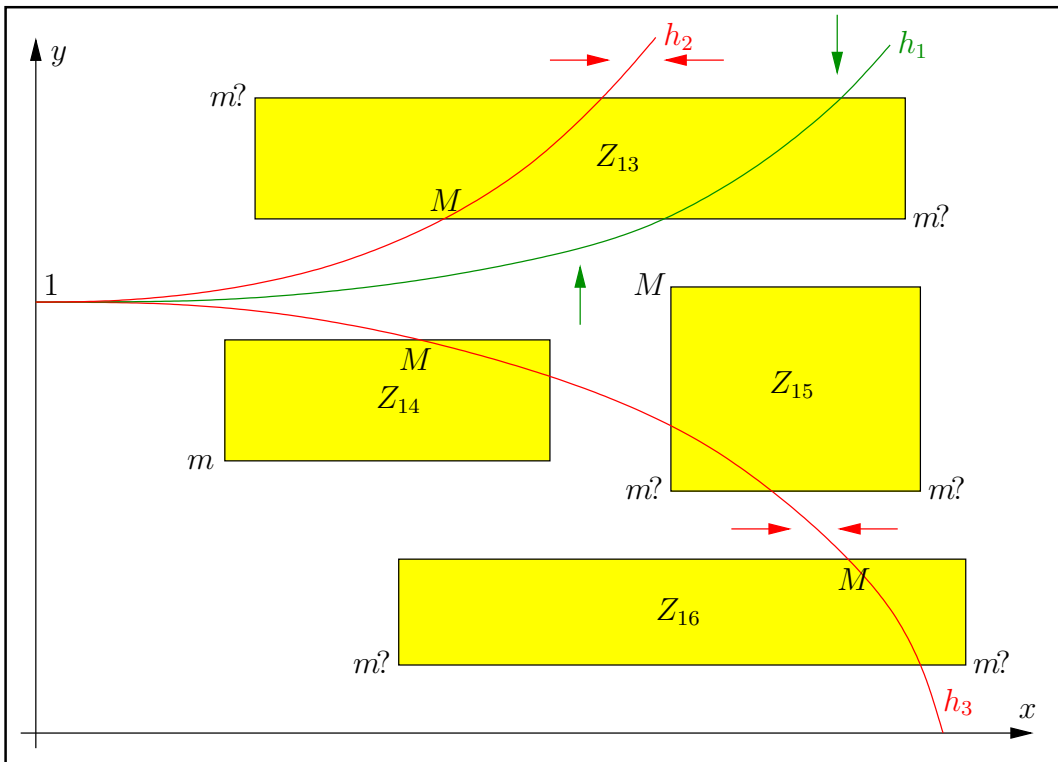


Abbildung B.23.: Rechtecke  $Z_i$  und Extremalkurven  $h_1(x), h_2(x), h_3(x)$  im 1. Quadranten

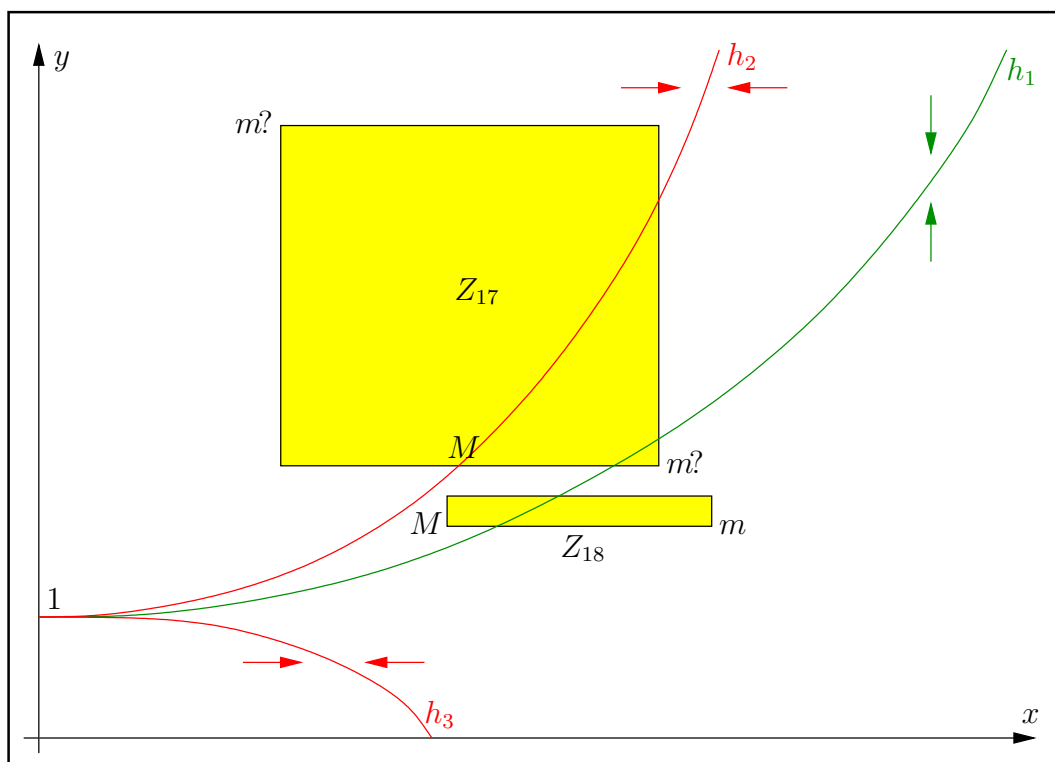


Abbildung B.24.: Rechtecke  $Z_i$  und Extremalkurven  $h_1(x), h_2(x), h_3(x)$  im 1. Quadranten

### B.7.1. Maximumbestimmung

Die Berechnung des Maximums von  $f(x, y)$  über einem beliebig vorgegebenen Rechteck kann reduziert werden auf die Berechnung des Maximums über einem Rechteck  $Z$ , welches nur im ersten Quadranten liegt. Auf Seite 233 haben wir gezeigt, wie vorzugehen ist, wenn das ursprüngliche Rechteck z.B. nur in der oberen Halbebene liegt. Wir betrachten jetzt also nur Rechtecke  $Z$ , die im 1. Quadranten liegen und den Punkt  $(0, +1)$  nicht enthalten. Bevor jedoch der maximale Funktionswert von  $f(x, y)$  auf dem Rand von  $Z$  berechnet werden kann, muss der Punkt  $M$  auf dem Rand von  $Z$  bestimmt werden, auf dem der maximale Funktionswert angenommen wird, und genau dies ist das eigentliche Problem, das für jede neue Funktion  $f(x, y)$  immer wieder neu gelöst werden muss. In unserem Fall benutzen wir den folgenden

**Satz B.1** (Die linke Parallele von  $Z$  schneidet  $h_1(x)$ ).

Schneidet die linke Parallele von  $Z$  den Graphen von  $h_1(x)$ , so ist dieser Schnittpunkt gleich  $M$ .

In den obigen Abbildungen wird dies bestätigt durch die Rechtecke  $Z_7, Z_9, Z_{11}, Z_{12}$ . Der Algorithmus zur Bestimmung von  $M$  kann daher konstruiert werden durch die folgenden drei Fälle, die sich selbst gegenseitig ausschließen:

1. Die linke obere Ecke von  $Z$  liegt unterhalb von  $h_1(x)$ , vgl.  $Z_5, Z_6, Z_8, Z_{14}, Z_{15}, Z_{16}$ .
2. Die linke untere Ecke von  $Z$  liegt oberhalb von  $h_1(x)$ , vgl.  $Z_1, Z_2, Z_3, Z_4, Z_{10}, Z_{13}, Z_{17}, Z_{18}$ .
3. Die linke Parallele von  $Z$  schneidet den Graphen von  $h_1(x)$ , vgl.  $Z_7, Z_9, Z_{11}, Z_{12}$ .

Auf der folgenden Seite findet man nach diesen Vorüberlegungen einen Grobalgorithmus zur Bestimmung des Maximumpunktes  $M$ .

```

if (y2 <= h1(x1)) // Linke obere Ecke von Z liegt unterhalb von h1(x):
{
  if (y2 <= h3(x2))
    M(x2,y2);
  else
  {
    if (y2 <= h3(x1)) // M: Schnittpunkt der oberen Parallelen mit h3(x)
      M(xS,y2) mit y2 = h3(xS); vgl. Z14
    else M(x1,y2);
  }
}
else
{
  if (y1 > h1(x1)) // Linke untere Ecke von Z liegt oberhalb von h1(x)
  {
    if (y1 > h2(x2)) M(x2,y1);
    else
    {
      if (y1 > h2(x1)) // M ist S.P. der unteren Parallelen mit h2(x)
        M(xS,y1) mit y1 = h2(xS);
      else M(x1,y1);
    }
  }
  else // Die linke Parallele von Z schneidet h1(x)
    M ist dieser Schnittpunkt mit h1(x);
}

```

Den Algorithmus zur Bestimmung des maximalen Funktionswertes findet man in der Funktion

```
MPFR::MpfrClass fM_Im_r1pz2(const MpfrClass& X, const MpfrClass& Y)
```

die in `mpfrclass.cpp` definiert ist. Wir betrachten jetzt noch die Fälle, in denen  $M$  der Schnittpunkt vom Rand von  $Z$  mit einer der Extremalkurven  $h_1(x), h_2(x), h_3(x)$  ist.

**Z<sub>14</sub>**: Die obere Parallele von  $Z_{14}$  schneidet den Graphen von  $h_3(x)$  im Punkt  $M(x_S, y_2)$ , mit  $y_2 = h_3(x_S) = \sqrt{1 + x_S^2 - 2x_S \cdot \sqrt{1 + x_S^2}}$ . Um  $f(x_S, y_2)$  berechnen zu können, muss zunächst  $x_S = x_S(y_2)$  als Funktion von  $y_2 < 1$  berechnet werden:  $x_S = \sqrt{(-1 - y_2^2 + 2\sqrt{1 - y_2^2 + y_2^4})/3}$ . Für  $f(x_S, y_2)$  erhält man schließlich, wenn  $y_2$  durch  $y$  ersetzt wird:

$$f(x_S(y), y) = h_4(y) := \frac{0.75 \cdot y \cdot \sqrt{-3 - 3y^2 + 6 \cdot \sqrt{1 - y^2 + y^4}}}{1 + y^4 + \sqrt{1 - y^2 + y^4} + y^2 \cdot (-4 + \sqrt{1 - y^2 + y^4})}, \quad 0 \leq y \leq 1/2.$$

Man erkennt, dass der obige Term für  $y \rightarrow 1$  die Form  $0/0$  besitzt und deshalb umgeformt werden muss:

$$f(x_S(y), y) = h_4(y) := \frac{(1 + y^2) \cdot \sqrt{1 - y^2 + y^4} - 1 + 4y^2 - y^4}{4y \cdot (1 - y^2) \cdot \sqrt{1 + y^2 + 2 \cdot \sqrt{1 - y^2 + y^4}}}, \quad 1/2 < y < 1.$$

Die Auswertung von  $h_4(x)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfrClass Im_r1pz2_fh4(const MPFR::MpfrClass& x)
```

die in `mpfrclass.cpp` definiert ist. Beachten Sie bitte, dass bei der Auswertung von  $h_4(y)$  in den beiden Teilbereichen  $0 \leq y \leq 1/2$  und  $1/2 < y < 1$  keine Auslöschung auftreten kann.

**Z17**: Die untere Parallele von  $Z_{17}$  schneidet den Graphen von  $h_2(x)$  im Punkt  $M(x_S, y_1)$ , mit  $y_1 = h_2(x_S) = \sqrt{1 + x_S^2 + 2x_S \cdot \sqrt{1 + x_S^2}}$ . Um  $f(x_S, y_1)$  berechnen zu können, muss zunächst  $x_S = x_S(y_1)$  als Funktion von  $y_1 > 1$  berechnet werden:  $x_S = \sqrt{(-1 - y_1^2 + 2\sqrt{1 - y_1^2 + y_1^4})/3}$ . Für  $f(x_S, y_1)$  erhält man, wenn  $y_1$  durch  $y$  ersetzt wird:

$$f(x_S(y), y) = h_6(y) := \frac{(1 + y^2) \cdot \sqrt{1 - y^2 + y^4} - 1 + 4y^2 - y^4}{4y \cdot (y^2 - 1) \cdot \sqrt{1 + y^2 + 2 \cdot \sqrt{1 - y^2 + y^4}}}, \quad 1 < y \leq 8.$$

Für  $y \rightarrow \infty$  liefert der obige Zähler starke Auslöschung, so dass weiter umgeformt wird:

$$h_6(y) := \frac{9 \cdot y \cdot (y^2 - 1)}{4 \cdot \sqrt{1 + y^2 + 2 \cdot \sqrt{1 - y^2 + y^4}} \cdot \left(1 - 4y^2 + y^4 + (1 + y^2) \cdot \sqrt{1 - y^2 + y^4}\right)}, \quad y > 8.$$

Für  $y \rightarrow \infty$  gilt  $h_6(y) \sim (9/(8\sqrt{3}))/y^2$ , so dass ein vorzeitiger Überlauf zu vermeiden ist. Mit  $u := (1/y)^2$  erhält man für  $y > 8$  die Darstellung

$$h_6(y) = g(u) := u \cdot \frac{9 \cdot (1 - u)}{4 \cdot \sqrt{1 + u + 2 \cdot \sqrt{1 - u + u^2}} \cdot \left(1 - 4u + u^2 + (1 + u) \cdot \sqrt{1 - u + u^2}\right)}.$$

Die Auswertung von  $h_6(x)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfiClass Im_r1pz2_fh6(const MPFR::MpfrClass& x)
```

die in `mpficlass.cpp` definiert ist. Beachten Sie bitte, dass bei der Auswertung von  $h_6(y)$  in den beiden Teilbereichen  $1 < y \leq 8$  und  $8 < y < \infty$  keine Auslöschung auftreten kann.

**Z12**: Die linke Parallele von  $Z_{12}$  schneidet den Graphen von  $h_1(x)$  im Punkt  $M(x_1, h_1(x_1))$ , mit  $h_1(x_1) = \sqrt{(1 - x_1^2 + 2 \cdot \sqrt{1 + x_1^2 + x_1^4})/3}$ . Für  $f(x, h_1(x))$  erhält man den Ausdruck

$$f(x, h_1(x)) = h_5(x) := \frac{3x \cdot \sqrt{3 - 3x^2 + 6 \cdot \sqrt{1 + x^2 + x^4}}}{4 \cdot \left(1 + 4x^2 + x^4 + (x^2 - 1) \cdot \sqrt{1 + x^2 + x^4}\right)}, \quad 2 \leq x < \infty.$$

Die Einschränkung  $2 \leq x < \infty$  ist notwendig, da der Nenner für  $x \rightarrow 0$  starke Auslöschung verursacht. Für  $x \rightarrow \infty$  gilt  $h_5(x) \sim 3\sqrt{3}/(8x^2)$ , so dass ein vorzeitiger Überlauf vermieden werden muss. Mit  $u := (1/x)^2$  erhält man nach einfachen Umformungen

$$h_5(x) = g(u) := u \cdot \frac{0.75 \cdot \sqrt{3u - 3 + 6 \cdot \sqrt{1 + u + u^2}}}{1 + 4u + u^2 + (1 - u) \cdot \sqrt{1 + u + u^2}}, \quad u = (1/x)^2, \quad 2 \leq x < \infty.$$

Im Fall  $x \rightarrow 0$  muss der obige Ausdruck für  $f(x, h_1(x))$  weiter umgeformt werden zu:

$$h_5(x) := \frac{3 \cdot (1 + x^4 + \sqrt{1 + x^2 + x^4}) \cdot \sqrt{3 - 3x^2 + 6 \cdot \sqrt{1 + x^2 + x^4}}}{4x \cdot \left(4 + 2x^2 + 5x^4 + x^6 + (5 + x^4) \cdot \sqrt{1 + x^2 + x^4}\right)}, \quad 0 < x \leq 2.$$

Beim Auswerten der Funktionen  $h_5(x)$  und  $g(u)$  tritt jetzt in beiden Teilbereichen keinerlei Auslöschung auf. Die nahezu optimale Einschließung der Funktionswerte von  $h_5(x)$  und  $g(u)$  erfolgt auf der Maschine mit Hilfe der Funktion

```
MpfiClass Im_r1pz2_fH1(const MPFR::MpfrClass& x)
```

die in `mpficlass.cpp` definiert ist.

Damit haben wir alle Fälle behandelt, in denen der Maximumpunkt  $M$  ein Schnittpunkt des Randes von  $Z$  mit einer der Extremalkurven  $h_1(x), h_2(x), h_3(x)$  ist.

## B.7.2. Minimumbestimmung

In Abbildung B.24 geben die Pfeile die Richtung wachsender Funktionswerte von  $f(x, y)$  an, so dass ein Minimumpunkt  $m$  nicht auf einer der Extremalkurven  $h_1(x), h_2(x), h_3(x)$  liegen kann. Daher kann  $m$  nur einer der vier Eckpunkte eines gegebenen Rechtecks  $Z$  im 1. Quadranten sein, so dass bei einem echten Rechteck vier Funktionswertberechnungen und drei Vergleiche notwendig sind. Bei einem entarteten Rechteck oder bei einem zweidimensionalen Punktintervall sind entsprechend nur zwei bzw. ein Funktionswert zu berechnen. Die Bestimmung des minimalen Funktionswertes von  $f(x, y)$  über einem der Eckpunkte von  $Z$  erfolgt mit Hilfe der Funktion

`MpfrClass fm_Im_r1pz2(const MpfiClass& X, const MpfiClass& Y),`

die in `mpficlass.cpp` definiert ist.

In nachfolgender Tabelle sind alle notwendigen numerischen Funktionen zusammengestellt, um den maximalen oder minimalen Funktionswert  $f(x, y)$  über dem Rand eines vorgegebenen Rechtecks  $Z$  aus dem 1. Quadranten berechnen zu können.

Nr.	Funktionsterm, $u := (1/x)^2$	Funktionsaufruf
1.	$f(x, y) = \frac{2xy}{4x^2y^2 + (1 + x^2 - y^2)^2}$ , mit Skalierung	<code>Im_r1pz2(x, y, rnd)</code>
2.	$h_1(x) = \sqrt{(1 - x^2 + 2\sqrt{1 + x^2 + x^4})/3}$ , $0 \leq x \leq 16$	<code>Re_r1pz2_H1(x)</code>
3.	$h_1(x) = x \cdot \sqrt{(u - 1 + 2\sqrt{1 + u + u^2})/3}$ , $x > 16$	<code>Re_r1pz2_H1(x)</code>
4.	$h_2(x) = \sqrt{1 + x^2 + 2x \cdot \sqrt{1 + x^2}}$ , $0 \leq x \leq 16$	<code>Re_r1pz2_H2(x)</code>
5.	$h_2(x) = x \cdot \sqrt{1 + u + 2\sqrt{1 + u}}$ , $x > 16$	<code>Re_r1pz2_H2(x)</code>
6.	$h_3(x) = \sqrt{1 + x^2 - 2x \cdot \sqrt{1 + x^2}}$ , $0 \leq x \leq 1/\sqrt{3}$	<code>Re_r1pz2_H3(x)</code>
7.	$h_4(x) = \frac{0.75 \cdot x \cdot \sqrt{-3 - 3x^2 + 6 \cdot \sqrt{1 - x^2 + x^4}}}{1 + x^4 + \sqrt{1 - x^2 + x^4} + x^2(\sqrt{1 - x^2 + x^4} - 4)}$ , $0 \leq x \leq 1/2$	<code>Im_r1pz2_fH4(x)</code>
8.	$h_4(x) = \frac{(1 + x^2) \cdot \sqrt{1 - x^2 + x^4 - 1 + 4x^2 - x^4}}{4x \cdot (1 - x^2) \cdot \sqrt{1 + x^2 + 2 \cdot \sqrt{1 - x^2 + x^4}}}$ , $1/2 < x < 1$	<code>Im_r1pz2_fH4(x)</code>
9.	$h_5(x) = u \cdot \frac{0.75 \cdot \sqrt{3u - 3 + 6 \cdot \sqrt{1 + u + u^2}}}{1 + 4u + u^2 + (1 - u) \cdot \sqrt{1 + u + u^2}}$ , $x > 2$	<code>Im_r1pz2_fH1(x)</code>
10.	$h_5(x) = \frac{3(1 + x^4 + \sqrt{1 + x^2 + x^4}) \cdot \sqrt{3 - 3x^2 + 6\sqrt{1 + x^2 + x^4}}}{4x \cdot (4 + 2x^2 + 5x^4 + x^6 + (5 + x^4) \cdot \sqrt{1 + x^2 + x^4})}$ , $x \leq 2$	<code>Im_r1pz2_fH1(x)</code>
11.	$h_6(x) = \frac{(1 + x^2) \cdot \sqrt{1 - x^2 + x^4 - 1 + 4x^2 - x^4}}{4x \cdot (x^2 - 1) \cdot \sqrt{1 + x^2 + 2 \cdot \sqrt{1 - x^2 + x^4}}}$ , $1 < x \leq 8$	<code>Im_r1pz2_fH6(x)</code>
12.	$h_6(x) = u \cdot \frac{9 \cdot (1 - u) / (1 - 4u + u^2 + (1 + u) \cdot \sqrt{1 - u + u^2})}{4 \cdot \sqrt{1 + u + 2 \cdot \sqrt{1 - u + u^2}}}$ , $x > 8$	<code>Im_r1pz2_fH6(x)</code>

Obige Tabelle zeigt die große Vielfalt der notwendigen numerischen Funktionen, um den minimalen und maximalen Funktionswert  $f(x, y)$  über dem Rand eines beliebigen Rechtecks  $Z$  aus dem 1. Quadranten berechnen zu können. Realisiert werden muss zusätzlich ein Algorithmus, um diese Berechnungen für ein beliebig vorgegebenes Rechteck  $Z$ , z.B. aus der oberen Halbebene, nur für ein entsprechendes Rechteck aus dem 1. Quadranten durchführen zu können. Außerdem benötigt man einen Algorithmus, der für ein beliebiges solches Rechteck aus dem 1. Quadranten die Punkte  $m, M$  ermittelt, über denen der minimale bzw. maximale Funktionswert  $f(x, y)$  angenommen wird.





**B.8.**  $x/\sqrt{1+x^2}$ 

Für das vorgegebene reelle Maschinenintervall  $X \ni x$  liefert die Intervallfunktion

`MpfiClass xdsqrt1px2 (MpfiClass& X)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x) = \frac{x}{\sqrt{1+x^2}}, \quad \text{mit: } x \in X.$$

Da  $f(x)$  im ganzen Definitionsbereich  $D_f = \mathbb{R}$  monoton wächst, ist die Implementierung recht einfach, da  $f(x)$  nur auf den Intervallrändern `Inf(X)` und `Sup(X)` mit entsprechender Rundung auszuwerten ist.

**B.9.**  $x/\sqrt{1-x^2}$ 

Für das vorgegebene reelle Maschinenintervall  $X \ni x$  liefert die Intervallfunktion

`MpfiClass xdsqrt1mx2 (MpfiClass& X)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x) = \frac{x}{\sqrt{1-x^2}}, \quad \text{mit: } x \in X.$$

Da  $f(x)$  im ganzen Definitionsbereich  $D_f = (-1, +1)$  monoton wächst, ist die Implementierung recht einfach, da  $f(x)$  nur auf den beiden Intervallrändern `Inf(X)` und `Sup(X)` mit entsprechender Rundung auszuwerten ist.

**B.10.**  $x/\sqrt{x^2-1}$ 

Für das vorgegebene reelle Maschinenintervall  $X \ni x$  liefert die Intervallfunktion

`MpfiClass xdsqrtx2m1 (MpfiClass& X)`

eine fast optimale Einschließung aller Funktionswerte

$$f(x) = \frac{x}{\sqrt{x^2-1}}, \quad \text{mit: } x \in X \subseteq D_f := \{x \in \mathbb{R} \mid x < -1 \text{ or } x > +1\}.$$

Da  $f(x)$  im ganzen Definitionsbereich  $D_f$  monoton fällt, ist die Implementierung recht einfach, da  $f(x)$  nur auf den beiden Intervallrändern `Inf(X)` und `Sup(X)` mit entsprechender Rundung auszuwerten ist.



## C. Elementarfunktionen für komplexe Punkt- und Intervallargumente

Die Aufgabenstellung wird wie folgt beschrieben:

Zu einem vorgegebenen Intervallargument

$$(C.1) \quad Z = X + i \cdot Y; \quad X = [x_1, x_2], \quad Y = [y_1, y_2];$$

mit reellen Intervallen  $X, Y$  ist zu einer gegebenen komplexwertigen Funktion  $w = f(z)$ , mit  $z, w \in \mathbb{C}$ , eine möglichst optimale Einschließung aller Funktionswerte  $w$  gesucht, wenn  $z \in Z$ , d.h. einzuschließen ist

$$(C.2) \quad W_{f,Z} := \{w \mid w = f(z) \wedge z \in Z\}, \quad z = x + i \cdot y.$$

Zu beachten ist, dass das Bild  $W = f(Z)$  des Argumentintervalls  $Z$  i.a. kein achsenparalleles Rechteck ist, so dass die Einschließung von  $W$  durch ein Intervall  $F(Z) = U(Z) + i \cdot V(Z)$  auch dann i.a. erhebliche Überschätzungen liefert, wenn das Rechteck  $F$  die Menge  $W$  **optimal** einschließt. Für  $Z = [-1, 3] + i \cdot [0.2, 1]$  werden dies Überschätzungen in Abb. C.1 am Beispiel der komplexen sin-Funktion, d.h. für  $f(z) = \sin(z)$ ,  $z \in Z$ , überzeugend dargestellt. So gehört z.B. der ganze von  $F(Z)$  eingeschlossene Teil des 3. Quadranten **nicht** zur Wertemenge  $W_{f,Z}$ , [38, S.16].

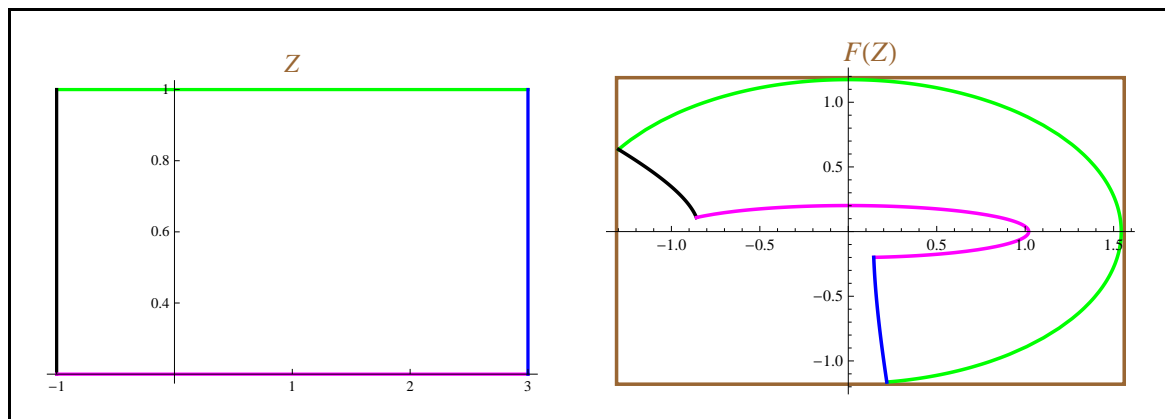


Abbildung C.1.: Intervall-Funktion

Jetzt muss noch geklärt werden, wie die reellen Intervalle  $U(Z)$  und  $V(Z)$  zu berechnen sind.

Dazu betrachten wir zunächst den einfachsten Fall, dass mit  $f(z) = u + i \cdot v$  die Funktionen  $u(x, y)$  und  $v(x, y)$  separabel sind, d.h. für z.B.  $u$  soll im Falle einer Multiplikation gelten:

$$u(x, y) = u_1(x) \cdot u_2(y), \quad x \in X, \quad y \in Y.$$

Eine optimale Einschließung von  $u(x, y)$  erhält man aber nur dann, wenn in den auszuwertenden Ausdrücken  $u_1$  und  $u_2$  die Variablen  $x$  bzw.  $y$  jeweils **nur einmal** vorkommen und wenn das Produkt  $u_1(x) \cdot u_2(y)$  intervallmäßig ausgewertet wird, [5]. Es gilt dann auch für  $v(x, y)$

$$\begin{aligned} u(x, y) &= u_1(x) \cdot u_2(y) \in u_1(X) \diamond u_2(Y) =: U(Z), \\ v(x, y) &= v_1(x) \cdot v_2(y) \in v_1(X) \diamond v_2(Y) =: V(Z). \end{aligned}$$

Ein Beispiel für eine separable Funktion ist die komplexe Exponentialfunktion, bei ihr gilt:

$$(C.3) \quad u(x, y) = e^x \cdot \cos(y), \quad v(x, y) = e^x \cdot \sin(y), \quad x \in X, \quad y \in Y.$$

Weitere Beispiele für separable Funktionen sind:  $\sin(z)$ ,  $\cos(z)$ ,  $\sinh(z)$  und  $\cosh(z)$ , wobei deren Implementierung besonders einfach ist, da die dabei benötigten reellen Intervall-Funktionen  $\sin(X)$ ,  $\cos(X)$ ,  $\cosh(Y)$  und  $\sinh(Y)$  bereits in `mpficlass.cpp` definiert sind, vgl. auch die Tabelle auf Seite 68.

Wir kommen jetzt zu den **nicht-separablen** Funktionen  $f(z) = u(x, y) + i \cdot v(x, y)$ , bei denen die Berechnung von  $U(Z) \ni u(x, y)$  und  $V(Z) \ni v(x, y)$  ausnahmslos sehr viel schwieriger ist. Eine Basis zur Berechnung der reellen Intervalle  $U, V$  liefert jetzt der folgende Satz, [7].

Schreibt man mit  $z = x + i \cdot y \in \mathbb{C}$  die Funktion  $w = f(z)$  in der Form

$$w = f(z) = u(x, y) + i \cdot v(x, y), \quad z = x + i \cdot y \in Z \subset \mathbb{C}, \quad \text{so gilt:}$$

Ist  $f = u + i \cdot v : Z \subset G \rightarrow \mathbb{C}$  holomorph im Gebiet  $G$ , so nehmen sowohl  $u(x, y)$  als auch  $v(x, y)$  als harmonische Funktionen ihr Maximum und Minimum auf dem Rand von  $Z$  an.

Die Extrema von  $u(x, y)$  und  $v(x, y)$  müssen also nur auf dem **Rand** von  $Z$  gesucht werden.

Wir betrachten auch jetzt wieder den einfachsten Fall, wenn die Koordinaten der Randpunkte  $m$  und  $M$ , in denen das Minimum bzw. Maximum angenommen wird, Maschinenzahlen des durch `prec` bestimmten Zahlenformats sind. In der folgenden Abb. A.2 sind für verschiedene Argumentintervalle  $Z$  die Punkte  $m$  und  $M$  angegeben, in denen für die Realteilfunktion  $u(x, y)$  der  $\arcsin(z)$ -Funktion das Minimum bzw. Maximum angenommen wird.

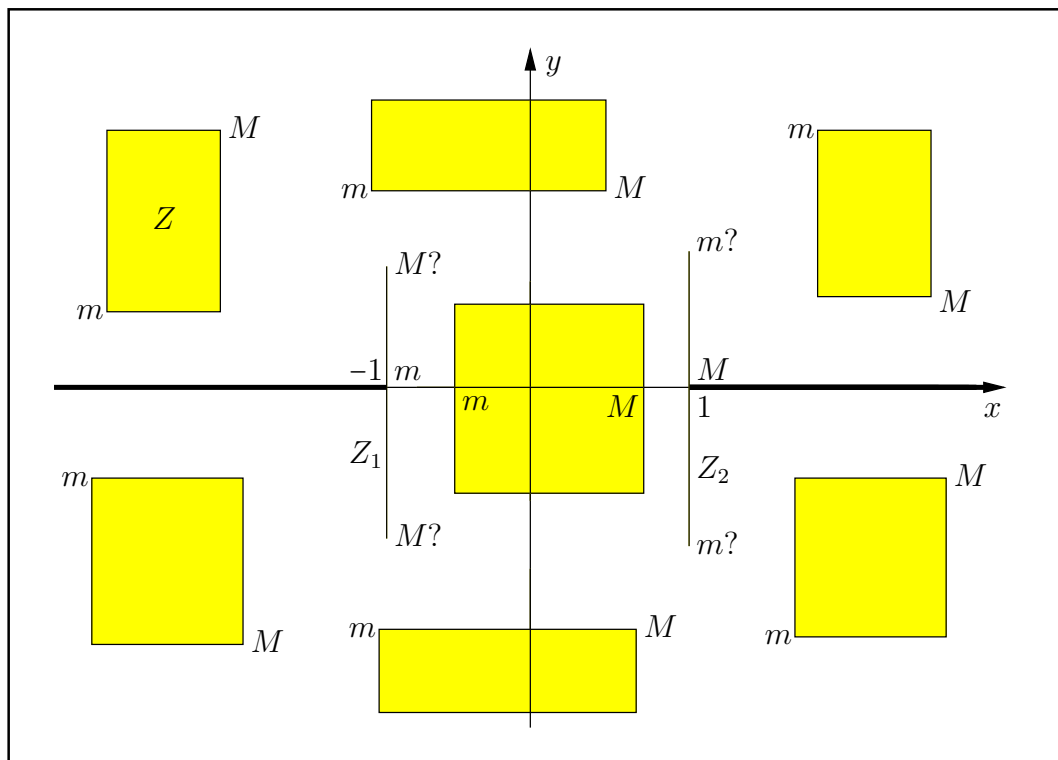


Abbildung C.2.: Die Lage der Punkte  $m, M$  auf  $Z$  beim Realteil von  $\arcsin(z)$ .

Bei den mit '?' gekennzeichneten Punkten ist deren  $y$ -Koordinate als Maximum von  $-y_1$  und  $y_2$  definiert, mit  $Y := [y_1, y_2]$ . Für die durch  $-1$  bzw.  $+1$  laufenden Argumentintervalle  $Z_1, Z_2$  gilt:

$$m = (1, \max(-y_1, y_2)); \quad M = (-1, \max(-y_1, y_2));$$

Aus Abb. C.2 erkennt man, dass die Koordinaten von  $m$  und  $M$  alles Maschinenzahlen sind. Bedeutet  $u(x, y)$  die Realteilstfunktion der  $\arcsin(z)$ -Funktion und ist  $u_d(m)$  der abgerundete und  $u_u(M)$  der aufgerundete Funktionswert, jeweils an den Punkten  $m$  und  $M$ , so ist  $U(Z)$  gegeben durch

$$U(Z) := [u_d(m), u_u(M)] \ni u(x, y), \quad z = x + i \cdot y \in Z.$$

Es stellt sich jetzt noch die Frage, wie die Werte  $u_d(m)$  und  $u_u(M)$  zu berechnen sind. Dazu gibt es im Wesentlichen drei Möglichkeiten:

**Methode 1.** Nach [17],[18],[38],[39] wird eine garantierte Fehlerschranke für  $u(x, y)$  berechnet, mit der dann die Werte  $u_d(m)$  und  $u_u(M)$  bestimmt werden können. Dies ist jedoch programmiertechnisch sehr aufwendig, so dass meist das folgende Verfahren zur Anwendung kommt.

**Methode 2.** Wenn zur intervallmäßigen Auswertung von  $u(x, y)$  die entsprechenden Intervallfunktionen bereits implementiert sind, so kann  $u(x, y)$  an den Stellen  $m$  und  $M$  durch die Intervalle  $\mathbf{u}(m)$  und  $\mathbf{u}(M)$  eingeschlossen werden, und  $U(Z)$  ist gegeben durch, [54],[55]

$$U(Z) = [\mathbf{Inf}(\mathbf{u}(m)), \mathbf{Sup}(\mathbf{u}(M))].$$

Dieses Verfahren kann programmiertechnisch vergleichsweise einfach realisiert werden, es erzeugt aber wegen der notwendigen Intervall-Auswertungen deutlich größere Laufzeiten.

**Methode 3.** Wenn  $u(x, y)$  für punktförmige Maschinenzahlen  $x, y$  so implementiert ist, dass man mit  $\mathbf{u}(x, y, \text{RoundDown})$  bzw.  $\mathbf{u}(x, y, \text{RoundUp})$  garantiert ab- bzw. aufgerundete Funktionswerte im Punkt  $P(x, y)$  erhält, so kann  $U(Z)$  realisiert werden durch

$$U(Z) = [\mathbf{u}(m, \text{RoundDown}), \mathbf{u}(M, \text{RoundUp})].$$

Der Nachteil ist jetzt, dass  $\mathbf{u}(x, y, \text{rnd})$  mit den Rundungsparametern  $\text{rnd} = \text{RoundDown}$  und  $\text{rnd} = \text{RoundUp}$  mit Hilfe der MPFR- und MPFI-Bibliotheken sehr sorgfältig zu implementieren ist, [23],[24],[52],[60]. Der Vorteil ist aber die nahezu optimale Laufzeit, da Intervallauswertungen jetzt nicht zur Anwendung kommen.

Bei der Realisierung von  $U(Z)$  und  $V(Z)$  sollte daher nur die Methode 3. zur Anwendung kommen. Nur in wirklich komplizierten Fällen wird man daher auf die Methode 2. zurückgreifen.

Abschließend noch einige Bemerkungen zur Lage der achsenparallelen Argumentintervalle  $Z$  in der komplexen Ebene.

- Ist  $f(z)$  holomorph in der ganzen komplexen Ebene, so kann  $Z = X + i \cdot Y$  beliebig gewählt werden, wobei aber die Randpunkte der reellen Intervalle  $X$  und  $Y$  Maschinenzahlen sind, so dass damit das komplexe Intervall  $Z$  stets endlich sein wird.
- Der Hauptwert der komplexen Funktion  $f(z) = \arcsin(z)$  besitzt nach Abb. C.2 zwei Verzweigungsschnitte auf der reellen Achse von 1 bis  $+\infty$  und von  $-\infty$  bis  $-1$ . In diesem Fall darf ein Verzweigungsschnitt das Intervall  $Z$  nicht im Innern durchlaufen, vielmehr darf  $Z$  einen Verzweigungsschnitt höchstens von oben oder von unten berühren. In der CoStLy-Bibliothek von M. Neher, [54],[55], darf  $Z$  keinen Punkt der beiden Verzweigungsschnitte enthalten; lediglich Argumentintervalle  $Z$  der Breite Null durch die Verzweigungspunkte  $P_1(-1, 0)$  und  $P_2(+1, 0)$  sind zugelassen.
- Bei mehrdeutigen Funktionen werden stets Einschließungen des jeweiligen Hauptwertes berechnet, wobei die erlaubten bzw. verbotenen Lagen der Argumentintervalle  $Z$  stets anzugeben sind.

Wir kommen jetzt zum kompliziertesten Fall, wenn die Extrempunkte  $m, M$  nicht mehr Eckpunkte oder Schnittpunkte von  $Z$  mit den Achsen sind. In diesem Fall ist dann i.a. nur noch eine Koordinate von  $m$  oder  $M$  eine Maschinenzahl, so dass die Auswertung der Real- bzw. Imaginärteildfunktion  $u(x, y)$  bzw.  $v(x, y)$  auf der Maschine nicht mehr so ohne Weiteres möglich ist. Als Beispiel betrachten wir in Abb. C.3 für den Hauptwert der  $\arctan(z)$ -Funktion in der oberen Halbebene einige Intervalle  $Z$ , bei denen nur die  $y$ -Koordinate des jeweiligen Punktes  $m$  eine Maschinenzahl ist, [38].

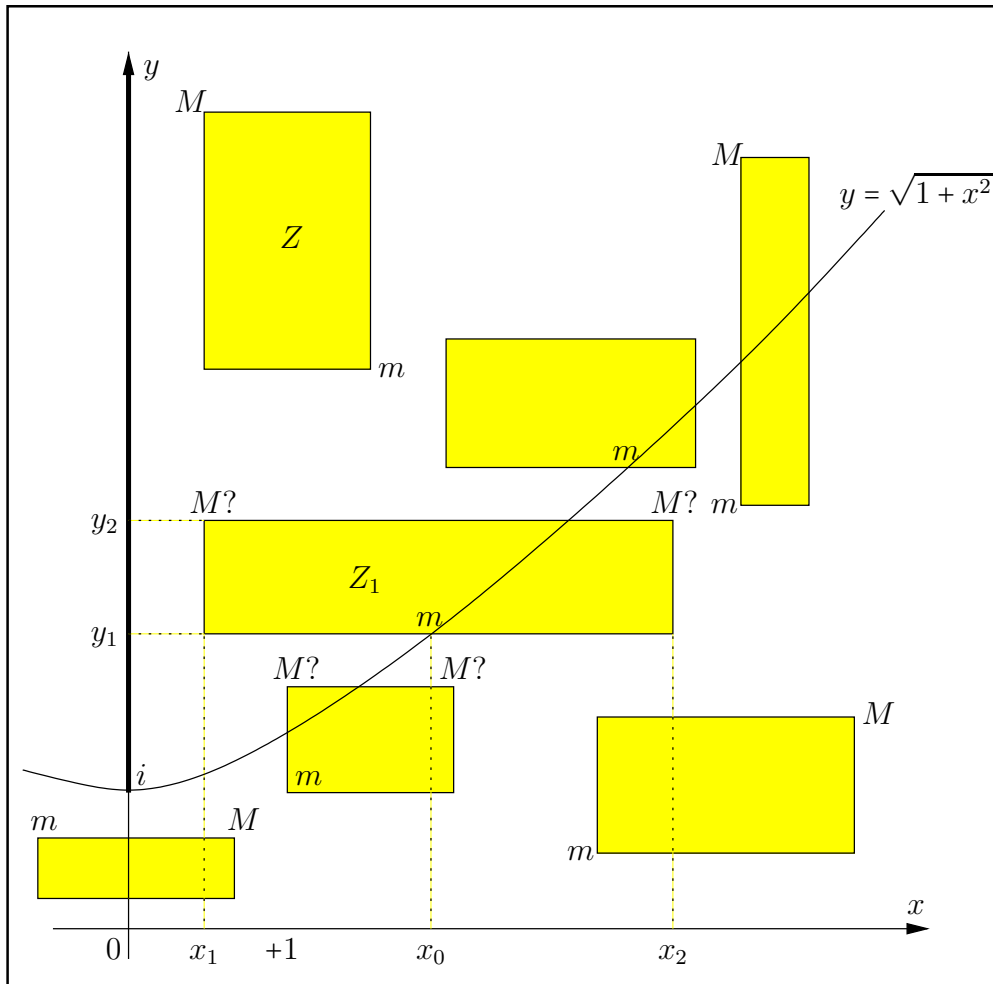


Abbildung C.3.: Die Lage der Punkte  $m, M$  auf  $Z$  beim Realteil von  $\arctan(z)$ .

Wir betrachten das Argumentintervall  $Z_1$  im 1. Quadranten außerhalb des Einheitskreises um 0. Für die Realteildfunktion  $u(x, y)$  von  $\arctan(z)$  gilt in diesem Bereich

$$(C.4) \quad u(x, y) := \frac{1}{2} \arctan \frac{2x}{1-x^2-y^2} + \frac{\pi}{2},$$

und wenn für die  $x$ -Koordinate  $x_0$  von  $m$  gelten soll  $x_1 < x_0 < x_2$ , so muss für ein lokales Minimum in  $m$  die notwendige Bedingung

$$\frac{\partial u(x, y)}{\partial x} = 0 \quad \Longleftrightarrow \quad y = +\sqrt{1+x^2}$$

erfüllt sein, d.h. es muss gelten  $y_1 = +\sqrt{1+x_0^2}$ , wobei  $x_0 = \sqrt{y_1^2 - 1}$  i.a. keine Maschinenzahl ist. Glücklicherweise wird aber  $x_0$  nicht explizit benötigt, da zur Einschließung von  $u(x, y)$  in  $Z$  ein abgerundeter Wert von  $u(x_0, y_1)$  benötigt wird. Dazu setzt man in (C.4)  $x = \sqrt{y_1^2 - 1}$  und  $y = y_1$

und erhält

$$(C.5) \quad u(x_0, y_1) := -\frac{1}{2} \arctan \frac{1}{\sqrt{y_1^2 - 1}} + \frac{\pi}{2} =: u(y_1),$$

wobei jetzt  $u(y_1)$  mit der Maschinenzahl  $y_1$  problemlos ausgewertet werden kann. Ein abgerundeter Wert von  $u(y_1)$  kann jetzt wieder analog zu Seite 245 mit drei verschiedenen Methoden berechnet werden.

**Methode 1.** Nach [17],[18],[38],[39] wird eine garantierte Fehlerschranke für  $u(y)$  berechnet, mit der dann der abgerundeter Wert von  $u(y_1)$  bestimmt werden kann. Dies ist jedoch programmiertechnisch wieder sehr aufwendig, so dass meist das folgende Verfahren zur Anwendung kommt.

**Methode 2.** Da zur intervallmäßigen Auswertung von  $u(y)$  die entsprechenden Intervallfunktionen bereits implementiert sind, kann  $u(y_1)$  durch das Intervall  $\mathbf{u}(y_1)$  eingeschlossen werden, und der gesuchte abgerundete Wert ist gegeben durch, [54],[55]

$$\mathbf{Inf}(\mathbf{u}(y_1)) \leq u(y_1) \leq u(x, y), \quad z = x + i \cdot y \in Z..$$

Dieses Verfahren kann programmiertechnisch vergleichsweise einfach realisiert werden, es erzeugt aber wegen der notwendigen Intervall-Auswertungen deutlich größere Laufzeiten.

**Methode 3.** Wenn  $u(y)$  für die punktförmigen Maschinenzahlen  $y$  so implementiert ist, dass man mit  $u(y, \text{RoundDown})$  einen garantiert abgerundeten Funktionswert erhält, so ist dieser die gesuchte Unterschranke

$$u(y_1, \text{RoundDown}) \leq u(y_1) \leq u(x, y), \quad z = x + i \cdot y \in Z.$$

Der Nachteil ist jetzt, dass  $u(y, \text{RoundDown})$  mit Hilfe der MPFR- und MPFI-Bibliotheken sehr sorgfältig zu implementieren ist, [23],[24],[52],[60]. Der Vorteil ist aber die nahezu optimale Laufzeit, da Intervallauswertungen jetzt nicht zur Anwendung kommen.

### Anmerkungen:

1. Da die Bedingung  $\partial u(x_0, y_1)/\partial x = 0$  im Punkt  $m$  für ein lokales Minimum nur notwendig ist, muss die Existenz des lokales Minimums zusätzlich abgesichert werden. In [38, Seite 145] erfolgt dies durch eine allgemeine Monotoniebetrachtung der Funktion  $u(x, y)$ .
2. Nach Abb. C.3 schneidet die Hyperbel  $y = +\sqrt{1+x^2}$  auch die obere Parallele von  $Z_1$ , so dass auch dort ein lokales Minimum auftreten könnte. Wegen  $0 < y_1 < y_2$  gilt jedoch nach (C.5)  $u(y_1) < u(y_2)$ , so dass das lokale Minimum der Realteilmfunktion  $u(x, y)$  wirklich in  $m$  angenommen wird.
3. Für die Realteilmfunktion  $u(x, y)$  liefert  $\partial u(x, y)/\partial x = 0$  den geometrischen Ort aller derjenigen Punkte  $P(x, y)$ , in denen auf einer **Parallelen** zur  $x$ -Achse, z.B. durch  $y = y_1$ , ein lokaler Extremwert auftreten kann. Dieser geometrische Ort wird auch als **Extremalkurve** bezeichnet und ist bei der  $\arctan(z)$ -Funktion für den Realteil die Hyperbel  $y^2 - x^2 = 1$ . Schneidet also einer der beiden Hyperbeläste  $y = \pm\sqrt{1+x^2}$  einen zur reellen Achse parallelen Rand von  $Z$ , so können relative Extrema auf diesen Rändern nur in diesen Schnittpunkten angenommen werden. Ganz entsprechend liefert ein Schnittpunkt der **Extremalkurve**  $\partial u(x, y)/\partial y = 0$  mit einer Randparallelen von  $Z$  zur imaginären Achse einen Punkt, in dem ein lokales Extremum auf dieser Parallelen zur  $y$ -Achse auftreten kann.
4. Für die Imaginärteilmfunktion  $v(x, y)$  sind die beiden **Extremalkurven**  $\partial v(x, y)/\partial x = 0$  und  $\partial v(x, y)/\partial y = 0$  und ihre Schnittpunkte mit dem Rand von  $Z$  ganz entsprechend zu betrachten.

## C.1. Elementarfunktionen für komplexe Punktargumente

Alle nachfolgenden Funktionen mit Punktargumenten  $z$  vom Typ `MpfcClass` können mit einem zusätzlichen Rundungsparameter `rnd` aufgerufen werden, wobei nur die Rundungen

$$\text{rnd} = \text{RoundNearest}, \quad \text{rnd} = \text{RoundDown}, \quad \text{rnd} = \text{RoundUp}$$

zur Verfügung stehen. Die obigen Rundungen werden aber nicht immer optimal ausgeführt, d.h. mit z.B. `rnd = RoundUp` wird nicht die nächst-größere Rasterzahl berechnet, sondern nur eine der benachbarten Rasterzahlen, die garantiert rechts vom exakten Funktionswert liegt. Wird `rnd` nicht gesetzt, so wird ebenfalls nur in die unmittelbare Nähe des exakten Funktionswertes gerundet. Die für die Implementierung der neuen Funktionen benötigten Bedingungen für das korrekte Runden bei den vier Grundoperationen werden ausführlich beschrieben im Abschnitt A.1 ab Seite 179. Die Elementarfunktionen für komplexe Punktargumente sind in der Tabelle auf Seite 92 zusammengestellt. Im ersten sehr einfachen Beispiel betrachten wir nur den Realteil der komplexen Exponentialfunktion, wobei nur die korrekte Rundung bei den Operanden und bei der Multiplikation zu beachten ist.

### C.1.1. Exponentialfunktion, Realteil

Nach (C.3) von Seite 244 gilt mit  $z = x + i \cdot y \in \mathbb{C}$  für den Realteil  $u(x, y)$  der Exponentialfunktion

$$(C.6) \quad \Re(e^z) := u(x, y) = e^x \cdot \cos(y).$$

Zunächst sollen nur **aufgerundete** Funktionswerte von  $u(x, y)$  berechnet werden. Nach (A.15) und (A.16) von Seite 181 muss dazu der Funktionswert  $\cos(y)$  stets aufgerundet werden. Der nachfolgende Code zeigt die entsprechenden Anweisungen.

```
1  MpfcClass re(0), tmp(0);
2  // Calculating the real part:
3  if (rnd==MPFR_RNDU)
4  {
5      mpfr_cos(tmp.GetValue(), z.mpfr_im, rnd); // tmp = cos(Im(z))
6      if (tmp < 0)
7          mpfr_exp(re.GetValue(), z.mpfr_re, MPFR_RNDU);
8      else mpfr_exp(re.GetValue(), z.mpfr_re, rnd);
9  };
10 mpfr_mul(re.GetValue(), re.GetValue(), tmp.GetValue(), rnd);
```

In Zeile 5 wird der aufgerundete Funktionswert  $\cos(y)$  berechnet und in `tmp` gespeichert. Nach (A.15) und (A.16) muss  $e^x$  im Fall `tmp < 0` abgerundet und andernfalls aufgerundet werden, was in den Zeilen 7 und 8 realisiert wird. Nach (A.15) wird noch verlangt

$$\text{tmp} > 0 \implies \cos(y) > 0 \quad \text{und} \quad \text{re} > 0 \implies e^x > 0.$$

Die beiden Bedingungen sind wegen der optimalen Rundung der Funktionen aus der MPFR-Bibliothek sicher erfüllt. Mit (A.15) wird zusätzlich noch verlangt, dass nur einer der beiden aufgerundeten Funktionswerte verschwindet. Wegen der positiven Exponentialfunktion ist auch diese Forderung sicher erfüllt. Nach (A.15) und (A.16) ist bei der Multiplikation der gerundeten Werte `tmp` und `re` jeweils aufzurunden, was in Zeile 10 realisiert wird, wobei der aufgerundete Realteil von  $e^z$  in `re` gespeichert wird. Die Berechnung eines abgerundeten Realteils erfolgt nach (A.11) und (A.12) ganz analog. Der entsprechende Quelltext, auch für die Berechnung des Imaginärteils, befindet sich in der Datei `mpfcclass.cpp`.



### C.1.2. $\sin(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist der komplexe Sinus definiert durch

$$\sin(z) = \sin(x) \cdot \cosh(y) + i \cdot \cos(x) \cdot \sinh(y).$$

Die Berechnung der gerundeten Realteilwerte erfolgt ganz analog zum Realteil der Exponentialfunktion, da einer der Faktoren stets positiv ist. Im Vergleich dazu ist die Auswertung des Imaginärteils etwas komplizierter, da neben  $\cos(x)$  auch der zweite Faktor  $\sinh(y)$  sowohl positiv als auch negativ werden kann. Beachten Sie, dass die vier obigen Funktionen aus der MPFR-Bibliothek direkt zur Verfügung stehen und auf der Maschine **optimal** gerundet werden. Daher gelten für diese Funktionen die Aussagen (A.36) bis (A.41) auf Seite 183. Weitere Einzelheiten findet man in der Funktion `MpfcClass sin(const MpfcClass& z, RoundingMode rnd)`; die in der Datei `mpfcclass.cpp` definiert ist.

### C.1.3. $\cos(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist der komplexe Cosinus definiert durch

$$\cos(z) = \cos(x) \cdot \cosh(y) - i \cdot \sin(x) \cdot \sinh(y).$$

Die Berechnung der gerundeten Real- und Imaginärteilwerte erfolgt analog zur  $\sin(z)$ -Funktion, wobei zusätzlich das Minus-Zeichen beim Imaginärteil zu beachten ist. Weitere Einzelheiten findet man in der Funktion `MpfcClass cos(const MpfcClass& z, RoundingMode rnd)`; die in der Datei `mpfcclass.cpp` definiert ist.

### C.1.4. $\tan(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist der komplexe Tangens definiert durch

$$\tan(z) = \frac{\sin(2x)}{2 \cdot \{\cos^2(x) + \sinh^2(y)\}} + i \cdot \frac{\sinh(2y)}{2 \cdot \{\cos^2(x) + \sinh^2(y)\}}.$$

Die Argumente  $2x$  und  $2y$  werden rundungsfehlerfrei berechnet, falls kein Überlauf eintritt. Der Nenner  $2 \cdot \{\cos^2(x) + \sinh^2(y)\}$  wird wegen der Quadrate ohne Auslöschung berechnet, wobei jedoch durch  $\sinh^2(y)$  ein Überlauf erzeugt werden kann. Dadurch treten bei den gerundeten Real- und Imaginärteilwerten starke Überschätzungen ein. Weitere Einzelheiten findet man in der Funktion `MpfcClass tan(const MpfcClass& z, RoundingMode rnd)`; die in der Datei `mpfcclass.cpp` definiert ist.

### C.1.5. $\cot(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist der komplexe Cotangens definiert durch

$$\cot(z) = \frac{\sin(2x)}{2 \cdot \{\sin^2(x) + \sinh^2(y)\}} + i \cdot \frac{\sinh(-2y)}{2 \cdot \{\sin^2(x) + \sinh^2(y)\}}.$$

Es gelten die gleichen Überlegungen wie bei der  $\tan(z)$ -Funktion. Real- und Imaginärteil sind implementiert in den Funktionen

```
MpfrClass Re_cot(const MpfcClass& z, RoundingMode rnd);  
MpfrClass Im_cot(const MpfcClass& z, RoundingMode rnd);
```

die in der Datei `mpfcclass.cpp` definiert sind.

### C.1.6. $\arg(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  wird das Argument von  $z$ , d.h.  $\arg(z)$ , definiert durch das Bogenmaß des Winkels, den der Fahrstrahl vom Ursprung nach  $z$  mit der positiven reellen Achse einschließt. Die komplexe Ebene ist dabei längs der negativen reellen Achse von 0 bis  $-\infty$  aufgeschnitten.

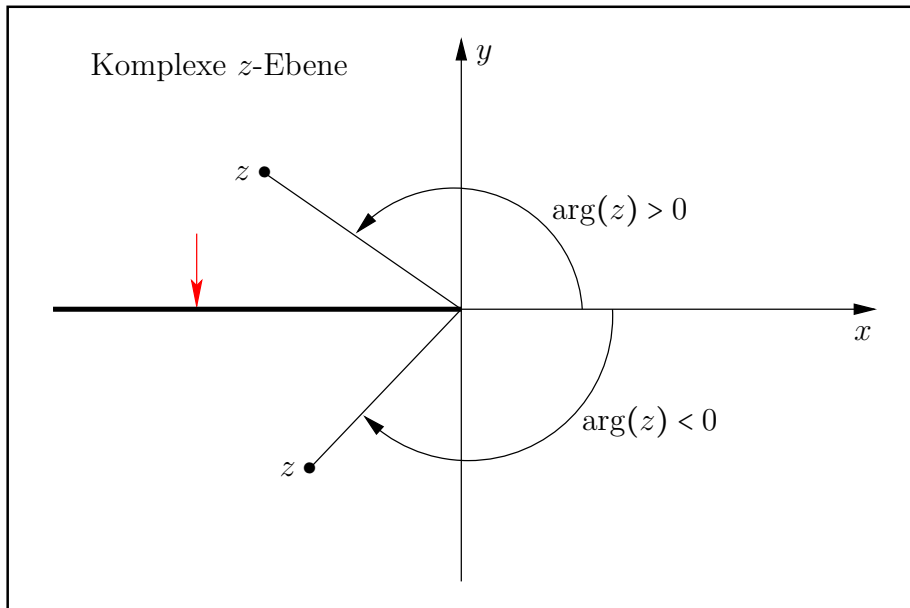


Abbildung C.4.: Verzweigungsschnitt von  $\arg(z)$ ,  $z \in \mathbb{C}$ .

Der **Pfeil** zeigt die Richtung an, aus der  $\arg(z)$  auf den Verzweigungsschnitt stetig ergänzt wird, d.h. es gilt z.B.  $\arg(-1 + 0 \cdot i) = +\pi$ .  $\arg(z)$  wird intern mithilfe der Funktion `mpfr_atan2(...)` aus der MPFR-Bibliothek direkt ausgewertet, so dass die gewünschten Rundungen sogar optimal berechnet werden können. Der gerundete Funktionswert  $\arg(z)$  wird bestimmt mit dem Funktionsaufruf

```
MpfrClass arg(const MpfcClass& z, RoundingMode rnd, PrecisionType prec);
```

Wird `prec` gesetzt, so wird  $\arg(z)$  mittels `rnd` in die gewählte Präzision gerundet und dann zurückgegeben. Ohne `prec` wird mittels `rnd` in die Current-Präzision gerundet. Ohne `prec` und ohne `rnd` wird  $\arg(z)$  mit dem Current-Rundungsmodus in die Current-Präzision gerundet.

In der folgenden Tabelle sind für einige spezielle  $z = x + i \cdot y$ -Werte die Funktionswerte  $\arg(z)$  zusammengestellt. Dabei bedeuten die angegebenen Werte wie  $\pi/2$  oder  $-3\pi/4$  die bezüglich des voreingestellten Current-Rundungsmodus entsprechend gerundeten Werte.

$x$	$y$	$\arg(z)$	$x$	$y$	$\arg(z)$	$x$	$y$	$\arg(z)$
0	0	0	0	$+\infty$	$\pi/2$	$-\infty$	+1	$\pi$
+1	0	0	0	-1	$-\pi/2$	$+\infty$	$+\infty$	$+\pi/4$
$+\infty$	0	0	0	$-\infty$	$-\pi/2$	$+\infty$	$-\infty$	$-\pi/4$
$+\infty$	+1	0	-1	0	$\pi$	$-\infty$	$+\infty$	$+3\pi/4$
$+\infty$	-1	0	$-\infty$	0	$\pi$	$-\infty$	$-\infty$	$-3\pi/4$
0	+1	$\pi/2$	$-\infty$	-1	$-\pi$	$x$ oder $y = \text{NaN}$		NaN

Abbildung C.5.: Wertetabelle für  $\arg(z)$ ,  $z \in \mathbb{C}$ .

### C.1.7. $\text{argp1}(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist das Argument von  $1 + z$ , d.h.  $\text{argp1}(z)$ , definiert durch das Bogenmaß des Winkels, den der Fahrstrahl vom Ursprung nach  $1 + z$  mit der positiven reellen Achse einschließt. Die komplexe Ebene ist dabei längs der negativen reellen Achse von  $-1$  bis  $-\infty$  aufgeschnitten.

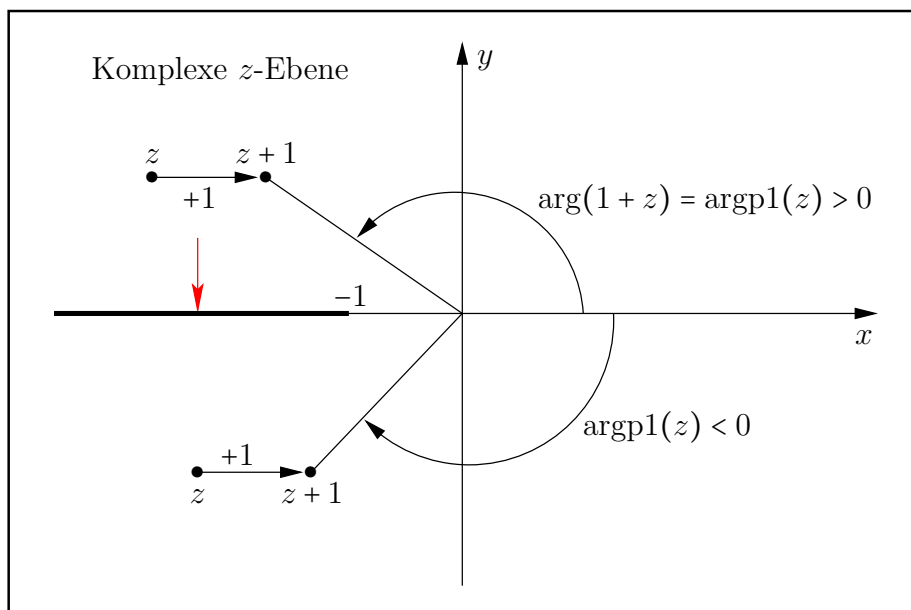


Abbildung C.6.: Verzweigungsschnitt von  $\text{argp1}(z)$ ,  $z \in \mathbb{C}$ .

Der **Pfeil** zeigt die Richtung an, aus der  $\text{argp1}(z)$  auf den Verzweigungsschnitt stetig ergänzt wird, d.h. es gilt z.B.  $\text{argp1}(-2 + 0 \cdot i) = +\pi$ . Der zu rundende exakte Funktionswert  $\arg(1 + z)$  wird bestimmt mit der Funktion

```
MpfrClass argp1(const MpfcClass& z, RoundingMode rnd, PrecisionType prec);
```

Wird `prec` gesetzt, so wird  $\arg(1 + z)$  mittels `rnd` in die gewählte Präzision gerundet und dann zurückgegeben. Ohne `prec` wird mittels `rnd` in die Current-Präzision gerundet. Ohne `prec` und ohne `rnd` wird  $\arg(1 + z)$  mit dem Current-Rundungsmodus in die Current-Präzision gerundet.

In der folgenden Tabelle sind für einige spezielle  $z = x + i \cdot y$ -Werte die Funktionswerte  $\arg(1 + z)$  zusammengestellt. Dabei bedeuten die angegebenen Werte wie  $\pi/2$  oder  $-3\pi/4$  die bezüglich des voreingestellten Current-Rundungsmodus (**RoundNearest**) entsprechend gerundeten Werte.

$x$	$y$	$\arg(1 + z)$	$x$	$y$	$\arg(1 + z)$	$x$	$y$	$\arg(1 + z)$
-1	-1	0	-1	$+\infty$	$\pi/2$	$-\infty$	+1	$\pi$
+1	0	0	-1	-1	$-\pi/2$	$+\infty$	$+\infty$	$+\pi/4$
$+\infty$	0	0	-1	$-\infty$	$-\pi/2$	$+\infty$	$-\infty$	$-\pi/4$
$+\infty$	+1	0	-2	0	$\pi$	$-\infty$	$+\infty$	$+3\pi/4$
$+\infty$	-1	0	$-\infty$	0	$\pi$	$-\infty$	$-\infty$	$-3\pi/4$
-1	+1	$\pi/2$	$-\infty$	-1	$-\pi$	$x$ oder $y = \text{NaN}$		NaN

Abbildung C.7.: Wertetabelle für  $\arg(1 + z)$ ,  $z \in \mathbb{C}$ .

### C.1.8. $|z|$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist der Betrag definiert durch

$$|z| := \sqrt{x^2 + y^2}.$$

Die obige Quadratwurzel wird mit der Funktion `mpfr_hypot(...)` aus der MPFR-Bibliothek direkt ausgewertet, so dass die gewünschten Rundungen sogar optimal berechnet werden können. Der Betrag von  $z$  wird bestimmt mit dem Funktionsaufruf

```
MpfrClass abs(const MpfcClass& z, RoundingMode rnd, PrecisionType prec);
```

Wird `prec` gesetzt, so wird  $|z|$  mittels `rnd` in die gewählte Präzision gerundet und dann zurückgegeben. Ohne `prec` wird mittels `rnd` in die Current-Präzision gerundet. Ohne `prec` und ohne `rnd` wird  $|z|$  mit dem Current-Rundungsmodus in die Current-Präzision gerundet.

### C.1.9. $\log(z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $\log(z)$  definiert durch

$$\log(z) := \log\left(\sqrt{x^2 + y^2}\right) + i \cdot \arg(z),$$

wobei der Realteil  $\log\left(\sqrt{x^2 + y^2}\right)$  mit der vordefinierten Funktion

```
MpfrClass ln_sqrtx2y2(const MpfrClass&x, const MpfrClass&y, RoundingMode rnd);
```

aus `mpfrclass.cpp` direkt ausgewertet werden kann.

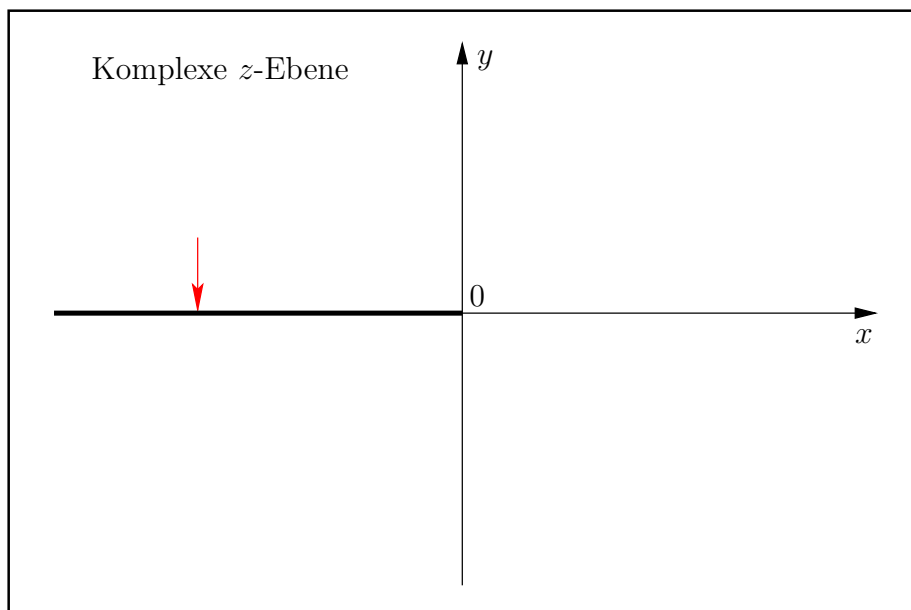


Abbildung C.8.: Verzweigungsschnitt von  $\log(z)$ ,  $z \in \mathbb{C}$ .

Der **Pfeil** gibt die Richtung an, aus der  $\log(z)$  auf den Verzweigungsschnitt analytisch fortgesetzt wird. Die Auswertung von  $\log(z)$  erfolgt mit

```
MpfcClass ln(const MpfcClass& z, RoundingMode rnd);
```

aus `mpfcclass.cpp`. In C-XSC wird die Logarithmusfunktion zur Basis  $e$  traditionsgemäß mit `ln(...)` bezeichnet.

### C.1.10. $1/z$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $1/z$  definiert durch

$$(C.7) \quad f(z) := \frac{1}{z} = u(x, y) + i \cdot v(x, y) = \frac{x}{x^2 + y^2} + i \cdot \frac{-y}{x^2 + y^2}, \quad x \neq 0 \wedge y \neq 0.$$

Die Funktionswerte  $f(z)$  werden nahezu optimal berechnet mit

```
MpfcClass reci(const MpfcClass& z, RoundingMode rnd);
```

Für `rnd` sind die folgenden Rundungsmodi zugelassen: `RoundUp`, `RoundDown`, `RoundNearest`. Wird `rnd` nicht gesetzt, so wird mit dem Current-Rundungsmodus gerundet.

#### C.1.10.1. Realteil

Für die Realteilmfunktion  $u(x, y)$  gilt mit  $x \neq 0 \wedge y \neq 0$ :

$$(C.8) \quad u(x, y) = \begin{cases} 0, & x = 0 \\ 1/x, & y = 0 \\ \frac{x}{x^2 + y^2}, & \text{sonst.} \end{cases}$$

Die Auswertung von  $u(x, y)$  erfolgt mit Hilfe der Funktion

```
MpfrClass x_div_x2py2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfrclass.cpp` definiert ist. Eine ausführliche Beschreibung der Funktion `x_div_x2py2` findet man auf Seite 185.

#### C.1.10.2. Imaginärteil

Für die Imaginärteilmfunktion  $v(x, y)$  gilt mit  $x \neq 0 \wedge y \neq 0$ :

$$(C.9) \quad v(x, y) = \begin{cases} 0, & y = 0 \\ 1/y, & x = 0 \\ \frac{-y}{x^2 + y^2}, & \text{sonst.} \end{cases}$$

Die Auswertung von  $v(x, y)$  erfolgt mit Hilfe der gleichen Funktion `x_div_x2py2(...)` durch den Aufruf

```
x_div_x2py2(-Im(z), Re(z), rnd);
```

Dabei ist zu beachten, dass der obige Aufruf nicht durch `-x_div_x2py2(Im(z), Re(z), rnd)` ersetzt werden darf, da das Vorzeichen `-` erst zum Schluss zur Anwendung kommt und daher die Rundungen bez. `rnd` vorher falsch berechnet wird.

#### Numerisches Beispiel:

Mit  $x = 2$  und  $y = -3$ , d.h. mit  $z = 2 - 3 \cdot i$ , und `prec = 400` liefern `reci(z, RoundUp)` und `reci(z, RoundDown)` den nahezu optimal auf- bzw. abgerundeten Wert von  $1/z$ :

```
(0.153846153846153846153846153846153846153846153846153846153846...153846169072,  
0.230769230769230769230769230769230769230769230769230769230769...230769254157)
```

mit gemeinsamen 120 Dezimalstellen bei Real- und Imaginärteil.

### C.1.11. $1/z^2$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $1/z^2$  definiert durch

$$(C.10) \quad f(z) := \frac{1}{z^2} = u(x, y) + i \cdot v(x, y) = \frac{x^2 - y^2}{(x^2 + y^2)^2} + i \cdot \frac{-2xy}{(x^2 + y^2)^2}, \quad x \neq 0 \wedge y \neq 0.$$

Die Funktionswerte  $f(z)$  werden nahezu optimal berechnet mit

```
MpfcClass reci_z2(const MpfcClass& z, RoundingMode rnd);
```

Für `rnd` sind die folgenden Rundungsmodi zugelassen: `RoundUp`, `RoundDown`, `RoundNearest`. Wird `rnd` nicht gesetzt, so wird mit dem Current-Rundungsmodus gerundet.

#### C.1.11.1. Realteil

Für die Realteilmfunktion  $u(x, y)$  gilt mit  $x \neq 0 \wedge y \neq 0$ :

$$(C.11) \quad u(x, y) = \frac{x^2 - y^2}{(x^2 + y^2)^2}.$$

Die Auswertung von  $u(x, y)$  erfolgt mit Hilfe der Funktion

```
MpfrClass Re_rz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfrclass.cpp` definiert ist. Eine ausführliche Beschreibung der Funktion `Re_rz2(...)` findet man auf Seite 186.

#### C.1.11.2. Imaginärteil

Für die Imaginärteilmfunktion  $v(x, y)$  gilt mit  $x \neq 0 \wedge y \neq 0$ :

$$(C.12) \quad v(x, y) = \frac{-2xy}{(x^2 + y^2)^2}$$

Die Auswertung von  $v(x, y)$  erfolgt mit Hilfe der Funktion

```
MpfrClass Im_rz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfrclass.cpp` definiert ist. Eine ausführliche Beschreibung der Funktion `mIm_rz2(...)` findet man auf Seite 187, wobei zu beachten ist, dass `mIm_rz2(...)` den Term  $+2xy/(x^2 + y^2)^2$  auswertet.

#### C.1.11.3. Numerische Ergebnisse

Im **1. Beispiel** wählen wir  $z = 1 + 2 \cdot i$  und erhalten für das exakte Ergebnis  $1/z^2 = -0.12 - 0.16 \cdot i$  mit `rnd = RoundDown` und der Präzision `prec = 300` den nahezu optimal abgerundeten Wert mit 93 ausgegebenen Dezimalstellen

$$(-1.200000000000 \dots 00000000002 \cdot 10^{-1}, -1.600000000000 \dots 00000000044 \cdot 10^{-1}).$$

Im **2. Beispiel** wählen wir  $z = 1 + 2 \cdot i$  und erhalten für das exakte Ergebnis  $1/z^2 = -0.12 - 0.16 \cdot i$  mit `rnd = RoundUp` und der Präzision `prec = 300` den nahezu optimal aufgerundeten Wert mit 93 ausgegebenen Dezimalstellen

$$(-1.199999999999 \dots 99999999941 \cdot 10^{-1}, -1.599999999999 \dots 99999999921 \cdot 10^{-1}).$$

### C.1.12. $1/(1+z^2)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $1/(1+z^2)$  definiert durch

$$(C.13) \quad f(z) := \frac{1}{1+z^2} = u+i \cdot v = \frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2} + i \cdot \frac{-2xy}{4x^2y^2+(1+x^2-y^2)^2}, \quad x \neq 0 \wedge y \neq \pm 1.$$

Die Funktionswerte  $f(z)$  werden nahezu optimal berechnet mit

```
MpfcClass reci_1pz2(const MpfcClass& z, RoundingMode rnd);
```

Für `rnd` sind die folgenden Rundungsmodi zugelassen: `RoundUp`, `RoundDown`, `RoundNearest`. Wird `rnd` nicht gesetzt, so wird mit dem Current-Rundungsmodus gerundet.

#### C.1.12.1. Realteil

Für die Realteilmfunktion  $u(x,y)$  gilt mit  $x \neq 0 \wedge y \neq \pm 1$ :

$$(C.14) \quad u(x,y) = \frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}.$$

Die Auswertung von  $u(x,y)$  erfolgt mit Hilfe der Funktion

```
MpfrClass Re_r1pz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfrclass.cpp` definiert ist. Eine ausführliche Beschreibung der Funktion `Re_r1pz2(...)` findet man auf Seite 189.

#### C.1.12.2. Imaginärteil

Für die Imaginärteilmfunktion  $v(x,y)$  gilt mit  $x \neq 0 \wedge y \neq \pm 1$ :

$$(C.15) \quad v(x,y) = \frac{-2xy}{4x^2y^2+(1+x^2-y^2)^2}$$

Die Auswertung von  $v(x,y)$  erfolgt mit Hilfe der Funktion

```
MpfrClass mIm_r1pz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfrclass.cpp` definiert ist. Eine ausführliche Beschreibung der Funktion

```
MpfrClass Im_r1pz2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

welche  $-v(x,y)$  auswertet, findet man auf Seite 188.

#### C.1.12.3. Numerische Ergebnisse

Im **1. Beispiel** wählen wir  $z = 2 + 3 \cdot i$ , und für den exakten Funktionswert  $1/(1+z^2) = -1/40 - 3/40 \cdot i$  erhalten wir mit `rnd = RoundDown` und der Präzision `prec = 300` den nahezu optimal abgerundeten Wert mit 93 ausgegebenen Dezimalstellen

$$(-2.500000000000 \dots 000000000030 \cdot 10^{-2}, -7.500000000000 \dots 000000000245 \cdot 10^{-2}).$$

Im **2. Beispiel** wählen wir  $z = 2 + 3 \cdot i$ , und für das exakte Ergebnis  $1/(1+z^2) = -1/40 - 3/40 \cdot i$  erhalten wir mit `rnd = RoundUp` und `prec = 300` den nahezu optimal aufgerundeten Wert mit 93 ausgegebenen Dezimalstellen

$$(-2.499999999999 \dots 99999999877 \cdot 10^{-2}, -7.499999999999 \dots 99999999631 \cdot 10^{-2}).$$

### C.1.13. $z^2$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $z^2$  definiert durch

$$z^2 := x^2 - y^2 + 2i \cdot x \cdot y.$$

Der Realteil  $x^2 - y^2$  wird direkt mit der vordefinierten Funktion

```
MpfrClass x2my2 (const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

aus `mpfrclass.cpp` ausgewertet. Um beim Imaginärteil einen vorzeitigen Unterlauf bei der Auswertung von  $x \cdot y$  zu vermeiden, wird zunächst das betragsmäßige Minimum von  $x, y$  mit 2 multipliziert, und erst dann erfolgt die Multiplikation mit dem betragsmäßig größeren zweiten Faktor. Um dies zu testen, wähle man z.B.  $x = \text{minfloat}()$  und  $y = 0.75$  und berechne damit den abgerundeten Wert  $z^2$  mit der Funktion

```
MpfcClass sqr (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist. Man erhält für den Imaginärteil von  $z^2$  den positiven Wert  $3.573\dots \cdot 10^{-323228497}$ , während schon das Produkt  $x \cdot y$  in  $2i \cdot (x \cdot y)$  einen Unterlauf verursacht.

### C.1.14. $\sqrt{z}$

Mit  $z = x + i \cdot y \in \mathbb{C}$  ist der Hauptwert der komplexen Quadratwurzel definiert durch

$$(C.16) \quad \sqrt{z} := \begin{cases} \sqrt{x} + i \cdot 0, & \text{falls } y = 0 \wedge x \geq 0 \\ 0 + i \cdot \sqrt{|x|}, & \text{falls } y = 0 \wedge x < 0 \\ \frac{\sqrt{2 \cdot (|z| + x)}}{2} + i \cdot \frac{y}{\sqrt{2 \cdot (|z| + x)}}, & \text{falls } y \neq 0. \end{cases}$$

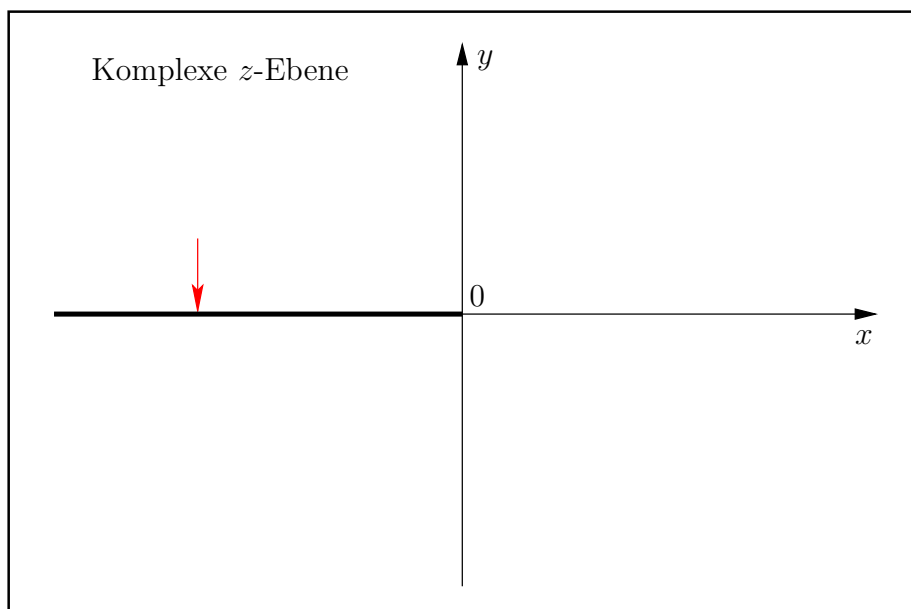


Abbildung C.9.: Verzweigungsschnitt von  $\sqrt{z}$ ,  $z \in \mathbb{C}$ .

Der **Pfeil** gibt die Richtung an, aus der  $\sqrt{z}$  auf den Verzweigungsschnitt analytisch fortgesetzt wird, d.h. z.B.  $\sqrt{-1 + 0 \cdot i} = +i$ .



Nach (C.16) können im Fall  $y \neq 0$  und  $x < 0$  bei der Auswertung der gerundeten Wurzel­ausdrücke  $\sqrt{2 \cdot (|z| + x)}$  wegen drohender Auslöschung starke Überschätzungen entstehen. Um diese Überschätzungen zu vermeiden, schreibt man im Fall negativer  $x$ -Werte  $x = -|x|$ , und das Erweitern mit  $(|z| + |x|)$  liefert

$$(C.17) \quad \Re(\sqrt{z}) = \begin{cases} \frac{\sqrt{2 \cdot (|z| + x)}}{2}, & \text{falls } y \neq 0 \wedge x \geq 0 \\ \frac{|y|}{\sqrt{2 \cdot (|z| + |x|)}}, & \text{falls } y \neq 0 \wedge x < 0. \end{cases}$$

$$(C.18) \quad \Im(\sqrt{z}) = \begin{cases} \frac{y}{\sqrt{2 \cdot (|z| + x)}}, & \text{falls } y \neq 0 \wedge x \geq 0 \\ \frac{\text{sign}(y)}{2} \cdot \sqrt{2 \cdot (|z| + |x|)}, & \text{falls } y \neq 0 \wedge x < 0. \end{cases}$$

Die reellen Wurzeln in (C.17) und (C.18) können jetzt ohne Auslöschung mit Hilfe der Funktion

```
MpfrClass Sqrt_zpx(const MpfrClass& x, const MpfrClass& y, const RoundingMode rnd);
```

aus der Datei `mpfcclass.cpp` ausgewertet werden, wobei mit dem Parameter  $x \geq 0$  jetzt keine negativen Werte übergeben werden. Es besteht aber noch ein weiteres Problem, denn mit der Abkürzung  $M := \text{MaxFloat}()$  gilt für die reelle Wurzel die Abschätzung

$$A := \sqrt{2 \cdot (|z| + x)} \leq \sqrt{2 \cdot (\sqrt{2}M^2 + M)} = \sqrt{2(\sqrt{2} + 1)} \cdot M < 3 \cdot \sqrt{M} < M,$$

so dass  $A$  im Zahlenformat stets darstellbar ist, während  $2 \cdot (|z| + x)$  durchaus einen vorzeitigen Überlauf verursachen kann. Um diesen Überlauf zu vermeiden, wird bei zu großem  $|x|$  oder  $|y|$  skaliert, d.h.  $x$  und  $y$  werden bei korrekter Rundung durch  $2^4$  dividiert, und die ausgewertete Wurzel wird dann am Ende zum Ausgleich wieder mit  $2^2 = 4$  rundungsfehlerfrei multipliziert, da wegen der obigen Abschätzung kein Überlauf eintreten kann.

Die Auswertung der gerundeten Werte von  $\sqrt{z}$  erfolgt mit Hilfe der Funktion

```
MpfcClass sqrt(const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist.

### C.1.15. $z^2 + a \cdot z + b$

Mit den Maschinenzahlen  $z, a, b \in \mathbb{C}$  berechnet die Funktion

```
MpfcClass poly2(const MpfcClass& z, const MpfcClass& a,
                const MpfcClass& b, RoundingMode rnd);
```

nahezu optimal gerundete Funktionswerte des Polynoms  $z^2 + a \cdot z + b$ , wobei für den Rundungsparameter `rnd` nur die Rundungsmodi `RoundDown`, `RoundUp`, `RoundNearest` erlaubt sind. Ohne `rnd` wird nach dem Current-Rundungsmodus gerundet; ist dieser jedoch nicht gesetzt, so wird mit `RoundNearest` gerundet. Die Implementierung erfolgt in `mpfcclass.cpp` mit Hilfe der gleichnamigen Intervallfunktion, wobei die komplexen Zahlen  $z, a, b$  als Punktintervalle übernommen werden. Weitere Einzelheiten und numerische Ergebnisse findet man auf Seite 291.

**C.1.16.**  $1/\sqrt{z}$

Mit  $z = x + i \cdot y \in \mathbb{C} \setminus \{0\}$  ist der Hauptwert der Funktion  $f(z) = 1/\sqrt{z}$  definiert durch

$$(C.19) \quad 1/\sqrt{z} := \begin{cases} \frac{1}{\sqrt{x}} + i \cdot 0, & \text{falls } y = 0 \wedge x > 0 \\ 0 - \frac{i}{\sqrt{-x}}, & \text{falls } y = 0 \wedge x < 0 \\ \frac{\sqrt{2 \cdot (|z| + x)}}{2 \cdot |z|} - i \cdot \frac{y}{|z| \cdot \sqrt{2 \cdot (|z| + x)}}, & \text{falls } y \neq 0. \end{cases}$$

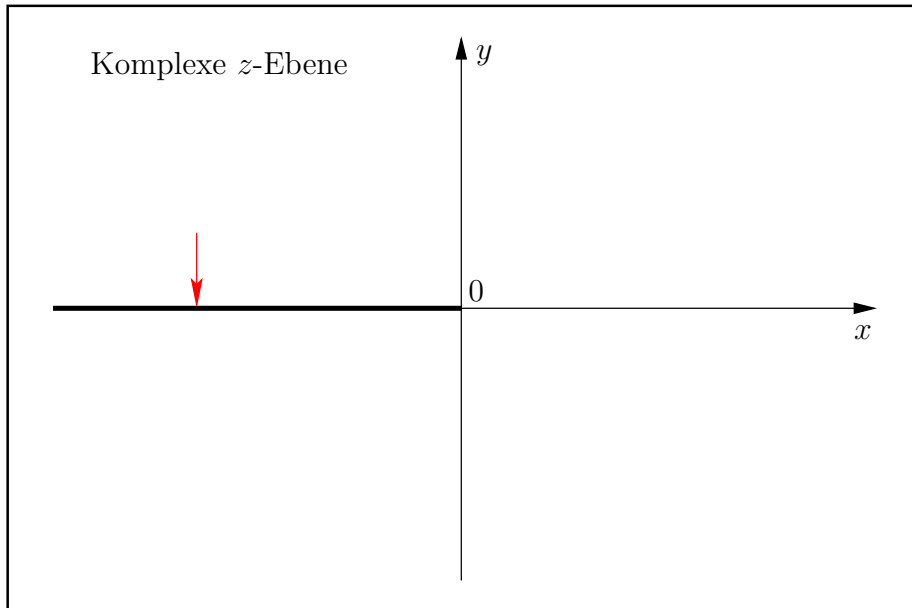


Abbildung C.10.: Verzweigungsschnitt von  $1/\sqrt{z}$ ,  $z \in \mathbb{C} \setminus \{0\}$ .

Der obige **Pfeil** gibt die Richtung an, aus der  $1/\sqrt{z}$  auf den Verzweigungsschnitt analytisch fortgesetzt wird, d.h. z.B.  $1/\sqrt{-1 + 0 \cdot i} = -i$ .

Nach (C.19) können im Fall  $y \neq 0$  und  $x < 0$  bei der Auswertung der gerundeten Wurzel-  
ausdrücke  $\sqrt{2 \cdot (|z| + x)}$  wegen drohender Auslöschung starke Überschätzungen entstehen. Um  
diese Überschätzungen zu vermeiden, schreibt man im Fall negativer  $x$ -Werte  $x = -|x|$ , und das  
Erweitern mit  $(|z| + |x|)$  liefert

$$(C.20) \quad \Re(1/\sqrt{z}) = \begin{cases} \frac{\sqrt{2 \cdot (|z| + x)}}{2 \cdot |z|}, & \text{falls } y \neq 0 \wedge x \geq 0 \\ \frac{|y|/|z|}{\sqrt{2 \cdot (|z| + |x|)}}, & \text{falls } y \neq 0 \wedge x < 0. \end{cases}$$

$$(C.21) \quad \Im(1/\sqrt{z}) = \begin{cases} -\frac{y/|z|}{\sqrt{2 \cdot (|z| + x)}}, & \text{falls } y \neq 0 \wedge x \geq 0 \\ -\frac{\text{sign}(y)}{2 \cdot |z|} \cdot \sqrt{2 \cdot (|z| + |x|)}, & \text{falls } y \neq 0 \wedge x < 0. \end{cases}$$

Die reellen Wurzeln in (C.20) und (C.21) können jetzt ohne Auslöschung und ohne vorzeitigen  
Über- oder Unterlauf mit Hilfe der Funktion

```
MpfrClass Sqrt_zpx(const MpfrClass& x, const MpfrClass& y, const RoundingMode rnd);
```

aus der Datei `mpfcclass.cpp` ausgewertet werden, wobei mit dem Parameter  $x \geq 0$  jetzt keine negativen Werte übergeben werden. Es besteht aber noch ein weiteres Problem, denn mit der Abkürzung  $M := \text{MaxFloat}()$  gilt für die reelle Wurzel die Abschätzung

$$A := \sqrt{2 \cdot (|z| + x)} \leq \sqrt{2 \cdot (\sqrt{2M^2} + M)} = \sqrt{2(\sqrt{2} + 1) \cdot M} < 3 \cdot \sqrt{M} < M,$$

so dass  $A$  im Zahlenformat stets darstellbar ist, während  $2 \cdot (|z| + x)$  durchaus einen vorzeitigen Überlauf verursachen kann. Um diesen Überlauf zu vermeiden, wird bei zu großem  $|x|$  oder  $|y|$  skaliert, d.h.  $x$  und  $y$  werden bei korrekter Rundung durch  $2^4$  dividiert, und die ausgewertete Wurzel wird dann am Ende zum Ausgleich wieder mit  $2^2 = 4$  rundungsfehlerfrei multipliziert, da wegen der obigen Abschätzung kein Überlauf eintreten kann. Ein weiterer möglicher Überlauf bei der Auswertung von C.20 bzw. C.21 wird vermieden, wenn man  $|x|, |y| < 2^{1073741821}$  verlangt, und dies ist wegen  $\text{expo}(\text{MaxFloat}()) = 1073741823$  keine wirkliche Einschränkung.

Die Auswertung der dann nahezu optimal gerundeten Werte von  $1/\sqrt{z}$  erfolgt mit Hilfe der Funktion

```
MpfcClass sqrt_r(const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` implementiert ist.

### C.1.17. $\sinh(z)$

Mit  $z = x + i \cdot y \in \mathbb{C}$  kann der hyperbolische Sinus wie folgt definiert werden

$$\sinh(z) := \mathfrak{I}(\sin(y + i \cdot x)) + i \cdot \mathfrak{R}(\sin(y + i \cdot x)).$$

Mit der bereits definierten  $\sin(z)$ -Funktion wird daher  $\sinh(z)$  wie folgt implementiert:

```
1  MpfcClass sinh (const MpfcClass& z, RoundingMode rnd)
2  // z = x + i*y;
3  {
4      MpfcClass z_(Im(z), Re(z), rnd, z.GetPrecision()); // z_ = y + i*x; Exakt!
5      z_ = sin(z_, rnd);
6      return MpfcClass( Im(z_), Re(z_) );
7  }
```

Zunächst wird in Zeile 4 das transformierte Argument  $z_ = y + i \cdot x$  mit der gleichen Präzision von  $z$  bestimmt. Mit diesem Argument wird in Zeile 5 zuerst der mittels `rnd` gerundete Wert  $\sin(z_)$  berechnet und in die Current-Präzision gerundet. Anschließend erfolgt die Wertübergabe an  $z_$  in der gleichen Current-Präzision. In Zeile 6 wird der gerundete Funktionswert  $\sinh(z)$  zurückgegeben.

### C.1.18. $\cosh(z)$

Mit  $z = x + i \cdot y \in \mathbb{C}$  kann der hyperbolische Cosinus wie folgt definiert werden

$$\cosh(z) := \cos(i \cdot z) = \cos(-y + i \cdot x).$$

Die Auswertung der gerundeten Werte von  $\cosh(z)$  erfolgt mit Hilfe der Funktion

```
MpfcClass cosh (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist.

### C.1.19. $\tanh(z)$

Mit  $z = x + i \cdot y \in \mathbb{C}$  kann der hyperbolische Tangens wie folgt definiert werden

$$\tanh(z) := \mathfrak{I}(\tan(y + i \cdot x)) + i \cdot \mathfrak{R}(\tan(y + i \cdot x)).$$

Die Auswertung der gerundeten Werte von  $\tanh(z)$  erfolgt mit Hilfe der Funktion

```
MpfcClass tanh (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist.

### C.1.20. $\coth(z)$

Mit  $z = x + i \cdot y \in \mathbb{C}$  kann der hyperbolische Cotangens wie folgt definiert werden

$$\coth(z) := -\mathfrak{I}(\cot(y + i \cdot x)) + i \cdot \mathfrak{R}(\cot(y + i \cdot x)).$$

Die Auswertung der gerundeten Werte von  $\coth(z)$  erfolgt mit Hilfe der Funktion

```
MpfcClass coth (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist.

### C.1.21. $\arcsin(z)$

Die mehrdeutige  $\arcsin$ -Funktion besitzt auf der reellen Achse die beiden Verzweigungspunkte  $(-1,0)$  und  $(+1,0)$ . Für den Hauptzweig gehen die Verzweigungsschnitte von  $-\infty$  bis  $-1$  und von  $+1$  bis  $+\infty$ .

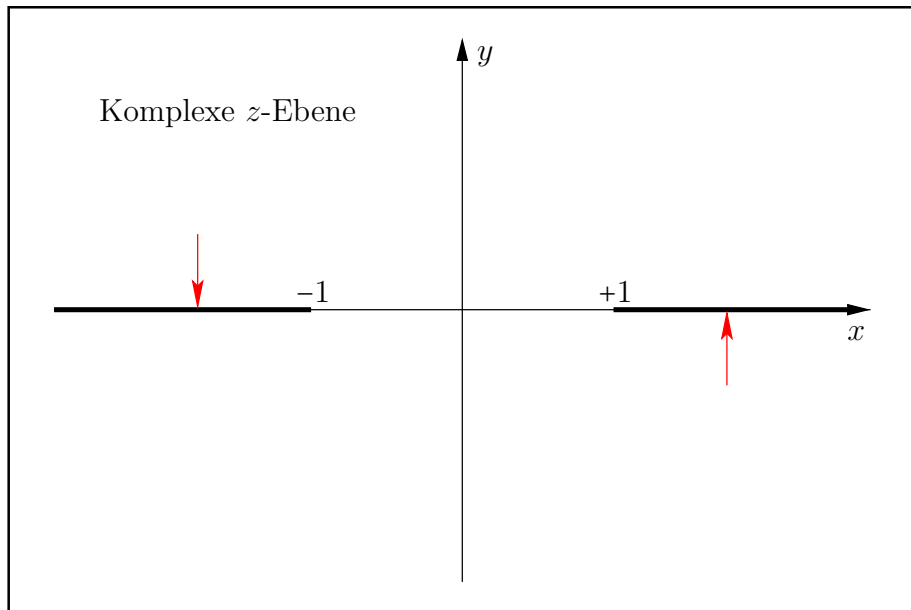


Abbildung C.11.: Verzweigungsschnitte von  $\arcsin(z)$ ,  $z \in \mathbb{C}$ .

Die **Pfeile** geben die jeweilige Richtung an, aus denen die Funktionswerte von  $\arcsin(z)$  auf den jeweiligen Verzweigungsschnitt analytisch fortgesetzt werden.

#### C.1.21.1. Realteil

Mit  $z = x + i \cdot y \in \mathbb{C}$ ,  $\arcsin(z) = u(x, y) + i \cdot v(x, y)$  und

$$(C.22) \quad T(x, y) := \frac{1}{2} \left( \sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2} \right)$$

gilt z.B. nach [38] für die Realteilmfunktion

$$(C.23) \quad \Re(\arcsin(z)) = u(x, y) = \arcsin \frac{x}{T}.$$

Da die reelle  $\arcsin$ -Funktion rechts in (C.23) streng monoton wächst, müssen für auf- bzw. abgerundete Funktionswerte  $u(x, y)$  die Argumente  $x/T(x, y)$  selbst auf- bzw. abgerundet werden. Um dies möglichst einfach realisieren zu können, wird  $T(x, y)$  zunächst mithilfe der Funktion<sup>1</sup>

```
MpfrClass asin_T(const MpfcClass& z, bool& scal, RoundingMode rnd);
```

so mit einem Rundungsparameter `rnd` implementiert, dass mit ihm auf- bzw. abgerundete Funktionswerte von  $T(x, y)$  berechnet werden können. Wegen

$$(C.24) \quad T(x, y) \equiv T(|x|, |y|) = \frac{1}{2} \left( \sqrt{(|x|+1)^2 + |y|^2} + \sqrt{(|x|-1)^2 + |y|^2} \right) \geq T(|x|, 0),$$

$$(C.25) \quad T(|x|, 0) = \frac{1}{2} \{ (|x|+1) + ||x|-1| \} = \begin{cases} |x|, & \text{falls } |x| \geq 1, \\ 1, & \text{falls } |x| \leq 1, \end{cases}$$

<sup>1</sup>Den Quellcode von `asin_T` findet man in `mpfcclass.cpp`.

folgt direkt  $T(x, y) \geq 1$ , und wegen  $T(x, y) \equiv T(|x|, |y|)$  kann  $T(x, y)$  mithilfe der Wurzelsumme in (C.24) ausgewertet werden. Die Voraussetzungen  $x \geq 0$  und  $y \geq 0$  stellen dabei für die korrekte Rundung von  $T(x, y)$  eine große Hilfe dar. Um einen vorzeitigen Überlauf bei der Berechnung der Wurzelsumme in (C.24) zu vermeiden, betrachten wir zunächst mit  $M := \max(|x|, |y|)$  die folgenden Abschätzungen:

$$\sqrt{(|x|+1)^2 + |y|^2} + \sqrt{(|x|-1)^2 + |y|^2} \leq 2\sqrt{(|x|+1)^2 + |y|^2} \leq 2\sqrt{(M+1)^2 + M^2} < 2\sqrt{2}(M+1).$$

Ein Überlauf wird also verhindert durch  $2\sqrt{2}(M+1) < \text{MaxFloat}()$  bzw. durch die Forderung

$$\text{expo}(M) < \text{expo}\left(\frac{\text{MaxFloat}()}{2\sqrt{2}} - 1\right) = 1073741821 =: p.$$

Der obige minimale Wert  $p = 1073741821$  wurde mit der kleinst-möglichen Current-Precision  $\text{prec} = 2$  berechnet. Schon mit  $\text{prec} \geq 4$  erhält man  $p_g = 1073741822$ . Um für alle Präzisionen einen Überlauf zu vermeiden, wird daher im Fall  $\text{expo}(M) \geq p$  geeignet skaliert, d.h. die Werte  $|y|$ ,  $(|x|+1)$ ,  $(|x|-1)$  werden durch 8 dividiert, wobei  $1/8$  selbst bei  $\text{prec} = 2$  exakt gespeichert wird. Im Fall  $\text{scal} == \text{true}$  muss daher der Rückgabewert von  $\text{asin}_T(\dots)$  noch mit 8 multipliziert werden, um den korrekten Wert von  $T(x, y)$  nach (C.24) zu erhalten. In (C.24) erfolgt die Auswertung der beiden Wurzeln mit den quadratischen Argumenten mithilfe der bereits implementierten Funktion  $\text{sqrtox2y2}(\dots)$ .

Die Auswertung der Realteilmfunktion  $u(x, y)$  erfolgt mithilfe von  $\text{asin}_T(\dots)$  in der Funktion

```
MpfrClass Re_asin (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist. Im Fall  $y = 0$  erhält man mit (C.23) und (C.25)

$$u(x, 0) = \begin{cases} \arcsin(x), & |x| \leq 1, \\ \arcsin(+1) = +\pi/2, & x \geq +1, \\ \arcsin(-1) = -\pi/2, & x \leq -1. \end{cases}$$

Im Fall  $y = 0$  kann  $u(x, 0)$  damit direkt, d.h. ohne die Funktion  $T(x, y)$ , ausgewertet werden. Im Fall  $y \neq 0$  wird in (C.23) zunächst der Quotient  $x/T(x, y)$  berechnet, wobei auf die korrekten Rundungen zu achten ist. Wenn bei der Auswertung von  $T(x, y)$  skaliert worden ist, so wird  $T$  jetzt nicht mit 8 multipliziert, weil dies zum Überlauf führen würde, sondern der Zähler  $x$  ist dafür durch 8 zu dividieren.

Mit (C.24) und (C.25) folgt für den Quotienten  $|x|/T(x, y)$  die Abschätzung

$$\frac{|x|}{T(x, y)} \leq \frac{|x|}{T(x, 0)} = \begin{cases} 1, & \text{falls } |x| \geq 1, \\ |x|, & \text{falls } |x| < 1, \end{cases} \implies \frac{|x|}{T(x, y)} \leq 1,$$

die natürlich erfüllt sein muss, wenn die reelle arcsin-Funktion in (C.23) mit diesem Quotienten ohne Fehlermeldung ausgewertet werden soll. Die gewünschten Rundungen bei der Auswertung des Quotienten können aber auf der Maschine  $|x|/T > 1$  zur Folge haben. Um die damit verbundenen Fehlermeldungen zu vermeiden, müssen daher vor Auswertung der reellen arcsin-Funktion die folgenden Abfragen erfolgen. Zur Abkürzung wird dabei der auf der Maschine berechnete und gerundete Quotient jetzt mit  $x$  bezeichnet.

```
if (x>1)
    x = 1;
if (x<-1)
    x = -1;
```

Alle weiteren Einzelheiten findet man im Quellcode der Funktion  $\text{Re\_asin}(\dots)$  in der Datei `mpfcclass.cpp`.

### C.1.21.2. Imaginärteil

Mit der in (C.22) bereits definierten Funktion  $T(x, y)$  gilt nach [38] für die Imaginärteilmfunktion

$$(C.26) \quad v(x, y) := \begin{cases} +\log(T + \sqrt{T^2 - 1}), & \text{falls } y > 0 \\ +\log(T + \sqrt{T^2 - 1}), & \text{falls } y = 0 \text{ und } x \leq -1 \\ 0, & \text{falls } y = 0 \text{ und } |x| \leq 1 \\ -\log(T + \sqrt{T^2 - 1}), & \text{falls } y = 0 \text{ und } x \geq +1 \\ -\log(T + \sqrt{T^2 - 1}), & \text{falls } y < 0. \end{cases}$$

Wegen  $T(x, y) \geq 1$  gilt zusätzlich

$$(C.27) \quad \operatorname{arcosh}(T) \equiv \log(T + \sqrt{T^2 - 1}),$$

so dass in (C.26) die Logarithmus-Funktionen mit dem komplizierten Argument  $T + \sqrt{T^2 - 1}$  im Bedarfsfall durch  $\operatorname{arcosh}(T)$  ersetzt werden können. Bei der Auswertung von  $\operatorname{arcosh}(T)$  treten jedoch zwei grundsätzliche Probleme auf:

1. Bei der Auswertung von  $T(x, y)$  mithilfe der Funktion `asin_T(...)` wird bei zu großem  $|x|$  oder  $|y|$  geeignet skaliert, so dass der Rückgabewert  $T$  noch mit 8 zu multiplizieren ist. In diesem Fall wäre  $\operatorname{arcosh}(8 \cdot T)$  auszuwerten, wodurch ebenfalls ein Überlauf entstehen würde. In diesem Fall muss man auf die rechte Seite von (C.27) zurückgreifen und den komplizierteren Ausdruck  $\log(8 \cdot T + \sqrt{64 \cdot T^2 - 1})$  weiter umformen, um den z.B. durch  $8 \cdot T$  drohenden Überlauf zu vermeiden.
2. Im Fall  $T(x, y) \rightarrow +1$  werden beide Funktionen in (C.27) in der Nähe ihrer gemeinsamen Nullstellen  $T_0(x, y) = 1$  ausgewertet, wobei wegen starker Auslöschung bei den Rundungen zu große Überschätzungen entstehen. Diese Überschätzungen werden vermieden, wenn man in  $\operatorname{arcosh}(T)$  das Argument  $T = 1 + r$  in zwei Summanden zerlegt und  $r := T - 1$  so umformt, dass bei seiner Auswertung nur minimale Überschätzungen entstehen können. Damit folgt

$$\operatorname{arcosh}(T) = \operatorname{arcosh}(1 + r) = \operatorname{acoshp1}(r),$$

wobei die letzte Funktion `acoshp1(r)` bereits implementiert ist, vgl. die Tabelle auf Seite 37.

Wir behandeln zunächst das **Problem 1**.

Im Skalierungsfall gilt mit  $T = 8 \cdot T$

$$\begin{aligned} \log(8 \cdot T + \sqrt{64 \cdot T^2 - 1}) &= \log(8 \cdot T + \sqrt{64 \cdot (T^2 - 1/64)}) \\ &= \log(8 \cdot T + 8 \cdot \sqrt{(T^2 - 1/64)}) \\ &= \log(8) + \log(T + \sqrt{(T^2 - 1/64)}) \\ &= 3 \cdot \log(2) + \log(T + \sqrt{T - 1/8} \cdot \sqrt{T + 1/8}). \end{aligned}$$

Jetzt kann die ganze letzte rechte Seite ohne Überlauf ausgewertet werden, wobei  $\log(2)$  als Konstante zur Verfügung steht.  $T$  ist der mit Skalierung berechnete Rückgabewert der Funktion `asin_T(...)`, und  $1/8$  kann auch noch bei der minimalen Current-Präzision `prec = 2` exakt gespeichert werden. Die Auswertung des letzten Ausdrucks oben rechts erfolgt mit der Funktion

```
MpfrClass acosh_T(const MpfcClass& z, RoundingMode rnd);
```

die in der Datei `mpfcclass.cpp` definiert ist.

Wir betrachten jetzt das **Problem 2**.

Der geometrische Ort aller Punkte mit  $T_0(x, y) = 1$  ist auf der  $x$ -Achse das Intervall  $[-1, +1]$ . In seiner Umgebung

$$U := \{(x, y) \in \mathbb{R}^2 \mid |x| < 1.125 \wedge |y| < 0.125\}$$

wird  $T$  zerlegt in  $T = 1 + r$ , womit die Beziehung  $\operatorname{arcosh}(T) = \operatorname{arcosh}(1 + r) = \operatorname{acoshp1}(r)$  zur Anwendung kommt. Die Schranken 1.125 und 0.125 wurden so gewählt, dass mit `prec = 53` beim Überschreiten der Umgebungsgrenze die Genauigkeit höchstens um eine Dezimalstelle abnimmt. Die nachfolgenden Rechnungen zeigen, dass die Auswertung innerhalb der Umgebung  $U$  deutlich aufwendiger sind als außerhalb. Die gewählten Schranken sollten daher möglichst klein sein, so dass die Wahl der obigen Werte einen guten Kompromiss darstellt. Bei größeren Präzisionen `prec > 53` wird beim Überschreiten der Umgebungsgrenze die Genauigkeit ebenfalls höchstens nur um eine Dezimalstelle kleiner.

Mit  $T = 1 + r$ ,  $r \geq 0$  gilt

$$\begin{aligned} 2r &= 2T(x, y) - 2 \\ &= \sqrt{(|x| + 1)^2 + y^2} + \sqrt{(|x| - 1)^2 + y^2} - 2 \\ &= (|x| + 1) \sqrt{1 + \left(\frac{y}{|x| + 1}\right)^2} - 2 + \sqrt{(|x| - 1)^2 + y^2} \\ &= (|x| + 1) \left\{ \sqrt{1 + \left(\frac{y}{|x| + 1}\right)^2} - 1 + 1 \right\} - 2 + \sqrt{(|x| - 1)^2 + y^2} \\ &= (|x| + 1) \left\{ \sqrt{1 + \left(\frac{y}{|x| + 1}\right)^2} - 1 \right\} + (|x| - 1) + \sqrt{(|x| - 1)^2 + y^2} \\ &= (|x| + 1) \cdot g\left(\left(\frac{|y|}{|x| + 1}\right)^2\right) + B, \quad g(t) := \sqrt{1 + t} - 1, \\ &= A + B. \end{aligned}$$

Der Ausdruck  $A$  wird mithilfe der Funktion `asin_rA(...)` ausgewertet, wobei die Funktion  $g(t)$  mithilfe der bereits implementierten, monoton wachsenden Funktion `sqrtp1m1(...)` direkt ausgewertet werden kann, vgl. die Tabelle auf Seite 37. Das Argument  $(|y|/(|x| + 1))^2$  kann jetzt mit nur ganz minimalen Überschätzungen auf- bzw. abgerundet werden. Der Ausdruck  $B$  wird am Anfang der Funktion `acoshp1_r(...)` wie folgt ausgewertet:

$$B = \begin{cases} |y|, & \text{falls } |x| = 1 \\ (|x| - 1) + \sqrt{(|x| - 1)^2 + y^2}, & \text{falls } |x| > 1 \\ (1 - |x|) \left\{ \sqrt{1 + \left(\frac{|y|}{1 - |x|}\right)^2} - 1 \right\}, & \text{falls } |x| < 1, \end{cases}$$

wobei auch jetzt der Ausdruck  $\{...\}$  mithilfe der Funktion `sqrtp1m1(...)` berechnet wird. Es ist zu beachten, dass in den Fällen  $|x| > 1$  und  $|x| < 1$  die entsprechenden Terme ohne Auslöschung berechnet werden. Lediglich im letzten Fall  $|x| < 1$  könnte man annehmen, dass für  $|x| \rightarrow +1$  das Argument  $(|y|/(1 - |x|))^2$  einen vorzeitigen Überlauf erzeugt. Erfreulicherweise ist diese Gefahr jedoch rein theoretisch, da sie nur eintritt, wenn die Current-Präzision in die Größenordnung von ca. `prec = 500000000` Bits kommt, was schon aus Laufzeitgründen völlig unrealistisch ist!



### C.1.22. $\arccos(z)$

Die mehrdeutige  $\arccos$ -Funktion besitzt auf der reellen Achse die beiden Verzweigungspunkte  $(-1, 0)$  und  $(+1, 0)$ . Für den Hauptzweig gehen die Verzweigungsschnitte wie bei der  $\arcsin$ -Funktion von  $-\infty$  bis  $-1$  und von  $+1$  bis  $+\infty$ , vgl. Abb. C.11 auf Seite 261. Für den Hauptwert gelten mit  $T(x, y)$  von Seite 261 die Beziehungen

$$(C.28) \quad \Re(\arccos(z)) = \arccos(x/T),$$

$$(C.29) \quad \Im(\arccos(z)) = -\Im(\arcsin(z)).$$

Der Realteil von  $\arccos(z)$  wird im Vergleich zu (C.23) ganz analog zum Realteil von  $\arcsin(z)$  berechnet. Im Gegensatz zur reellen  $\arcsin$ -Funktion ist die reelle  $\arccos$ -Funktion jedoch monoton fallend, so dass für die korrekten Rundungen das Argument  $x/T$  jetzt im Vergleich zur  $\arcsin$ -Funktion genau entgegengesetzt zu runden ist. Weitere Einzelheiten findet man in der Funktion `Re_acos(...)`, die in der Datei `mpfcclass.cpp` definiert ist.

Nach (C.29) unterscheidet sich der Imaginärteil von  $\arccos(z)$  vom Imaginärteil der  $\arcsin$ -Funktion nur durch das Vorzeichen, so dass für die korrekten Rundungen nach (A.3) und (A.4) von Seite 180 die Rundungsparameter bei den entsprechenden Funktionsaufrufen `Im_asin(...)` nur zu vertauschen sind. Weitere Einzelheiten findet man in der Funktion `Im_acos(...)`, die in der Datei `mpfcclass.cpp` definiert ist. Die Funktionswerte  $f(z) = \arccos(z)$  werden berechnet mit Hilfe der Funktion

```
MpfcClass acos (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist. Die Funktionswerte  $f(z)$ , mit  $z \in \mathbb{C} \setminus \{\pm 1\}$ , werden mit Hilfe des Parameters `rnd`  $\in \{\text{RoundDown}, \text{RoundUp}, \text{RoundNearest}\}$  entsprechend nahezu optimal gerundet. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus.

### C.1.23. $1/(1 - z^2)$

Setzt man  $\hat{z} := i \cdot z$ , so gilt  $f(z) := 1/(1 - z^2) = 1/(1 + \hat{z}^2)$ , d.h.  $f(z)$  kann mit Hilfe der bereits definierten Funktion `reci_r1pz2( $\hat{z}$ )` ausgewertet werden. Dies wird realisiert in der Funktion

```
MpfcClass reci_1mz2 (const MpfcClass& z, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist. Die Funktionswerte  $f(z)$ , mit  $z \in \mathbb{C} \setminus \{\pm 1\}$ , werden mit Hilfe des Parameters `rnd`  $\in \{\text{RoundDown}, \text{RoundUp}, \text{RoundNearest}\}$  entsprechend nahezu optimal gerundet. Wird `rnd` nicht gesetzt, so erfolgt die Rundung nach dem Current-Rundungsmodus.

### C.1.24. $\log(1+z)$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $\log(1+z)$  definiert durch

$$\log(1+z) := \log\left(\sqrt{(1+x)^2 + y^2}\right) + i \cdot \arg(1+z),$$

wobei der Realteil  $\log\left(\sqrt{(1+x)^2 + y^2}\right)$  mit der vordefinierten Funktion

```
MpfrClass ln_sqrtxp1_2y2(const MpfrClass& x, const MpfrClass& y,  
                        RoundingMode rnd);
```

aus `mpfrclass.cpp` direkt ausgewertet werden kann.

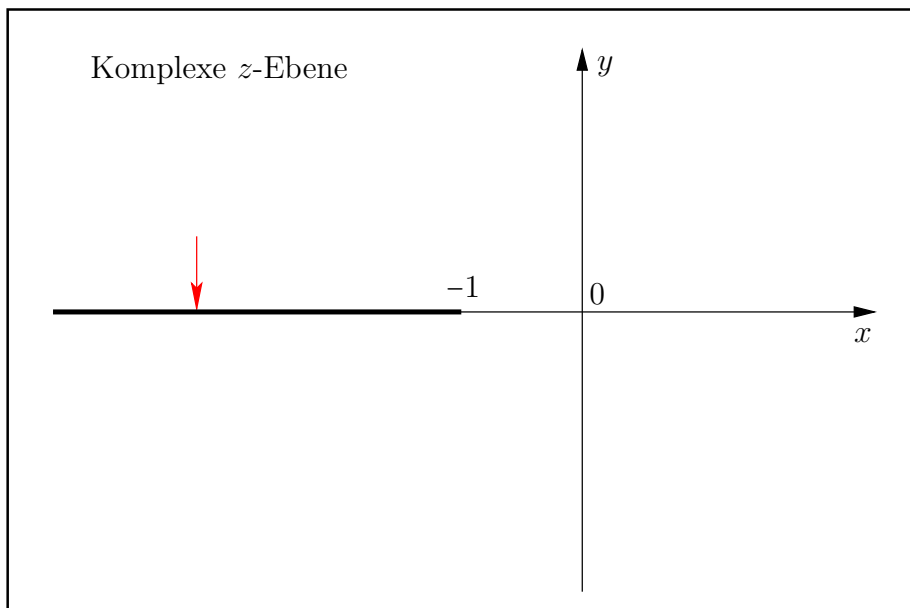


Abbildung C.12.: Verzweigungsschnitt von  $\log(1+z)$ ,  $z \in \mathbb{C}$ .

Der **Pfeil** gibt die Richtung an, aus der  $\log(1+z)$  auf den Verzweigungsschnitt analytisch fortgesetzt wird. Die Auswertung von  $\log(1+z)$  erfolgt mit der Funktion

```
MpfcClass lnp1(const MpfcClass& z, RoundingMode rnd);
```

aus `mpfcclass.cpp`. In C-XSC wird die Logarithmusfunktion zur Basis  $e$  traditionsgemäß mit `ln(...)` bezeichnet.

Mithilfe des Rundungsmodus `rnd` kann im Vergleich zum exakten Funktionswert auf- oder abgerundet bzw. in Richtung des exakten Funktionswertes gerundet werden. Diese Rundungen werden jedoch nicht immer optimal ausgeführt, d.h. der z.B. aufgerundete Wert ist nicht immer die erste Maschinenzahl rechts vom exakten Funktionswert<sup>2</sup>. Es wird jedoch garantiert, dass der zurückgegebene Wert rechts vom exakten Funktionswert liegt. Ganz analog wird mit `rnd = RoundNearest` nur garantiert, dass der Rückgabewert in der unmittelbaren Nähe des exakten Funktionswertes liegt. Beachten Sie, dass bei einer **optimalen** Rundung mit z.B. `rnd = RoundNearest` der Rückgabewert die zum exakten Funktionswert nächstgelegene Maschinenzahl ist. Beachten Sie außerdem, dass bei den beiden oberen Funktionen der Rundungsmodus `rnd` nicht gesetzt werden muss. In diesem Fall wird bez. des voreingestellten Current-Rundungsmodus gerundet.

<sup>2</sup>Bei diesem Beispiel wird angenommen, dass der exakte Funktionswert keine Maschinenzahl ist.

### C.1.24.1. Realteil

Die auszuwertende Realteilfunktion lautet

$$(C.30) \quad u(x, y) := \frac{1}{2} \cdot \ln \left( (1+x)^2 + y^2 \right) = \ln \left( \sqrt{(1+x)^2 + y^2} \right);$$

Zur Vermeidung eines vorzeitigen Überlaufs bei der Auswertung von  $\beta := (1+x)^2 + y^2$  berechnen wir die folgende Obergrenze

$$(1+x)^2 + y^2 \leq (1+|x|)^2 + |y|^2 < (1+|x|)^2 + (1+|y|)^2.$$

Mit  $M := \text{MaxFloat}(\text{prec})$  wird ein vorzeitiger Überlauf vermieden, wenn die beiden folgenden Bedingungen erfüllt sind:

$$(1+|x|)^2 < M/2 \wedge (1+|y|)^2 < M/2 \iff \\ |x| < \sqrt{M/2} - 1 \wedge |y| < \sqrt{M/2} - 1.$$

Wegen

$$(C.31) \quad 2^{k1} < \text{pred}(\sqrt{M/2}) \leq \sqrt{M/2} - 1, \quad k1 = 536870909$$

wird damit ein vorzeitiger Überlauf vermieden, wenn gilt

$$(C.32) \quad \text{expo}(x) < k1 \wedge \text{expo}(y) < k1, \quad k1 = 536870909.$$

Zu beachten ist, dass in (C.31)  $M$  eine von  $\text{prec}$  abhängige Größe ist und dass  $2^{k1}$  eine Unterschranke für alle  $\text{prec} \geq 2$  ist. Nach (C.32) kann für alle  $\text{prec} \geq 2$  ein Überlauf nur auftreten, wenn gilt:

$$(C.33) \quad \text{expo}(x) \geq k1 \text{ oder } \text{expo}(y) \geq k1, \quad k1 = 536870909.$$

Wir betrachten jetzt den Fall, dass mit (C.33) ein Überlauf bei der Auswertung von  $(1+x)^2 + y^2$  eintreten kann. Zur Vermeidung eines solchen Überlaufs betrachten wir die Umformung

$$\ln[(1+x)^2 + y^2] = \ln[2^{2k} \cdot 2^{-2k}((1+x)^2 + y^2)], \quad k \in \mathbb{N} \\ = 2k \cdot \ln(2) + \ln[(2^{-k} + x \cdot 2^{-k})^2 + (2^{-k} \cdot y)^2].$$

Für das Argument der letzten  $\ln$ -Funktion gilt die Abschätzung

$$\alpha := (2^{-k} + x \cdot 2^{-k})^2 + (2^{-k} \cdot y)^2 < (1+|x| \cdot 2^{-k})^2 + (1+|y| \cdot 2^{-k})^2.$$

Wählt man jetzt nach Seite 17 für  $x, y$  die für alle  $\text{prec} \geq 2$  gültige Obergrenze der größten positiven Zahl  $\text{MaxFloat}(\text{prec}) < 2^{1073741824}$ , so folgt

$$\alpha < 2 \cdot (1 + 2^{1073741824-k})^2.$$

Ein vorzeitiger Überlauf wird also vermieden, wenn man wieder nach Seite 17 verlangt

$$(C.34) \quad 2 \cdot (1 + 2^{1073741824-k})^2 < 2^{1073741823} \iff 1 + 2^{1073741824-k} < 2^{536870911},$$

und wegen  $1 + 2^{1073741824-k} < 2^{1073741825-k}$  ist die letzte Ungleichung in (C.34) erfüllt, wenn gilt

$$2^{1073741825-k} < 2^{536870911} \iff k > 536870914.$$

#### Zusammenfassung:

Im Fall (C.33) wird ein Überlauf bei der Auswertung von  $\beta$  vermieden, wenn  $u(x, y)$  mithilfe der Konstanten  $\text{Ln2}(\text{rnd})$  wie folgt berechnet wird, vgl. Seite 36:

$$(C.35) \quad u(x, y) = k \cdot \ln(2) + \frac{1}{2} \cdot \ln \left[ (2^{-k} + 2^{-k} \cdot x)^2 + (2^{-k} \cdot y)^2 \right], \quad k = 536870915.$$

Wir betrachten jetzt den Fall, dass bei der Berechnung des  $\ln$ -Arguments  $\beta := (1+x)^2 + y^2$  **kein Überlauf** eintreten kann. Bei der Auswertung von  $\ln(\beta)$  sind dann aber in Verbindung mit Abb. C.13 noch folgende Punkte zu beachten:

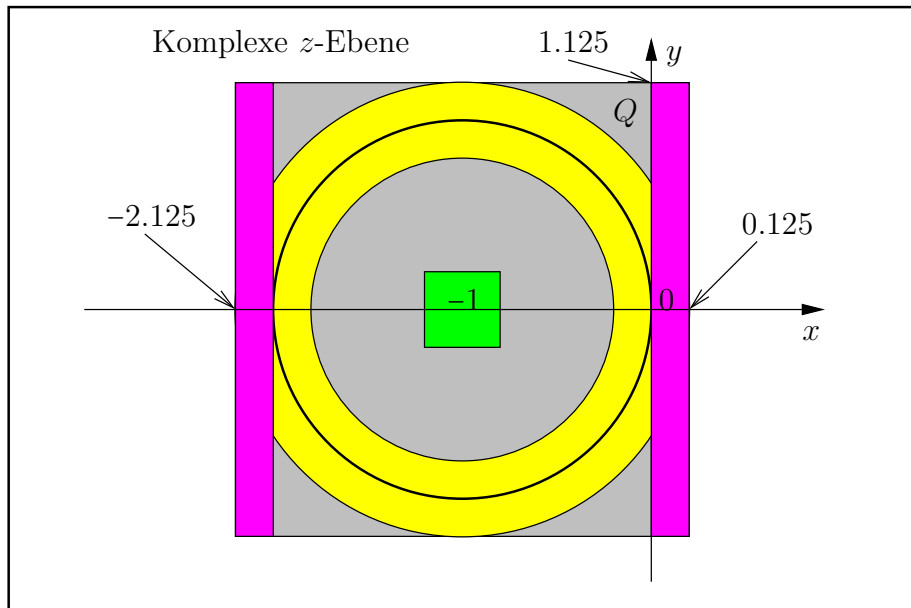


Abbildung C.13.: Verschiedene Bereiche zur Auswertung von  $u(x, y)$ .

1. Im grünen Quadrat mit dem Mittelpunkt  $M = (-1, 0)$  und der Kantenlänge  $2 \cdot 0,125 = 0,25$  gilt  $\ln((1+x)^2 + y^2)/2 = \ln(|1+x|) + \ln(1 + (y/(1+x))^2)/2$ , wobei in  $-1,25 \leq x \leq -0,75$  die Summe  $(1 \oplus x)$  für alle  $\text{prec} \geq 2$  **exakt** berechnet wird. Mit dem exakten Argument  $|1 \oplus x| \ll 1$  kann daher die Logarithmusfunktion in der Nähe ihrer Singularität 0 optimal ausgewertet werden. Würde man das Argument  $(1+x)^2 + y^2$  direkt auswerten, so müsste die Logarithmusfunktion in der Nähe ihrer Singularität 0 mit einem fehlerbehafteten Argument berechnet werden, was starke Überschätzungen zur Folge hätte. Die obige Summe der Logarithmusfunktionen führt nicht zu Auslöschungseffekten, da das grüne Quadrat vom gelben Kreisring weit genug entfernt ist.
2. Die in der Nähe der Nullstelle  $\beta = 1$  auftretende Auslöschung lässt sich vermeiden, wenn man  $\ln(\beta) = \ln(1 + \{x(2+x) + y^2\}) = \ln(1 + \delta)$  und  $\ln(1 + \delta)$  mit Hilfe der bereits implementierten Funktion  $\text{lnp1}(\delta)$  auswertet. Die Umgebung von  $\beta(x, y) = 1$  ist in Abb. C.13 der durch die beiden roten Rechtecke teilweise überdeckte gelbe Kreisring, der durch das schraffierte Quadrat  $Q$  mit dem Mittelpunkt  $M = (-1, 0)$  und der Kantenlänge 2,25 eingeschlossen wird. Die Punktmenge  $\beta(x, y) = 1$  ist der Kreis um  $M$  mit dem Radius 1 und berührt die  $y$ -Achse im Ursprung.
3. Der Fall  $x = 0$ , d.h.  $\delta = y^2$ , wird gesondert behandelt.
4. Im Fall  $|x| \ll 1 \wedge |y| \ll 1$  wird  $\delta$  intervallmäßig berechnet durch:  $\delta = x \cdot (2+x) + (y/x) \cdot y$ , um optimal gerundete  $\delta$ -Werte berechnen zu können.
5. In der Teilmenge  $-2 \leq x < 0 \wedge |y| < 1,125$  von  $Q$  besteht wegen  $x \cdot (2+x) < 0$  bei der Auswertung von  $\delta = x \cdot (2+x) + y^2$  die Gefahr starker Auslöschung, die durch schrittweise Verdoppelung der Current-Präzision bis zur exakten Berechnung von  $\delta$  vermieden wird.
6. In den beiden roten Teilmengen  $0 < x < 0,125 \wedge |y| < 1,125$  bzw.  $-2,125 < x < -2 \wedge |y| < 1,125$  von  $Q$  kann bei der Auswertung von  $\delta = x \cdot (2+x) + y^2$  keine Auslöschung auftreten, so dass  $\text{lnp1}(\delta)$  direkt ausgewertet werden kann. In beiden Teilmengen gilt:  $x \cdot (2+x) > 0$ .

7. Außerhalb des Quadrates  $Q$  und außerhalb des Überlaufbereiches wird  $\ln\{(1+x)^2 + y^2\}$  direkt ausgewertet. Weitere Einzelheiten findet man in der Datei `mpfrclass.cpp`.

### Numerische Ergebnisse:

Bezeichnungen:  $t = \ln(\sqrt{(1+x)^2 + y^2})$ ,  $\mathfrak{t} = \text{ln\_sqrtxp1\_2y2}(x, y, \text{rnd})$ ;

1. Mit `prec = 30000`, `rnd = RoundUp` und  $x = \text{succ}(-2) = -2 + 2^{-29999}$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Oberschranke mit 9030 Dezimalstellen:

$$t < \mathfrak{t} = -2.5186050871681829145830615704304081184503391605...47876326915e - 9031;$$

2. Mit `prec = 30000`, `rnd = RoundDown`,  $x = \text{succ}(-2) = -2 + 2^{-29999}$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Unterschranke mit 9030 Dezimalstellen:

$$t > \mathfrak{t} = -2.5186050871681829145830615704304081184503391605...47876326925e - 9031;$$

3. Mit `prec = 30000`, `rnd = RoundUp` und  $x = \text{pred}(-2) = -2 - 2^{-29998}$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Oberschranke mit 9030 Dezimalstellen:

$$t < \mathfrak{t} = 5.03721017433636582916612314086081623690067832100...95752653837e - 9031;$$

4. Mit `prec = 30000`, `rnd = RoundDown`,  $x = \text{pred}(-2) = -2 - 2^{-29998}$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Unterschranke mit 9030 Dezimalstellen:

$$t > \mathfrak{t} = 5.03721017433636582916612314086081623690067832100...95752653824e - 9031;$$

5. Mit `prec = 300000`, `rnd = RoundUp` und  $x = -1$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Oberschranke mit 90308 Dezimalstellen:

$$t < \mathfrak{t} = -7.4426111795489301787390319512589204791944768398...268012457556477e8;$$

6. Mit `prec = 300000`, `rnd = RoundDown` und  $x = -1$ ,  $y = \text{minfloat}(\text{prec})$  erhält man für  $t$  die nahezu optimale Unterschranke mit 90308 Dezimalstellen:

$$t > \mathfrak{t} = -7.4426111795489301787390319512589204791944768398...268012457556488e8;$$

7. Mit `prec = 300000`, `rnd = RoundUp` und  $x = -2^{-300001}$ ,  $y = 2^{-150000}$  erhält man für  $t$  die nahezu optimale Oberschranke mit 90308 Dezimalstellen:

$$t < \mathfrak{t} = 1.257510512341738730449956098921283616894008510311...06698175e - 180619;$$

8. Mit `prec = 300000`, `rnd = RoundDown` und  $x = -2^{-300001}$ ,  $y = 2^{-150000}$  erhält man für  $t$  die nahezu optimale Unterschranke mit 90308 Dezimalstellen:

$$t > \mathfrak{t} = 1.257510512341738730449956098921283616894008510311...06698174e - 180619;$$

Beachten Sie, dass bei den beiden letzten Beispielen wegen  $y^2 \approx -2x$  starke Auslöschung bei der Auswertung von  $x \cdot (2+x) + y^2$  auftritt, die Ober- und Unterschranken jedoch fast optimal berechnet werden. Bei zu kleinen absoluten Werten von z.B.  $x = -2^{-300001}$  und  $y = 2^{-150000}$  ist die optimale Einschließung bei zu kleiner Präzision jedoch nicht mehr möglich, da dann bei der Abrundung von  $2+x$  der zu kleine Wert 2 und damit die Unterschranke von  $t$  als 0 viel zu grob berechnet wird. Abhilfe schafft in diesem Fall nur eine ausreichend vergrößerte Präzision, mit der  $t$  dann wieder fast optimal eingeschlossen werden kann, vgl. Seite 298

In Ausnahmefällen kann also die optimale Einschließung eines Funktionswertes erst durch eine hinreichend große Präzision realisiert werden.

### C.1.24.2. Imaginärteil

Die Imaginärteilmfunktion  $v(x, y) := \arg(1 + z)$  wird ausführlich beschrieben ab Seite 251.

### C.1.24.3. Numerische Ergebnisse

Bezeichnungen:  $t = \ln(1 + z)$ ,  $z = x + i \cdot y \in \mathbb{C}$ ,  $\mathbf{t} = \text{lnp1}(z, \text{rnd})$ ;

1. Mit `prec = 30000`, `rnd = RoundDown` und  $x = -2$ ,  $y = \text{minfloat}()$  erhält man für  $t$  die komplexe Unterschranke<sup>3</sup>  $\mathbf{t}$ :

$$\mathbf{t} = (-0, 3.14159265358979323846264338327950288419716\dots536534940603402166544).$$

2. Mit `prec = 30000`, `rnd = RoundUp` und  $x = -2$ ,  $y = \text{minfloat}()$  erhält man für  $t$  die komplexe Oberschranke  $\mathbf{t}$ :

$$\mathbf{t} = (2.3825649048879\dots475741163e - 323228497, 3.14159265358979323\dots3402166545).$$

Die groben Unter- und Oberschranken des Realteils von  $t$  ergeben sich aus dem sehr kleinen Realteil:  $\Re(t) = \ln(1 + y^2)/2 \ll \text{minfloat}()$ . Mit  $x = 0$  erhält man für die Ober- und Unterschranken von  $\Re(t)$  die gleichen Ergebnisse.

3. Mit `prec = 30000`, `rnd = RoundDown` und  $x = \text{succ}(-1) = -1 + 2^{-\text{prec}}$ ,  $y = \text{succ}(+1) = 1 + 2^{-\text{prec}+1}$  erhält man für  $t$  die fast optimale komplexe Unterschranke

$$\mathbf{t} = (2.5186050871681\dots74787632691e - 9031, 1.5707963267948966\dots4703017010832719).$$

4. Mit `prec = 30000`, `rnd = RoundUp` und  $x = \text{succ}(-1) = -1 + 2^{-\text{prec}}$ ,  $y = \text{succ}(+1) = 1 + 2^{-\text{prec}+1}$  erhält man für  $t$  die fast optimale komplexe Oberschranke

$$\mathbf{t} = (2.5186050871681\dots74787632692e - 9031, 1.5707963267948966\dots4703017010832721).$$

5. Mit `prec = 30000`, `rnd = RoundDown` und  $x = -1$ ,  $y = -\text{minfloat}() = -2^{-1073741824}$  erhält man für  $t$  die fast optimale komplexe Unterschranke

$$\mathbf{t} = (-7.4426111795489301787\dots8304256133647e8, -1.5707963267948966\dots7010832724).$$

6. Mit `prec = 30000`, `rnd = RoundUp` und  $x = -1$ ,  $y = -\text{minfloat}() = -2^{-1073741824}$  erhält man für  $t$  die fast optimale komplexe Oberschranke

$$\mathbf{t} = (-7.4426111795489301787\dots8304256133634e8, -1.5707963267948966\dots7010832721).$$

7. Mit `prec = 30000`, `rnd = RoundDown` und  $x = y = 2 \cdot \text{minfloat}() = 2^{-1073741823}$  erhält man für  $t$  die fast optimale komplexe Unterschranke

$$\mathbf{t} = (4.765129809775902\dots32546e - 323228497, 4.765129809775902\dots32486e - 323228497).$$

8. Mit `prec = 30000`, `rnd = RoundUp` und  $x = y = 2 \cdot \text{minfloat}() = 2^{-1073741823}$  erhält man für  $t$  die fast optimale komplexe Oberschranke

$$\mathbf{t} = (4.765129809775902\dots32726e - 323228497, 4.765129809775902\dots32606e - 323228497).$$

---

<sup>3</sup>Die komplexe Unterschranke bedeutet die jeweilige Unterschranke für den Real- bzw. Imaginärteil.

### C.1.25. $e^z - 1$

Wir betrachten die Aufgabe, den Funktionsterm  $f(z) = e^z - 1$ ,  $z \in \mathbb{C}$ , für  $|z| \rightarrow 0$  so auszuwerten, dass sich die auf- und abgerundeten Funktionswerte bei einer vorgegebenen Präzision `prec` nur wenig unterscheiden. Wegen der Taylorentwicklung  $e^z - 1 = z + z^2/2! + \dots$  gilt  $f(z) \approx z$  für  $z \rightarrow 0$ , so dass mit  $e^z \approx 1$  bei naiver Auswertung von  $f(z)$  beim Realteil starke Auslöschung zu erwarten ist, die beim verlangten Auf- und Abrunden zu stark unterschiedlichen Funktionswerten führen wird. Mit  $z = x + i \cdot y$  betrachten wir dazu das Beispiel  $x = y = 2^{-400} = 3.8725 \dots \cdot 10^{-121}$ , wobei die verlangten Rundungen durch die naive Intervallauswertung von  $f(z)$  mit `prec = 53` realisiert werden. Man erhält für den Realteil die erwartete grobe Einschließung:

$$f(z) \in ([-1.11\dots16e - 16, 2.22\dots32e - 16], [3.87259191484931e - 121, 3.87259191484932e - 121]).$$

Um für den Realteil eine Einschließung mit einer Genauigkeit von 53 Bits berechnen zu können, muss  $e^z$  mit einer Präzision von mindestens `prec = 400 + 53 = 453` Bits ausgewertet werden, um nach der Subtraktion von 1 für die Differenz  $e^z - 1$  noch 53 korrekte Bits zu erhalten. Bei einer Ausgabe von 16 Dezimalstellen erhält man jetzt für den Realteil die viel bessere Einschließung:

$$f(z) \in ([3.87259\dots931e - 121, 3.87259\dots932e - 121], [3.87259\dots931e - 121, 3.87259\dots932e - 121]).$$

Mögliche, bei der Auswertung arithmetischer Ausdrücke auftretende, Auslöschungseffekte lassen sich mit einer hinreichend großen Präzision vermeiden, wobei jedoch deutlich höhere Laufzeiten in Kauf zu nehmen sind.

Wählen wir daher im oberen Beispiel jetzt  $x = y = 2^{-400000} = 1.0040016 \dots \cdot 10^{-120412}$ , so müsste  $f(z)$  mit einer Präzision von `prec = 400053` Bits ausgewertet werden, was etwa 120428 Dezimalstellen entspricht und eine sehr hohe Laufzeit erfordert.

Um diese hohen Laufzeiten zu vermeiden, wird der Realteil  $\Re(f(z)) = e^x \cdot \cos(y) - 1$  im Fall<sup>4</sup>  $\cos(y) > 0$  wie folgt berechnet:

$$(C.36) \quad \Re(f(z)) = e^x \cdot \cos(y) - 1 = e^{x + \ln(\cos(y))} - 1, \quad \cos(y) > 0,$$

wobei die obere rechte Seite mit der bereits implementierten Funktion `expm1(x + ln(cos(y)))` auszuwerten ist. `ln(cos(y))` wird mit der Funktion `ln_cos(y)` berechnet, und bei der Auswertung des Exponenten  $x + \ln(\cos(y))$  ist wieder mögliche Auslöschung zu beachten, die jedoch in fast allen Fällen mit einer **nur** doppelten Präzision vermieden werden kann!

Für Maschinenzahlen  $z = x + i \cdot y \in \mathbb{C}$  erhält man auf- bzw. abgerundete Funktionswerte  $e^z - 1$  mit Hilfe der Funktion

```
MpfcClass expm1(const MpfcClass& z, RoundingMode rnd);
```

wobei für `rnd` die Rundungen `RoundDown`, `RoundUp`, `RoundNearest` zur Verfügung stehen. Wird der Rundungsparameter `rnd` nicht gesetzt, so wird der voreingestellte Current-Rundungsmodus benutzt.

Numerische Beispiele findet man ab Seite 272.

---

<sup>4</sup>Für  $\cos(y) \leq 0$  kann bei der Auswertung von  $e^x \cdot \cos(y) - 1$  keine Auslöschung auftreten!

### C.1.25.1. Realteil

Mit  $z = x + i \cdot y \in \mathbb{C}$  lautet die auszuwertende Realteilmfunktion

$$(C.37) \quad \Re(e^z - 1) = \begin{cases} \cos(y) - 1 = -2 \cdot \sin^2(y/2), & \text{falls } x = 0, \\ e^x \cdot \cos(y) - 1, & \text{falls } \cos(y) \leq 0, \\ e^{x+\ln(\cos(y))} - 1, & \text{falls } \cos(y) > 0. \end{cases}$$

Im Fall  $\cos(y) > 0$  wird der Exponent  $x + \ln(\cos(y))$  intervallmäßig ausgewertet durch

$$(C.38) \quad x + \ln(\cos(y)) \in \mathbf{u} := [x] \oplus \ln\_cos([y]);$$

Zur Berechnung eines **abgerundeten** bzw. **aufgerundeten** Funktionswertes wird die reelle Intervallfunktion `expm1()` wegen der Monotonie der Exponentialfunktion wie folgt aufgerufen:

$$\begin{aligned} (e^{x+\ln(\cos(y))} - 1)_d &= \text{expm1}(\text{Inf}(\mathbf{u}), \text{RoundDown}); \\ (e^{x+\ln(\cos(y))} - 1)_u &= \text{expm1}(\text{Sup}(\mathbf{u}), \text{RoundUp}); \end{aligned}$$

Um eine mögliche Auslöschung bei der Auswertung der rechten Seite von (C.38) zu vermeiden, wird in der Hilfsfunktion `MpfiClass expm1_u(const MpfcClass& z)` die Current-Präzision schrittweise verdoppelt, bis der relative Durchmesser von  $\mathbf{u}$  den Wert  $2^{-\text{prec}}$  unterschritten hat, wobei `prec` die ursprüngliche Current-Präzision ist. Weitere Einzelheiten findet man in der Datei `mpfcclass.cpp` bei der Funktion `expm1_u(...)`.

### C.1.25.2. Imaginärteil

Wegen  $\Im(e^z - 1) = \Im(e^z) = e^x \cdot \sin(y)$  wird der Imaginärteil von  $e^z - 1$  wie bei der Exponentialfunktion berechnet.

Es soll noch auf ein mögliches Problem aufmerksam gemacht werden. Für  $x \rightarrow +\infty$  und  $y \rightarrow 0$  kann bei  $e^x$  ein vorzeitiger Overflow auftreten, während das Produkt  $e^x \cdot \sin(y)$  selbst noch im MPFR-Zahlenformat darstellbar ist. Man könnte das Problem im Fall  $\sin(y) > 0$  mit Hilfe der Darstellung

$$(C.39) \quad e^x \cdot \sin(y) = e^{x+\ln(\sin(y))}, \quad \sin(y) > 0$$

lösen, indem man  $x + \ln(\sin(y))$  intervallmäßig in höherer Präzision auswertet, um mögliche Auslöschung zu vermeiden. Damit hätte man dann zwar einen Überlauf beim Imaginärteil vermieden, aber ein Überlauf bei der Auswertung des Realteils  $e^x \cdot \cos(y) \approx e^x$  ist dann jedoch unvermeidbar, so dass die Auswertung nach (C.39) nicht zur Anwendung kommt.

### C.1.25.3. Numerische Ergebnisse

Nach (C.38) betrachten wir zunächst den Fall  $\cos(y) > 0$  und wählen mit  $y = 1$  für die Maschinenzahl  $x$  die 54-stellige Näherung

$$x \approx -\ln(\cos(1)) \approx \text{str} := 0.615626470386014262147037516408891863350935423946372834;$$

Mit der Präzision `prec = 180` wird dann `str` zur nächsten Maschinenzahl  $x$  gerundet. Für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = e^z - 1$  erhalten wir

$$\begin{aligned} f(z)_d &= (1.246714622829803823121588...205e - 55, 1.557407724654902230506974...003); \\ f(z)_u &= (1.246714622829803823121588...206e - 55, 1.557407724654902230506974...005). \end{aligned}$$

Man erkennt, dass die obigen 54-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 53 Dezimalstellen übereinstimmen.



Im nächsten Beispiel wählen wir mit  $x = y = 2^{-400000}$  extrem kleine, positive Werte, so dass bei der *naiven* Auswertung des Realteils von  $f(z) = e^z - 1$  starke Auslöschung auftreten muss. Mit der Präzision `prec = 180` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = e^z - 1$  die Ergebnisse

$$\begin{aligned} f(z)_d &= (1.00400160603542656243...179e - 120412, 1.00400160603542656243...179e - 120412); \\ f(z)_u &= (1.00400160603542656243...180e - 120412, 1.00400160603542656243...180e - 120412). \end{aligned}$$

Beachten Sie, dass wir jetzt positive und extrem kleine Real- und Imaginärteilwerte erhalten, die jedoch nahezu optimal gerundet werden konnten; vergleichen Sie dazu auch die Bemerkungen von Seite 271.

Im folgenden Beispiel wählen wir  $x = 40$  und  $y = 2$ , so dass  $\cos(y)$  negativ wird und damit  $\Re(e^z - 1) = e^x \cdot \cos(y) - 1$  ohne Auslöschung direkt ausgewertet werden kann. Mit der Präzision `prec = 180` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = e^z - 1$  die Ergebnisse

$$\begin{aligned} f(z)_d &= (-9.795483416403085235753036...643e16, 2.1403521744757874043519614...556e17); \\ f(z)_u &= (-9.795483416403085235753036...640e16, 2.1403521744757874043519614...557e17). \end{aligned}$$

Im letzten Beispiel gilt  $x = 0$  und  $y = 2^{400000}$ , so dass der Realteil von  $e^z - 1$  ausgewertet wird mit:  $\Re(e^z - 1) = -2 \cdot \sin^2(y/2)$ . Mit der Präzision `prec = 180` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = e^z - 1$  die Ergebnisse

$$\begin{aligned} f(z)_d &= (-1.475723472738204366761703790...877, -8.795948939630691868538814...020e - 1); \\ f(z)_u &= (-1.475723472738204366761703790...875, -8.795948939630691868538814...018e - 1). \end{aligned}$$

Auch hier erkennt man, dass die berechneten Werte fast optimal gerundet sind. Beachten Sie, dass die Auswertung von  $\sin(2^{399999})$  mit sehr hohen Laufzeiten (ca. 13 Minuten!) verbunden ist, da die notwendige Argumentreduktion mit Hilfe von  $\pi$  in sehr hoher Präzision erfolgen muss.

### C.1.26. $\sqrt{z+1} - 1$

Wir betrachten die Aufgabe,  $f(z) = \sqrt{z+1} - 1$ ,  $z = x + i \cdot y \in \mathbb{C}$ , insbesondere für  $|z| \rightarrow 0$  so auszuwerten, dass sich die auf- und abgerundeten Funktionswerte bei einer vorgegebenen Präzision `prec` möglichst wenig unterscheiden. Wenn man  $f(z)$  z.B. für  $z = 2^{-30000} + i \cdot 2^{-30000}$  bei einer Präzision `prec = 2000` naiv auswertet und dabei zur nächsten Rasterzahl rundet, so erhält man für den Maschinenwert  $\tilde{f}(z)$

$$\tilde{f}(z) = (0, 6.296512717920457\dots 7933 \cdot 10^{-9032}).$$

Vergleicht man mit der Taylorreihe  $f(z) = z/2 - z^2/8 \pm \dots$ , so erwartet man die Näherung

$$f(z) \approx (6.296512717920457\dots \cdot 10^{-9032}, 6.296512717920457\dots \cdot 10^{-9032}),$$

so dass bei der naiven Auswertung von  $f(z)$  der Realteil `0` viel zu grob ausfällt. Wird bei der naiven Berechnung jedoch aufgerundet, so erhält man für den Realteil den viel zu großen Wert  $1.7419619\dots \cdot 10^{-602}$ . Um nahezu optimale Ergebnisse zu erhalten, muss man daher einen anderen Algorithmus benutzen. Ausgehend von der Aufspaltung in Real- und Imaginärteil von  $\sqrt{z}$  in [17, 18] erhält man mit  $z = x + i \cdot y$  zunächst für den Hauptwert von  $\sqrt{z+1} - 1$  die Aufspaltung

$$(C.40) \quad \sqrt{z+1} - 1 = \begin{cases} \sqrt{x+1} - 1 + 0 \cdot i, & \text{falls } y = 0 \wedge x \geq -1 \\ -1 + i \cdot \sqrt{|x|-1}, & \text{falls } y = 0 \wedge x < -1 \\ \sqrt{\frac{|z+1|+x+1}{2}} - 1 + i \cdot \text{sign}(y) \cdot \sqrt{\frac{|z+1|-x-1}{2}}, & \text{falls } y \neq 0. \end{cases}$$

Der Fall  $y = 0 \wedge x \geq -1$  ist trivial, da  $\sqrt{x+1} - 1$  mithilfe der bereits implementierten Funktion `sqrtp1m1()` direkt ausgewertet werden kann. Im obigen zweiten Fall  $y = 0 \wedge x < -1$  kann  $\sqrt{|x|-1}$  bez. der Rundungen ebenfalls einfach implementiert werden, da die Wurzelfunktion monoton wachsend ist und  $|x|-1$  natürlich ohne Auslöschung berechnet werden kann, da  $x$  als exakte Maschinenzahl anzusehen ist.

Im letzten Fall  $y \neq 0$  wird Auslöschung möglich, wenn im Real- und Imaginärteil  $x+1$  negativ bzw. positiv gewählt werden. Die obige Darstellung in (C.40) ist jedoch vorteilhaft, wenn das Monotonieverhalten mithilfe der partiellen Ableitungen auf den Rändern der Rechteckintervalle für Real- und Imaginärteil zu untersuchen ist. Um im Fall  $y \neq 0$  die auftretenden Wurzelausdrücke ohne Auslöschung auswerten zu können, erweitert man analog zu Seite 257 und erhält die Ergebnisse

$$(C.41) \quad \Re(\sqrt{z+1} - 1) = \begin{cases} \frac{\sqrt{2 \cdot (|z+1| + x + 1)}}{2} - 1, & \text{falls } y \neq 0 \wedge x \geq -1 \\ \frac{|y|}{\sqrt{2 \cdot (|z+1| + |x+1|)}} - 1, & \text{falls } y \neq 0 \wedge x < -1. \end{cases}$$

$$(C.42) \quad \Im(\sqrt{z+1} - 1) = \begin{cases} \frac{y}{\sqrt{2 \cdot (|z+1| + x + 1)}}, & \text{falls } y \neq 0 \wedge x \geq -1 \\ \frac{\text{sign}(y)}{2} \cdot \sqrt{2 \cdot (|z+1| + |x+1|)}, & \text{falls } y \neq 0 \wedge x < -1, \end{cases}$$

wobei die auftretenden Wurzeln jetzt ohne Auslöschung berechnet werden können. Beachten Sie aber, dass beim Realteil wegen der Subtraktion der Eins immer noch Auslöschung entstehen kann, die gegebenenfalls durch schrittweise Verdoppelung der internen Präzision vermieden werden kann.

Man könnte auch versuchen, den Realteil von  $\sqrt{z+1} - 1$  mit Hilfe des äquivalenten Terms

$$(C.43) \quad \sqrt{z+1} - 1 = \frac{z}{\sqrt{z+1} + 1}$$

intervallmäßig auszuwerten. Dabei übersieht man jedoch, dass bei der auftretenden komplexen Intervalldivision mögliche Auslöschungen bei der Realteilberechnung ebenfalls auftreten, so dass die Berechnung des Realteils nach (C.41) aus Laufzeitgründen viel vorteilhafter ist, weil, wie in (C.43), der Imaginärteil hier nicht automatisch mitberechnet werden muss.

Im nächsten Schritt wird untersucht, wo in der komplexen  $z$ -Ebene bei der Berechnung des Realteils nach (C.41) mit Auslöschungseffekten zu rechnen ist. Dazu setzt man in den Fällen  $y \neq 0$  mit  $x \geq -1$  und  $x < -1$  die entsprechenden Terme in (C.41) gleich Null und erhält nach einfachen Rechnungen für den geometrischen Ort, an dem Auslöschung bei der Berechnung des Realteils von  $\sqrt{z+1}-1$  eintritt, das Ergebnis

$$(C.44) \quad y^2 = -4x, \quad x \leq 0.$$

Natürlich wird Auslöschung nicht nur auf der Menge  $\alpha := \{(x, y) \mid y^2 = -4x \wedge x \leq 0\}$  selbst, sondern auch in einer ganzen Umgebung von  $\alpha$  auftreten. In der folgenden Abbildung ist diese **Umgebung** angegeben, in der bei der Berechnung des Realteils nach (C.41) Auslöschung durch schrittweise Verdoppelung der Präzision zu vermeiden ist. Der **Pfeil** gibt an, aus welcher Richtung die Funktionswerte auf den **Verzweigungsschnitt** analytisch fortgesetzt werden.

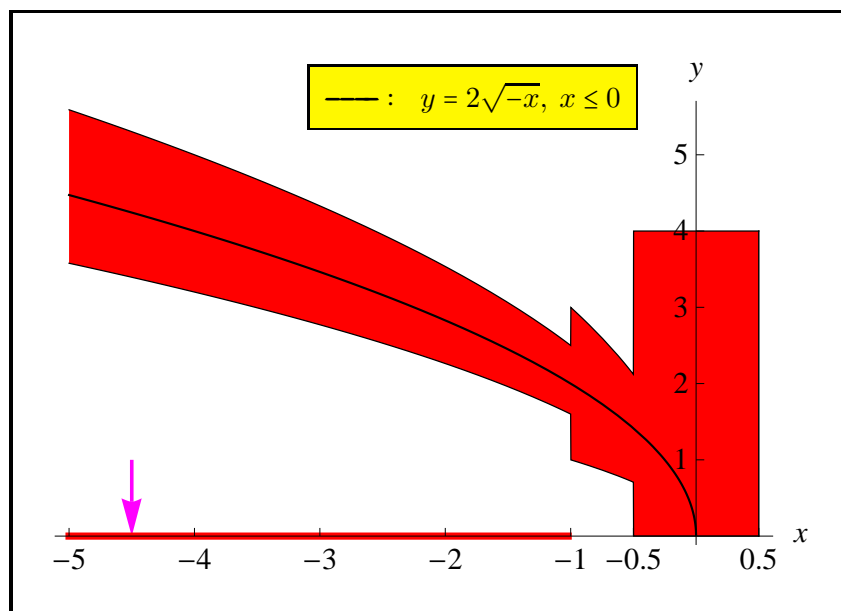


Abbildung C.14.: **Verzweigungsschnitt** und **Auslöschungsbereich** für den Realteil

Die **obere** und **untere** Umrandung des roten **Bereichs** ist mit  $p(x) := 2\sqrt{-x}$  gegeben durch:

$$g(x) := \begin{cases} 1.25 \cdot p(x), & -\infty \leq x < -1, \\ 1.5 \cdot p(x), & -1 \leq x < -0.5, \\ 4, & -0.5 \leq x \leq +0.5; \end{cases} \quad u(x) := \begin{cases} 0.8 \cdot p(x), & -\infty \leq x < -1, \\ 0.5 \cdot p(x), & -1 \leq x < -0.5, \\ 0, & -0.5 \leq x \leq +0.5; \end{cases}$$

Beide Funktionen sind definiert in der booleschen Funktion `bool Cancellation(x,y)`, die zu einem gegebenen Real- und Imaginärteil  $x, y$  entscheidet, ob bei der Realteilberechnung die Präzision entsprechend zu vergrößern ist, d.h. ob  $z = x + i \cdot y$  im roten Bereich liegt. In diesem Fall erfolgt für  $x < -1 \wedge y \neq 0$  die Berechnung des Realteils nach (C.41) mithilfe der Funktion

```
MpfrClass sqrtp1m1H2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd)
```

und für  $-1 \leq x \leq +0.5 \wedge y \neq 0$  mit

```
MpfrClass sqrtp1m1H1(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd).
```

Es ist zu beachten, dass man sich nach (C.41) bei der Berechnung des Realteils auf die obere Halbebene, d.h. auf  $y \geq 0$  beschränken kann.

Für  $|x| < 0.5 \wedge y < 0.5$ , d.h. im Rechteck der Breite 1 und der Höhe 0.5, das ganz im roten Bereich liegt, kann der Realteil trotz hinreichender Verdoppelung der Präzision beim Abrunden nicht optimal berechnet werden, wenn die Ausgangspräzision  $\text{pm}$  nicht groß genug gewählt wird. Wir betrachten dazu das Beispiel mit  $x = 0$ ,  $y = 2^{-1000}$  und wählen als Ausgangspräzision  $\text{pm} = 200$ . Nach (C.41) ist auszuwerten

$$s := |z + 1| + x + 1 = \sqrt{1 + y^2} + 1 = \sqrt{1 + 2^{-2000}} + 1.$$

Beim Abrunden erhält man daher auch bei doppelter Präzision  $2 \cdot \text{pm} = 400$  das Ergebnis  $s_d = 2$ , so dass  $\Re(\sqrt{z + 1} - 1)$  den viel zu stark abgerundeten Wert 0 erhält. Um einen positiven, abgerundeten Wert zu erhalten, müsste in diesem Beispiel die Ausgangspräzision mindestens den Wert  $\text{pm} = +1000$  erhalten, denn bei der ersten Präzisionsverdoppelung wäre die Current-Präzision dann  $\text{prec} = 2000$ . Der folgende Algorithmus liefert mit  $\text{ex}_x = \text{expo}(x)$  und  $\text{ex}_y = \text{expo}(y)$  den geeigneten Wert  $\text{pm}$ .

```

if ( (ex_x <= -1) && (ex_y <= -1) )
{
    if ((ex_x-1)/2 >= ex_y-1)
        pm = -(ex_x-1)/2 + 1;
    else pm = -(ex_y-1);
    if (pm < 2) pm = 2;
}
else pm = 2;

```

Ist die intern zu verwendende Präzision  $\text{prec}$  also kleiner als  $\text{pm}$ , so ist  $\text{prec} = \text{pm}$  zu setzen. Für das obige Beispiel erhält man dann mit  $\text{ex}_y = -999$  genau den gewünschten Wert  $\text{pm} = +1000$ .

Bei der Auswertung des Wurzelausdrucks in (C.41) besteht noch ein weiteres Problem, denn mit der Abkürzung  $M := \text{MaxFloat}()$  gilt die Abschätzung

$$A := \sqrt{2 \cdot (|z + 1| + x + 1)} < \sqrt{2 \cdot (\sqrt{2(M+1)^2 + M+1} + M+1)} = \sqrt{2(\sqrt{2} + 1) \cdot (M+1)} < 3 \cdot \sqrt{M+1} < M,$$

so dass  $A$  im Zahlenformat darstellbar ist, während  $2 \cdot (|z + 1| + x + 1)$  durchaus einen vorzeitigen Überlauf verursachen kann. Um diesen Überlauf zu vermeiden, wird bei zu großen Werten von  $|x|$  oder  $|y|$  skaliert, d.h.  $x, y, 1$  werden bei korrekter Rundung durch  $2^4$  dividiert, und die ausgewertete Wurzel wird dann am Ende zum Ausgleich wieder mit  $2^2 = 4$  rundungsfehlerfrei multipliziert, da wegen der obigen Abschätzung für  $A$  kein Überlauf eintreten kann. Wir zeigen noch, dass die genannte Skalierung den vorzeitigen Überlauf tatsächlich verhindert:

$$\begin{aligned}
2 \cdot \left( \sqrt{2^{-8} \cdot (x+1)^2 + 2^{-8} \cdot y^2} + 2^{-4} \cdot (x+1) \right) &= 2^{-3} \cdot \left( \sqrt{(x+1)^2 + y^2} + x+1 \right) \\
&\leq 2^{-3} \cdot \left( \sqrt{(x+1)^2 + (y+1)^2} + x+1 \right) \\
&\leq 2^{-3} \cdot (1 + \sqrt{2}) \cdot (M+1) \leq 0.3 \cdot (M+1) < M \blacksquare
\end{aligned}$$

Die Berechnung von  $\sqrt{2 \cdot (|z + 1| + x + 1)}$  erfolgt mithilfe der Funktion

```
MpfrClass Sqrt_zpx1 (const MpfrClass& x, const MpfrClass& y, const RoundingMode rnd)
```

die in der Datei `mpfclass.cpp` definiert ist.

Der vergleichsweise einfache komplexwertige Term  $\sqrt{z + 1} - 1$  zeigt, wie relativ kompliziert die Auswertung werden kann, wenn in der komplexen Ebene bezüglich der verlangten drei Rundungen in allen Bereichen eine fast optimale Genauigkeit erreicht werden soll.

### C.1.26.1. Numerische Ergebnisse

Im **1. Beispiel** wählen wir mit  $z = x + i \cdot y \in \mathbb{C}$ ,  $x = y = \text{MaxFloat}() = 2.0985 \dots 10^{+323228496}$ ,  $\text{prec} = 800$  und erhalten für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (1.59160566708640011463\dots630e161614248, 6.59264653257064147275\dots455e161614247); \\ f(z)_u &= (1.59160566708640011463\dots631e161614248, 6.59264653257064147275\dots457e161614247). \end{aligned}$$

Man erkennt, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen. Beachten Sie bitte, dass bei den internen Berechnungen kein vorzeitiger Überlauf auftritt.

Im **2. Beispiel** wählen wir  $x = 0.5$ ,  $y = 4.0$  und liegen damit nach Abbildung C.14 im oberen rechten Eckpunkt des roten Bereiches. Mit  $\text{prec} = 800$  erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (6.9882339762830638733342779\dots5497e - 1, 1.1772854098855480133550677\dots6233088); \\ f(z)_u &= (6.9882339762830638733342779\dots5498e - 1, 1.1772854098855480133550677\dots6233089). \end{aligned}$$

Man erkennt, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen.

Im **3. Beispiel** wählen wir  $x = 0.5$ ,  $y = 4.125$  und liegen damit nach Abbildung C.14 außerhalb des roten Bereiches. Mit  $\text{prec} = 800$  erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (7.1599278356872168635256989\dots0391e - 1, 1.2019281314870409954317952\dots1811408); \\ f(z)_u &= (7.1599278356872168635256989\dots0392e - 1, 1.2019281314870409954317952\dots1811409). \end{aligned}$$

Man erkennt, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen. Beachten Sie bitte, dass wir jetzt außerhalb des roten Bereichs praktisch die gleiche Genauigkeit erhalten wie innerhalb dieses Bereiches, in dem Auslöschung eintreten kann.

Im **4. Beispiel** wählen wir  $x = y = 2^{-5000000}$  und liegen damit in der rechten Halbebene sehr dicht am Ursprung innerhalb des roten Bereiches. Mit  $\text{prec} = 800$  erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  nach ca. zwei Minuten das Ergebnis:

$$\begin{aligned} f(z)_d &= (5.25593636371215653698\dots757e - 1505151, 5.25593636371215653698\dots757e - 1505151); \\ f(z)_u &= (5.25593636371215653698\dots758e - 1505151, 5.25593636371215653698\dots758e - 1505151). \end{aligned}$$

Man erkennt wieder, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen.

Im **5. Beispiel** wählen wir  $x = y = 0$  und liegen damit im Ursprung der komplexen Ebene. Mit  $\text{prec} = 800$  erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (0, 0); \\ f(z)_u &= (0, 0). \end{aligned}$$

Der exakte Funktionswert  $f(z) = 0$  wird also in beiden Rundungsmodi genau berechnet.

Im **6. Beispiel** wählen wir  $x = -2^{-5000000}$ ,  $y = 2 \cdot \sqrt{-x}$  und liegen damit nach Abbildung C.14 in der linken Halbebene genau auf der Parabel, auf der der Realteil von  $f(z) = \sqrt{z+1} - 1$  verschwindet. Mit `prec = 800` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z)$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (0, 1.0252742426992064582214716765280281147651738191817188703...012e - 752575); \\ f(z)_u &= (0, 1.0252742426992064582214716765280281147651738191817188703...013e - 752575). \end{aligned}$$

Man erkennt, dass der Realteil sogar exakt berechnet wird und dass der 241-stellige Imaginärteil fast optimal gerundet wurde, da die beiden gerundeten Werte auf den ersten 240 Dezimalstellen übereinstimmen.

Im **7. Beispiel** wählen wir  $x = -2^{-5000000}$ ,  $y = 2 \cdot \sqrt{-x}$  und mit  $y = \text{succ}(y)$  liegen wir dann nur ganz knapp oberhalb der Parabel  $y = 2 \cdot \sqrt{-x}$ , so dass der Realteil von  $f(z) = \sqrt{z+1} - 1$  nicht verschwinden wird. Mit `prec = 800` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z)$  nach etwa zwei Minuten das Ergebnis:

$$\begin{aligned} f(z)_d &= (3.15292440747892372470...656e - 1505391, 1.02527424269920645822...012e - 752575); \\ f(z)_u &= (3.15292440747892372470...657e - 1505391, 1.02527424269920645822...013e - 752575). \end{aligned}$$

Man erkennt wieder, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen. Wählt man  $x$  betragsmäßig deutlich kleiner, z.B.  $x = -2^{-50000000}$ , so wird bei gleicher Präzision die Laufzeit erheblich vergrößert.

Im **8. Beispiel** wählen wir  $x = 0$ ,  $y = 2^{-1000}$ . Mit `prec = 200` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (1.08872622702715208444...708043e - 603, 4.66631809251609439495...310856e - 302); \\ f(z)_u &= (1.08872622702715208444...708044e - 603, 4.66631809251609439495...310858e - 302). \end{aligned}$$

Man erkennt wieder, dass die obigen 60-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 59 Dezimalstellen übereinstimmen. Beachten Sie zu diesem Beispiel auch die entsprechenden Anmerkungen auf Seite 276.

Im **9. Beispiel** wählen wir  $x = -1.0625$ ,  $y = 2.5625$  und liegen damit nach Abbildung C.14 innerhalb des roten Bereiches. Mit `prec = 800` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (1.1820438243668915965034779...8317e - 1, 1.1458102115536487868975726...8878406); \\ f(z)_u &= (1.1820438243668915965034779...8318e - 1, 1.1458102115536487868975726...8878407). \end{aligned}$$

Man erkennt, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen.

Im **10. Beispiel** wählen wir  $x = -1.062500$ ,  $y = 2.57812500$  und liegen damit nach Abbildung C.14 außerhalb des roten Bereiches. Mit `prec = 800` erhalten wir für den **abgerundeten** und **aufgerundeten** Funktionswert  $f(z) = \sqrt{z+1} - 1$  das Ergebnis:

$$\begin{aligned} f(z)_d &= (1.2169123748123625205390267...3628e - 1, 1.1492133101570774719438668...4033545); \\ f(z)_u &= (1.2169123748123625205390267...3630e - 1, 1.1492133101570774719438668...4033546). \end{aligned}$$

Man erkennt, dass die obigen 241-stelligen Funktionswerte fast optimal gerundet wurden, da sie auf den ersten 240 Dezimalstellen übereinstimmen.

### C.1.27. $\sqrt{1+z^2}$

Wir betrachten die Aufgabe, für eine vorgegebene Präzision  $\text{prec} \geq 2$  nahezu optimal gerundete Näherungen der Funktion  $f(z) = \sqrt{1+z^2}$ ,  $z = x+i \cdot y \in \mathbb{C}$  zu berechnen. Die Verzweigungsschnitte von  $f(z)$  liegen auf der imaginären Achse von  $+i$  bis  $+i\infty$  und von  $-i$  bis  $-i\infty$ .

#### C.1.27.1. Realteil

Analog zu (C.16) auf Seite 256 erhält man für den Realteil von  $f(z) = \sqrt{1+z^2}$  den Ausdruck

$$(C.45) \quad \Re\{f(z)\} = u(x, y) := \sqrt{\left(\sqrt{(x^2 - y^2 + 1)^2 + 4x^2y^2} + (x^2 - y^2 + 1)\right) / 2}.$$

Aus (C.45) ergibt sich direkt, dass man sich auf  $x, y \geq 0$  beschränken kann. Speziell gilt:

$$(C.46) \quad u(0, y) = \begin{cases} 0, & x = 0 \wedge |y| \geq 1, \\ \sqrt{1 - y^2}, & x = 0 \wedge 0 \leq |y| < 1; \end{cases}$$

$$(C.47) \quad u(x, \pm 1) = \sqrt{|x|} \cdot \sqrt{\sqrt{1 + \left(\frac{x}{2}\right)^2} + \frac{|x|}{2}}, \quad \text{falls } |y| = 1,$$

wobei  $\sqrt{1-y^2}$  und  $\sqrt{1+(x/2)^2}$  mit den bereits in `mpfrclass.cpp` vordefinierten Funktionen `sqrt1mx2(...)` und `sqrt1px2(...)` optimal berechnet werden können.

Bei der Auswertung von  $u(x, y)$  spielt die Summe  $s := x^2 - y^2 + 1$  eine zentrale Rolle. Im Fall  $x \geq y \geq 0$  kann  $x^2 - y^2$  mithilfe von `x2my2(...)`, definiert in `mpfrclass.cpp`, ohne Auslöschung optimal berechnet werden, sonst kann in der Umgebung von  $y = \sqrt{1+x^2}$  bei der Auswertung von  $s = x^2 + (1+y)(1-y)$  Auslöschung auftreten, die jedoch bei einer Rechnung mit doppelter Präzision vermieden werden kann. Außerhalb dieser Umgebung wird  $s = x^2 + (1+y)(1-y)$  in einfacher Präzision ausgewertet.

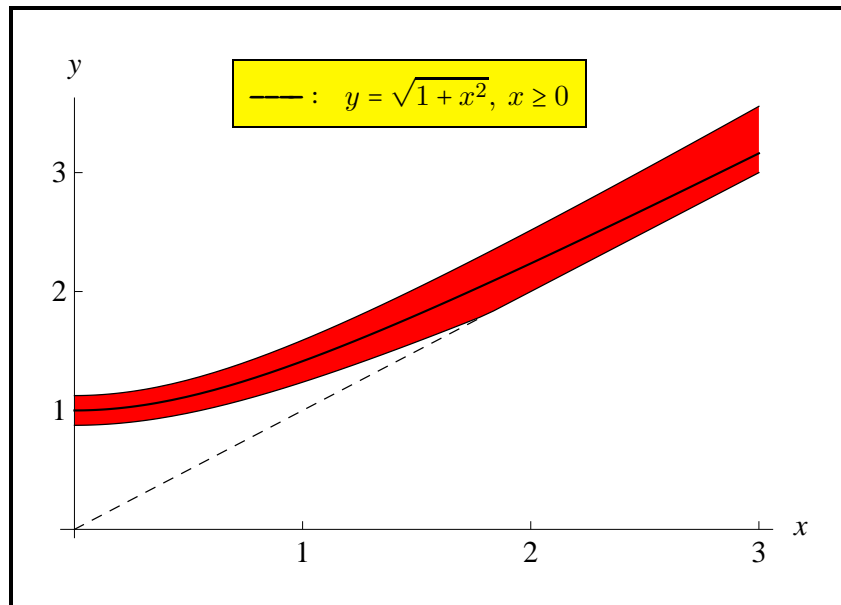


Abbildung C.15.: **Auslöschungsbereich** für  $s = x^2 + (1+y)(1-y)$

Der **Auslöschungsbereich** wird von oben begrenzt durch  $g_1(x) := 1.125 \cdot \sqrt{1+x^2}$  und von unten für  $0 \leq x \leq 1.8$  durch  $g_2(x) := 0.875 \cdot \sqrt{1+x^2}$  bzw. für  $x > 1.8$  durch  $g_3(x) = x$ . Die Berechnung der Funktionen  $g_1(x), g_2(x)$  ist nicht kritisch und erfolgt für  $x < 100$  mit Hilfe der schnellen C-XSC-Funktion `sqrt1px2(...)`. Für  $x \geq 100$  wird  $\sqrt{1+x^2} \approx x$  benutzt. Innerhalb des roten Bereiches ist der Graph von  $y = \sqrt{1+x^2}$  gezeichnet, auf dem  $s(x, y) = x^2 - y^2 + 1$  verschwindet.

Der Algorithmus zur Berechnung von  $s = x^2 - y^2 + 1$  ist gegeben durch:

```

if (x>=y)
  s = (x^2-y^2) + 1;
else
  if (y<=1.125*sqrt{1+x^2} && y>=0.875*sqrt{1+x^2})
    s = x^2 + (y-1)*(y+1);  (Doppelte Praezision)
  else
    s = x^2 + (y-1)*(y+1);  (Einfache Praezision)

```

Beachten Sie, dass es beim obigen Algorithmus für  $x \rightarrow \infty$  oder für  $y \rightarrow \infty$  bei der Auswertung von  $s = x^2 + (y-1) * (y+1)$  zu einem vorzeitigen Überlauf kommen kann. Dies ist jedoch kein wirklicher Nachteil, da im Fall  $|z| \rightarrow +\infty$  zur Berechnung von Real- und Imaginärteil die sehr einfache Näherung  $\sqrt{1+z^2} \approx z + 1/(2z)$  benutzt wird. Bis auf die genannten Einschränkungen  $x, y \rightarrow \infty$  erfolgt die nahezu optimale Berechnung von  $s$  mit Hilfe der Funktion

```
MpfrClass x2my2p1(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in `mpfcclass.cpp` definiert ist, aber nicht zur allgemeinen Anwendung zur Verfügung steht.

Wir kommen jetzt in (C.45) zur Berechnung des Ausdrucks

$$(C.48) \quad \alpha := \sqrt{(x^2 - y^2 + 1)^2 + 4x^2y^2} + (x^2 - y^2 + 1) = \sqrt{s^2 + (2xy)^2} + s,$$

der auch bei der Auswertung des Imaginärteils in (C.57) auf Seite 282 eine zentrale Rolle spielt. Die Wurzel  $\sqrt{s^2 + (2xy)^2}$  wird mit der in `pmfrclass.cpp` definierten Funktion `sqrtx2y2(...)` berechnet, und Auslöschung kann in (C.48) nur für  $s \leq 0$ , d.h. für  $y \geq \sqrt{1+x^2} \geq 1$  auftreten, und diese Auslöschung wird vermieden, wenn  $\alpha$  wie folgt umgeformt wird:

$$(C.49) \quad \sqrt{s^2 + (2xy)^2} + s = \frac{(\sqrt{s^2 + (2xy)^2} + s)(\sqrt{s^2 + (2xy)^2} - s)}{\sqrt{s^2 + (2xy)^2} - s} = \frac{(2xy)^2}{\sqrt{s^2 + (2xy)^2} - s},$$

wobei jetzt der letzte Nenner als Summe positiver Summanden problemlos ausgewertet werden kann. Wir fassen zusammen:

$$(C.50) \quad \alpha = \begin{cases} \sqrt{s^2 + (2xy)^2} + s, & s > 0 \\ \frac{((2x) \cdot y)^2}{\sqrt{s^2 + ((2x) \cdot y)^2} - s}, & s \leq 0, \quad \text{wobei } y \geq 1 \text{ automatisch erfüllt ist.} \end{cases}$$

Wir betrachten den Fall  $s \leq 0$ :

Wenn  $\alpha$  aufzurunden ist, so ist im Zähler  $(2x) \cdot y$  aufzurunden und unter der Wurzel ist  $(2x) \cdot y$  abzurunden, so dass es sinnvoll ist,  $(2x) \cdot y$  **intervallmäßig** auszuwerten; dies gilt auch dann, wenn  $\alpha$  abzurunden ist. Wenn  $\alpha$  aufzurunden ist, so muss neben der Wurzel im Nenner selbst wegen  $s \leq 0$  auch die nicht-negative Summe  $-s \geq 0$  abgerundet werden, so dass  $s \leq 0$  selbst aufzurunden ist. Wenn andererseits  $\alpha$  abzurunden ist, so muss neben der Wurzel im Nenner selbst wegen  $s \leq 0$  auch die nicht-negative Summe  $-s \geq 0$  aufgerundet werden, so dass  $s \leq 0$  selbst abzurunden ist.

Wir betrachten den Fall  $s > 0$ :

Um  $\alpha = \sqrt{s^2 + (2xy)^2} + s$  z.B. aufzurunden, müssen neben den Grundoperationen auch  $s$  und  $2xy$  selbst aufgerundet werden und für das Abrunden gelten ganz analoge Aussagen.

Die Auswertung von  $\alpha$  nach (C.50) erfolgt mithilfe der Funktion

```
MpfrClass sqrt1px2_alpha(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd)
```

die in der Datei `mpfcclass.cpp` definiert und nur für den internen Gebrauch bestimmt ist.



Im folgenden Abschnitt muss noch geklärt werden, wie bei der Auswertung von  $\alpha$  nach (C.50) ein vorzeitiger Überlauf verhindert werden kann. Gesucht ist dabei z.B. eine Schranke  $c > 0$ , so dass für  $0 \leq x < c$  und  $0 \leq y < c$  bei der Auswertung des Zählers  $(2xy)^2$  kein Überlauf entsteht, d.h. wir verlangen  $(2xy)^2 < M_0 := 1.573934 \cdot 10^{+323228496}$ , wobei  $M_0$  eine Unterschranke von  $\text{MaxFloat}(2) \leq \text{MaxFloat}(\text{prec})$  ist. Um  $\alpha$  nach (C.50) ohne vorzeitigen Überlauf berechnen zu können, müssen folgende Forderungen erfüllt sein:

1.  $4x^2y^2 < M_0 \quad (\implies \quad 2xy < M_0)$ ,
2.  $|s| = |x^2 - y^2 + 1| < M_0$ ,
3.  $\sqrt{s^2 + (2xy)^2} + |s| < M_0$ ,
4.  $\frac{4x^2y^2}{\sqrt{s^2 + (2xy)^2} + |s|} < M_0$ .

Wegen  $\sqrt{s^2 + (2xy)^2} + |s| \geq 2xy$  ist die letzte Forderung erfüllt, wenn gilt  $2xy < M_0$ , und diese Ungleichung ist erfüllt, wenn die Forderung 1 erfüllt ist. Danach sind nur noch die **drei ersten** Forderungen zu beachten.

Um die erste Forderung zu erfüllen, gilt mit  $0 \leq x < c_1$  und  $0 \leq y < c_1$  die Abschätzung  $4x^2y^2 < 4 \cdot c_1^4$ , und die erste Forderung ist erfüllt für  $c_1 = \sqrt[4]{M_0/4} = 7.9201176703\dots \cdot 10^{+80807123}$ .

Mit  $0 \leq x < c_2$  und  $0 \leq y < c_2$  ist wegen  $|x^2 - y^2 + 1| \leq x^2 + y^2 + 1 < 2c_2^2 + 1$  die zweite Forderung erfüllt, wenn gilt:  $c_2 = \sqrt{(M_0 - 1)/2}$ , und wegen  $c := c_1 < c_2$  sind die beiden ersten Forderungen für  $0 \leq x, y < c = \sqrt[4]{M_0/4}$  erfüllt, so dass jetzt unter der Voraussetzung  $0 \leq x, y < c$  nur noch die 3. Bedingung zu untersuchen ist. Mit der recht groben Abschätzung

$$\begin{aligned} \sqrt{s^2 + (2xy)^2} + |s| &\leq \sqrt{(2c^2 + 1)^2 + 4c^4} + (2c^2 + 1) < \sqrt{(3c^2)^2 + 4c^4} + 3c^2 \\ &= \sqrt{9c^4 + 4c^4} + 3c^2 < 7c^2 \end{aligned}$$

ist dann die 3. Bedingung erfüllt, wenn gilt:  $7c^2 < M_0 \iff 7 \cdot \sqrt{M_0/4} < M_0$ , wobei die letzte Ungleichung wegen  $M_0 \gg 1$  offensichtlich richtig ist. ■

#### Anmerkung:

Bei der obigen Berechnung von  $c$  wurde nicht berücksichtigt, dass der Maschinenwert von  $4x^2y^2$  bei minimaler Präzision  $\text{prec} = 2$  den exakten Wert etwa um den Faktor 2 übertreffen kann. Wenn man diesen Effekt bei der Berechnung von  $c$  aus Sicherheitsgründen mit dem Faktor 8 berücksichtigt, so erhält man für  $c$  den etwas kleineren Wert:  $c = 4.7093\dots \cdot 10^{+80807123}$ . Im Programm benutzen wir für die positiven Werte  $x, y \geq 0$  die noch etwas kleinere Oberschranke  $c_0 := 2^{268435454} < c = \sqrt[4]{M_0/(4 \cdot 8)} = 4.70933014\dots \cdot 10^{+80807123}$ .

#### Zusammenfassung:

Unter der Voraussetzung  $0 \leq x, y < c_0 = 0.5 \cdot 2^{268435455}$  wird  $\alpha$  nach (C.50) auf der Maschine für alle  $\text{prec} \geq 2$  ohne einen vorzeitigen Überlauf berechnet.

Es muss jetzt noch geklärt werden, wie  $\Re(f(z)) = u(x, y)$  im Fall  $x > c_0 \vee y > c_0$  möglichst optimal berechnet werden kann. Mit  $r := 1/z$  und

$$(C.51) \quad R := \sum_{k=3}^{\infty} \frac{1 \cdot 3 \cdot 5 \cdots (2k-3)(-1)^{k+1}}{2 \cdot 4 \cdot 6 \cdots (2k-2) \cdot (2k)} \cdot r^{2k-1}, \quad |r| = \frac{1}{|z|} < 1 \quad \text{erhält man}$$

$$(C.52) \quad f(z) = \sqrt{1+z^2} = \begin{cases} z + \frac{1}{2z} - \frac{1}{8z^3} + R, & x > 0, \\ -\left(z + \frac{1}{2z} - \frac{1}{8z^3} + R\right), & x < 0. \end{cases}$$

$$(C.53) \quad \Re(f(z)) = 0, \quad \text{falls } x = 0 \wedge |y| \geq 1.$$

Um  $f(z)$  für  $|z| \gg 1$  optimal einschließen zu können, benötigen wir zunächst eine Abschätzung für  $|R|$ . Mithilfe der geometrischen Reihe findet man wegen  $|r| < c_0^{-1}$  unmittelbar:

$$(C.54) \quad |R| \leq \frac{|r|^5}{16} \cdot (1 + |r| + |r|^2 + \dots) = \frac{|r|^5}{16} \cdot \frac{1}{1 - |r|} < \frac{|r|^5}{8} < \frac{c_0^{-5}}{8} < \Delta := \text{minfloat}().$$

Um für  $x > c_0 \vee y > c_0$ , d.h. für  $|z| > c_0$  möglichst optimal auf- bzw. abgerundete Werte von  $f(z)$  berechnen zu können, wird zunächst der Ausdruck  $t := 1/(2z) - 1/(8z^3) \in \mathbb{C}$  intervallmäßig möglichst optimal durch  $T \ni t$  eingeschlossen und eine Einschließung für  $f(z)$  erhält man dann durch

$$(C.55) \quad f(z) \in z \diamond (T \diamond ([-\Delta, +\Delta], [-\Delta, +\Delta])), \quad \Delta = \text{minfloat}(),$$

wobei  $([-\Delta, +\Delta], [-\Delta, +\Delta])$  als komplexes Intervall mit jeweils gleichem Real- und Imaginärteilintervall  $[-\Delta, +\Delta]$  aufzufassen ist.  $t = 1/(2z) - 1/(8z^3)$  wird optimal eingeschlossen durch

$$(C.56) \quad t \in T = \{(2 \diamond (1 \diamond z)) \diamond (2 \diamond (1 \diamond z))\} \diamond z \diamond 8,$$

wobei  $(1 \diamond z)$  aus Laufzeitgründen durch  $\text{reci}(z)$  realisiert wird. Die Berechnung der Realteilfunktion  $u(x, y) = \Re\{\sqrt{1 + z^2}\}$  erfolgt mit Hilfe der Funktion

```
MpfrClass Re_sqrt1px2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in der Datei `mpfcclass.cpp` definiert ist.

### C.1.27.2. Imaginärteil

Analog zu (C.16) auf Seite 256 erhält man für den Imaginärteil von  $f(z) = \sqrt{1 + z^2}$

$$(C.57) \quad \Im\{f(z)\} = v(x, y) := \begin{cases} \frac{\text{sign}(x) \cdot \text{sign}(y) \cdot |x| \cdot |y|}{\sqrt{(\sqrt{(x^2 - y^2 + 1)^2 + 4x^2y^2} + (x^2 - y^2 + 1))} / 2}, & x \neq 0, \\ \sqrt{y^2 - 1}, & x = 0 \wedge |y| \geq 1, \\ 0, & x = 0 \wedge |y| < 1 \end{cases}$$

Speziell gilt für  $y = \pm 1$

$$v(x, \pm 1) = \frac{\text{sign}(x) \cdot \text{sign}(y) \cdot \sqrt{|x|}}{\sqrt{\sqrt{1 + \left(\frac{x}{2}\right)^2} + \frac{|x|}{2}}}, \quad \text{falls } |y| = 1,$$

wobei die Wurzel im Nenner mithilfe der Funktion `sqrt1px2(...)` ausgewertet wird. Die Berechnung der Imaginärteilfunktion  $v(x, y) = \Im\{\sqrt{1 + z^2}\}$  erfolgt mit Hilfe der Funktion

```
MpfrClass Im_sqrt1px2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

die in der Datei `mpfcclass.cpp` definiert ist.

### C.1.28. $\sqrt{1-z^2}$

Wir betrachten die Aufgabe, für eine vorgegebene Präzision  $\text{prec} \geq 2$  nahezu optimal gerundete Näherungen der Funktion  $f(z) = \sqrt{1-z^2}$ ,  $z = x + i \cdot y \in \mathbb{C}$  zu berechnen. Die Verzweigungsschnitte von  $f(z)$  liegen auf der reellen Achse von  $+1$  bis  $+\infty$  und von  $-1$  bis  $-\infty$ . Auf dem rechten Verzweigungsschnitt werden die Funktionswerte von  $f(z)$  von unten und auf dem linken Verzweigungsschnitt von obenher analytisch fortgesetzt. In der folgenden Abbildung wird dies durch die beiden Pfeile gekennzeichnet.

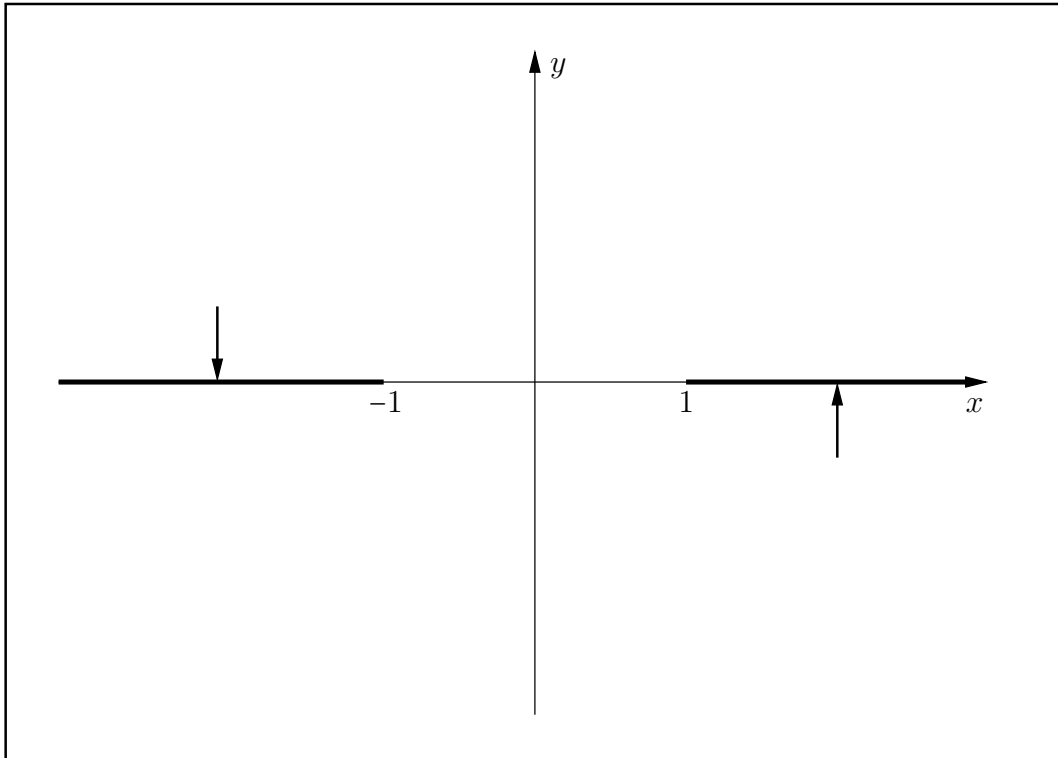


Abbildung C.16.: Verzweigungsschnitte der Funktion  $f(z) = \sqrt{1-z^2}$

Die gerundeten Funktionswerte werden berechnet mithilfe der Funktion

```
MpfcClass sqrt1mx2(const MpfcClass& z, RoundingMode rnd);
```

die in der Datei `mpfcclass.cpp` definiert ist. Mit dem Rundungsparameter `rnd` können dabei die Werte `RoundDown`, `RoundUp`, `RoundNearest` gewählt werden.

Mit  $z = x + i \cdot y \in \mathbb{C}$  und der Transformation

$$z = i \cdot z_T \quad \text{bzw.} \quad z_T = y + i \cdot (-x) \quad \text{erhält man} \quad \sqrt{1-z^2} = \sqrt{1+z_T^2},$$

so dass die Auswertung von  $f(z) = \sqrt{1-z^2}$  zurückgeführt werden kann auf die Auswertung von  $\sqrt{1+z_T^2}$ , wobei diese Funktion bereits durch `sqrt1px2(..)` in der Datei `mpfcclass.cpp` definiert wurde.

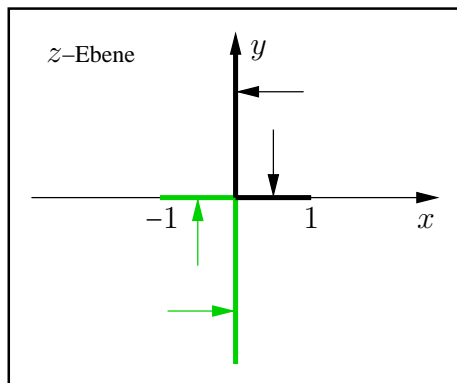
### C.1.29. $\sqrt{z^2 - 1}$

Wir betrachten die Aufgabe, für eine vorgegebene Präzision  $\text{prec} \geq 2$  nahezu optimal gerundete Näherungen der Funktion  $f(z) = \sqrt{z^2 - 1}$ ,  $z = x + i \cdot y \in \mathbb{C}$  zu berechnen. Die Verzweigungspunkte ergeben sich aus der Bedingung  $z^2 - 1 = 0$  zu  $z_{1,2} = \pm 1$ , und die Verzweigungsschnitte ergeben sich aus der Bedingung  $z^2 - 1 = -\alpha$ , mit  $\alpha \geq 0$ , bzw. aus dem Gleichungssystem

$$\begin{aligned} x^2 - y^2 - 1 &= -\alpha, & \alpha &\geq 0, \\ x \cdot y &= 0, & &\text{mit den zwei } (\pm) \text{ Lösungen} \end{aligned}$$

$$x = 0, \quad y = \pm\sqrt{\alpha - 1}, \quad \text{falls } \alpha \geq 1,$$

$$y = 0, \quad x = \pm\sqrt{1 - \alpha}, \quad \text{falls } 0 \leq \alpha \leq 1.$$



Im linken Bild sind die zwei Verzweigungsschnitte durch die Farben schwarz (+) und grün (-) gekennzeichnet. Die gleichfarbigen Pfeile geben die Richtungen an, aus denen die Funktionswerte von  $f(z) = \sqrt{z^2 - 1}$  auf den jeweiligen Verzweigungsschnitt analytisch fortgesetzt werden.

Zu beachten ist, dass die Funktionswerte der Realteilfunktion  $\hat{u}(x, y)$  auf beiden Verzweigungsschnitten nach (C.59) durch den Wert 0 stetig ergänzt werden können.

Bildet man jedoch auf beiden Verzweigungsschnitten die Grenzwerte der Imaginärteilfunktion  $\hat{v}(x, y)$ , und zwar jeweils von beiden Seiten her, so erhält man mit Ausnahme des Ursprungs stets **verschiedene** Grenzwerte.

Zur Berechnung der Real- und Imaginärteilfunktionen von  $f(z) = \hat{u}(x, y) + i \cdot \hat{v}(x, y)$  betrachten wir die formale Umformung

$$\sqrt{z^2 - 1} = \sqrt{(-1) \cdot (1 - z^2)} = i \cdot \sqrt{1 - z^2}.$$

Mit der Transformation  $z = i \cdot z_T$  und  $z^2 = -z_T^2$  bzw.  $z_T = y - i \cdot x$  erhält man

$$\begin{aligned} \sqrt{z^2 - 1} &= i \cdot \sqrt{1 + z_T^2} = i \cdot (u(y, -x) + i \cdot v(y, -x)) = -v(y, -x) + i \cdot u(y, -x), \quad \text{d.h.} \\ \hat{u}(x, y) &= -v(y, -x), & \hat{v}(x, y) &= u(y, -x). \end{aligned} \tag{C.58}$$

Mit den Ergebnissen von Seite 279 bzw. Seite 282 für  $u(x, y)$  bzw.  $v(x, y)$  erhält man schließlich mit

$$h(x, y) := \sqrt{\left(\sqrt{(y^2 - x^2 + 1)^2 + 4x^2y^2} + (y^2 - x^2 + 1)\right) / 2} \geq 0$$

für  $f(z) = \hat{u}(x, y) + i \cdot \hat{v}(x, y)$  die Beziehungen

$$\hat{u}(x, y) = \begin{cases} \frac{|x| \cdot |y|}{h(x, y)}, & y \neq 0, \\ \sqrt{x^2 - 1}, & y = 0 \wedge |x| \geq 1, \\ 0, & y = 0 \wedge |x| < 1. \end{cases} \tag{C.59}$$

$$\hat{v}(x, y) = \begin{cases} \text{sign}(x) \cdot \text{sign}(y) \cdot h(x, y), & x \cdot y \neq 0, \\ 0, & y = 0 \wedge |x| \geq 1, \\ \sqrt{1 - x^2}, & y = 0 \wedge |x| < 1, \\ \sqrt{1 + y^2}, & x = 0. \end{cases} \tag{C.60}$$

In C.59, C.60 sind die Vorzeichen so gewählt, dass  $f(0) = \sqrt{-1} = +i$  und  $f(x+i \cdot 0) = \sqrt{x^2-1}$  für  $|x| \geq 1$  gewährleistet sind. Nach C.59, C.60 können damit die Real- und Imaginärteilmfunktionen  $\hat{u}(x, y)$  und  $\hat{v}(x, y)$  mithilfe der entsprechenden Funktionen  $v(x, y)$  und  $u(x, y)$  von  $\sqrt{1+z^2} = u(x, y) + i \cdot v(x, y)$  berechnet werden, die bereits in der Datei `mpfcClass.cpp` definiert worden sind.

Die Funktionswerte  $f(z) = \sqrt{z^2-1}$  werden berechnet mithilfe der Funktion

```
MpfcClass sqrtx2m1(const MpfcClass& z, RoundingMode rnd)
```

wobei für `rnd` nur die Rundungsmodi `RoundDown`, `RoundUp`, `RoundNearest` zugelassen sind. Wird `rnd` nicht gesetzt, so kommt der aktuelle Rundungsmodus zur Anwendung.

### C.1.29.1. Realteil

Die Realteilmfunktion  $\hat{u}(x, y)$  wird nach (C.59) berechnet mithilfe der Funktion

```
MpfrClass Re_sqrtx2m1(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd)
```

wobei für `rnd` nur die Rundungsmodi `RoundDown`, `RoundUp`, `RoundNearest` zugelassen sind. Zu beachten ist, dass nach (C.58) auf die bereits implementierte Imaginärteilmfunktion  $v(x, y)$  von  $\sqrt{1+z^2}$  zurückgegriffen werden kann.

### C.1.29.2. Imaginärteil

Die Imaginärteilmfunktion  $\hat{v}(x, y)$  wird nach (C.60) berechnet mithilfe der Funktion

```
MpfrClass Im_sqrtx2m1(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd)
```

wobei für `rnd` nur die Rundungsmodi `RoundDown`, `RoundUp`, `RoundNearest` zugelassen sind. Zu beachten ist dabei, dass nach (C.58) auf die bereits implementierte Realteilmfunktion  $u(x, y)$  von  $\sqrt{1+z^2}$  zurückgegriffen werden kann.

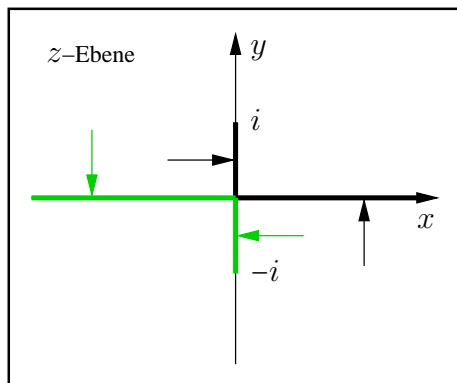
### C.1.30. $\sqrt{-z^2 - 1}$

Wir betrachten die Aufgabe, für eine vorgegebene Präzision  $\text{prec} \geq 2$  nahezu optimal gerundete Näherungen der Funktion  $f(z) = \sqrt{-z^2 - 1}$ ,  $z = x + i \cdot y \in \mathbb{C}$  zu berechnen. Die Verzweigungspunkte ergeben sich aus der Bedingung  $z^2 + 1 = 0$  zu  $z_{1,2} = \pm i$ , und die Verzweigungsschnitte ergeben sich aus der Bedingung  $-z^2 - 1 = -\alpha$ , mit  $\alpha \geq 0$ , bzw. aus dem Gleichungssystem

$$\begin{aligned} x^2 - y^2 + 1 &= \alpha, & \alpha \geq 0, \\ x \cdot y &= 0, & \text{mit den zwei } (\pm) \text{ Lösungen} \end{aligned}$$

$$x = 0, y = \pm\sqrt{1 - \alpha}, \quad \text{falls } 0 \leq \alpha \leq 1,$$

$$y = 0, x = \pm\sqrt{\alpha - 1}, \quad \text{falls } \alpha \geq 1.$$



Im linken Bild sind die zwei Verzweigungsschnitte durch die Farben schwarz (+) und grün (-) gekennzeichnet. Die gleichfarbigen Pfeile geben die Richtungen an, aus denen die Funktionswerte von  $f(z) = \sqrt{-z^2 - 1}$  auf den jeweiligen Verzweigungsschnitt analytisch fortgesetzt werden.

Zur Berechnung von  $f(z) = \sqrt{-z^2 - 1}$  betrachten wir die Transformation  $z = x + i \cdot y = i \cdot z_T$  und erhalten mit  $z_T = y - i \cdot x$  die einfache Beziehung

$$f(z) = \sqrt{-z^2 - 1} = \sqrt{z_T^2 - 1},$$

so dass  $f(z)$  mit Hilfe von  $z_T$  und der bereits implementierten Funktion `sqrtn2m1( $z_T$ )` direkt berechnet werden kann, vgl. Seite 284.

Die Auswertung von  $f(z) = \sqrt{-z^2 - 1}$  erfolgt mithilfe der Funktion

```
MpfcClass sqrtmx2m1(const MpfcClass& z, RoundingMode rnd)
```

die in `mpfcclass.cpp` definiert ist. Für `rnd` stehen die drei Rundungsmodi

`RoundDown`, `RoundUp`, `RoundNearest`

zur Verfügung.

## C.2. Elementarfunktionen für komplexe Intervallargumente

### C.2.1. $z^2$

Für  $z = x + i \cdot y \in \mathbb{C}$  ist  $z^2$  definiert durch

$$z^2 := x^2 - y^2 + 2i \cdot x \cdot y \equiv |x|^2 - |y|^2 + 2i \cdot x \cdot y.$$

Ist  $z$  jetzt ein Element aus einem komplexen Intervall  $Z = X + iY$  mit reellen Intervallen  $X, Y$ , so ist der Realteil von  $Z^2$  wie folgt zu berechnen:

$$\Re(Z^2) := \{(|x|^2 - |y|^2) \mid |x| \in |X| \wedge |y| \in |Y|\} = [\text{Inf}(|X|)^2 - \text{Sup}(|Y|)^2, \text{Sup}(|X|)^2 - \text{Inf}(|Y|)^2],$$

wobei z.B.  $\text{Inf}(|X|)^2 - \text{Sup}(|Y|)^2$  mithilfe der vordefinierten Funktion

```
MpfrClass x2my2 (const MpfrClass& x, const MpfrClass& y, RoundingMode rnd);
```

mit `rnd = RoundDown` auszuwerten ist.

Der Imaginärteil von  $Z^2$  ist wie folgt zu berechnen:

$$\Im(Z^2) := \{2 \cdot x \cdot y \mid x \in X \wedge y \in Y\} = 2 \cdot (X \cdot Y).$$

Um bei der Auswertung von  $(X \cdot Y)$  einen vorzeitigen Unterlauf zu vermeiden, berechnet man zuerst  $2 \cdot X$  bzw.  $2 \cdot Y$ , falls  $2 \cdot X$  zum Überlauf führt, und multipliziert dann mit dem Intervall  $Y$  bzw. mit  $X$ . Weitere Einzelheiten findet man in der Quelltext-Datei `mpfciclass.cpp`.

### C.2.2. $1/z$

Ist  $z = x + i \cdot y \in \mathbb{C}$  ein Element aus einem komplexen Intervall  $Z = X + iY$  mit reellen Intervallen  $X, Y$ , so berechnet die Funktion

```
Mpfciclass reci(const Mpfciclass& Z);
```

aus `mpfciclass.cpp` mit dem Aufruf `res = reci(Z);` eine optimale Einschließung der komplexen Menge

$$\{1/z \mid x \in X \wedge y \in Y\} \subseteq \text{res}, \quad 0 \notin Z.$$

#### Anmerkungen:

- 1  $\diamond$  Z und `reci(Z)` liefern bei gleicher Präzision `prec  $\geq$  2` **identische** komplexe Ergebnisintervalle.
- Im Vergleich zu 1  $\diamond$  Z liefert `reci(Z)` eine um den Faktor 13 geringere Laufzeit, so dass `reci(Z)` bevorzugt eingesetzt werden sollte!
- Ist W ein weiteres komplexes Intervall, so gilt i.a.  $W \diamond Z \subset W \diamond \text{reci}(Z)$ .  
Sind W und Z voneinander unabhängige komplexe Intervalle, so liefert die komplexe Intervalldivision  $W \diamond Z$  eine optimale Einschließung des Wertebereichs  $W := \{w \cdot z \mid w \in W \wedge z \in Z\}$ , und dies gilt in der komplexen Intervallarithmetik für alle vier Grundoperationen.

### C.2.3. $1/z^2$

Ist  $z = x + i \cdot y \in \mathbb{C}$  ein Element aus einem komplexen Intervall  $Z = X + iY$  mit reellen Intervallen  $X, Y$ , so berechnet die Funktion

```
Mpfciclass reci_z2(const Mpfciclass& Z);
```

aus `mpfciclass.cpp` mit dem Aufruf `res = reci_z2(Z);` eine optimale Einschließung der komplexen Menge

$$W_Z := \{1/z^2 \mid x \in X \wedge y \in Y\} \subseteq \text{res}, \quad 0 \notin Z.$$
$$1/z^2 = u(x, y) + i \cdot v(x, y); \quad u(x, y) := \frac{x^2 - y^2}{(x^2 + y^2)^2}, \quad v(x, y) := \frac{-2xy}{(x^2 + y^2)^2}.$$

Die optimale Einschließung der Menge  $\{u(x, y) \in \mathbb{R} \mid x \in X \wedge y \in Y\}$  mit Hilfe der Funktion

```
MpfiClass Re_rz2(const MpfiClass& X, const MpfiClass& Y)
```

wird ausführlich beschrieben ab Seite 210, und für den negativen Imaginärteil wird die optimale Einschließung der Menge  $\{-v(x, y) \in \mathbb{R} \mid x \in X \wedge y \in Y\}$  mit Hilfe der Funktion

```
MpfiClass mIm_rz2(const MpfiClass& X, const MpfiClass& Y)
```

ausführlich beschrieben ab Seite 216. Beachten Sie dabei, dass bei der Auswertung der Funktionen  $u(x, y), v(x, y)$  **kein** vorzeitiger Überlauf eintreten kann.

#### C.2.3.1. Numerische Ergebnisse

Im **1. Beispiel** wählen wir  $X = [2, 5]$  und  $Y = [1, 3]$  und berechnen mit dem Funktionsaufruf `res = reci_z2(Z)` und `prec = 53` eine optimale Einschließung der Menge  $W_Z$  durch

```
([-2.95857988165681e-2, 1.20000000000001e-1], [-1.62379763209583e-1, -1.47928994082840e-2])
```

Im **2. Beispiel** und **3. Beispiel** wählen wir für  $Z$  die gleichen reellen Intervalle  $X = [2, 5]$  und  $Y = [1, 3]$  und berechnen  $1/Z^2$  im **2. Beispiel** durch `sqr(1  $\diamond$  Z)` und im **3. Beispiel** durch `1  $\diamond$  (sqr(Z))`. Im folgenden werden die Einschließungen verglichen mit dem Ergebnis aus dem **1. Beispiel**:

```
([-2.95857988165681e-2, 1.20000000000001e-1], [-1.62379763209583e-1, -1.47928994082840e-2])  
([-4.08737024221454e-2, 1.58520710059172e-1], [-2.00000000000001e-1, -1.13122171945701e-2])  
([-1.25000000000000e-1, 1.25000000000000e-1], [-2.50000000000000e-1, -6.75675675675675e-3]).
```

Man erkennt, dass die Einschließungen des **2.** und **3. Beispiels** im Vergleich zum **1. Beispiel** deutlich schlechter ausfallen.

#### Zusammenfassung

Wertet man beliebige, reelle, äquivalente, arithmetische Intervallausdrücke, wie z.B.  $1 \diamond (\text{sqr}(X))$  oder  $\text{sqr}(1 \diamond X)$ , welche die Variable  $X$  nur einmal enthalten, naiv aus, so erhält man stets die gleiche, optimale Einschließung der Intervallausdrücke. Im Gegensatz dazu liefern entsprechende äquivalente, komplex Intervallausdrücke mit nur einer Variablen  $Z$  i.a. jeweils verschiedene und meist nicht-optimale Einschließungen. Man sollte daher auf eine möglichst große Anzahl von komplexen Funktionen zurückgreifen können, die optimale Einschließungen liefern, vgl. dazu die Tabelle auf Seite 115 und Abschnitt 6.12.4 auf Seite 119.



### C.2.4. $1/(1+z^2)$

Ist  $z = x + i \cdot y \in \mathbb{C}$  ein Element aus einem komplexen Intervall  $Z = X + iY$  mit reellen Intervallen  $X, Y$ , so berechnet die Funktion

```
Mpfciclass reci_1pz2(const Mpfciclass& Z);
```

aus `mpfciclass.cpp` mit dem Aufruf `res = reci_1pz2(Z)`; eine optimale Einschließung der komplexen Menge

$$W_Z := \{1/(1+z^2) \mid x \in X \wedge y \in Y\} \subseteq \text{res}, \quad (0 \pm i) \notin Z.$$

$$1/(1+z^2) = u + i \cdot v; \quad u(x, y) := \frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}, \quad v(x, y) := \frac{-2xy}{4x^2y^2+(1+x^2-y^2)^2}.$$

Die optimale Einschließung der Menge  $\{u(x, y) \in \mathbb{R} \mid x \in X \wedge y \in Y\}$  mit Hilfe der Funktion

```
MpfiClass Re_r1pz2(const MpfiClass& X, const MpfiClass& Y)
```

wird ausführlich beschrieben ab Seite 227. Für den Imaginärteil wird die optimale Einschließung der Menge  $\{v(x, y) \in \mathbb{R} \mid x \in X \wedge y \in Y\}$  mit Hilfe der mit  $(-1)$  zu multiplizierenden Funktion

```
MpfiClass Im_r1pz2(const MpfiClass& X, const MpfiClass& Y)
```

ausführlich beschrieben ab Seite 233. Beachten Sie dabei, dass bei der Auswertung der beiden Funktionen  $u(x, y), v(x, y)$  **kein** vorzeitiger Überlauf eintreten kann.

#### C.2.4.1. Numerische Ergebnisse

Im **1. Beispiel** wählen wir  $X = [1/8, 1/4]$  und  $Y = [1/8, 2]$  und berechnen mit dem Funktionsaufruf `res = reci_1pz2(Z)` und `prec = 53` eine optimale Einschließung der Menge  $W_Z$  durch

$$([-1.76556443707464, 2.26556443707464], [-3.98443603238417, -3.12195121951219e-2])$$

Im **2. Beispiel** wählen wir für  $Z$  die gleichen reellen Intervalle  $X = [1/8, 1/4]$  und  $Y = [1/8, 2]$  und berechnen eine Einschließung von  $W_Z$  durch die naive Intervallauswertung  $1 \diamond (1 \diamond \text{sqr}(Z))$ . Im folgenden wird diese Einschließungen verglichen mit dem Ergebnis aus dem **1. Beispiel**:

$$([-1.76556443707464, 2.26556443707464], [-3.98443603238417, -3.12195121951219e-2]) \\ ([-1.60000000000000e1, 1.60000000000000e1], [-3.20000000000000e1, -3.50829107852542e-3])$$

Man erkennt, dass die Einschließung von  $W_Z$  im **2. Beispiels** im Vergleich zum **1. Beispiel** deutlich schlechter ausfällt, d.h. dass die naive Intervallauswertung Überschätzungen bis etwa zum Faktor 10 verursachen kann.

Im **3. Beispiel** wählen wir für  $Z$  die gleichen reellen Intervalle  $X = [1/8, 1/4]$  und  $Y = [1/8, 2]$  und berechnen eine Einschließung von  $W_Z$  durch `reci_z2(sqrt1px2(Z))`. Im folgenden wird diese Einschließungen verglichen mit dem Ergebnis aus dem **1. Beispiel**:

$$([-1.76556443707464, 2.26556443707464], [-3.98443603238417, -3.12195121951219e-2]), \\ ([-6.01034495038369, 4.64224263704784e1], [-3.12306684752388e1, -2.91187528315533e-2]).$$

Man erkennt, dass auch diese Einschließung keine optimale Einschließung des Wertebereichs  $W_Z$  darstellt. Es genügt daher nicht, zuerst mit  $T = \text{sqrt1px2}(Z)$  eine optimale Einschließung von  $\{\sqrt{1+z^2} \mid x \in X \wedge y \in Y\} \subseteq T$  und anschließend mit  $V = \text{reci_z2}(T)$  eine Einschließung von  $W_Z \subseteq V$  zu berechnen, um mit  $V$  eine optimale Einschließung von  $W_Z$  zu erhalten. Eine solche optimale Einschließung berechnet man nur mit dem Aufruf `res = reci_1pz2(Z)` aus dem **1. Beispiel**.

### C.2.5. $1/(1 - z^2)$

Ist  $z = x + i \cdot y \in \mathbb{C}$  ein Element aus einem komplexen Intervall  $Z = X + iY$  mit reellen Intervallen  $X, Y$ , so berechnet die Funktion

```
Mpfciclass reci_1mz2(const Mpfciclass& Z);
```

aus `mpfciclass.cpp` mit dem Aufruf `res = reci_1mz2(Z)`; eine optimale Einschließung der komplexen Menge

$$W_Z := \{1/(1 - z^2) \mid x \in X \wedge y \in Y\} \subseteq \text{res}, \quad \pm 1 \notin Z.$$

Mit  $f(z) := 1/(1 - z^2)$ ,  $z \in Z$ , und  $z = i \cdot \hat{z}$  erhält man  $f(z) = 1/(1 + \hat{z}^2)$ , mit  $\hat{z} = y - i \cdot x$  und

$$W_Z = \{1/(1 + \hat{z}^2) \mid \Re(\hat{z}) \in Y \wedge \Im(\hat{z}) \in -X\}, \quad \hat{z} \neq \pm i.$$

Damit kann man die nahezu optimale Einschließung von  $W_Z$  zurückführen auf die Berechnung von `reci_1pz2( $\hat{Z}$ )`, mit  $\hat{Z} := Y - i \cdot X$ , und genau dies wird bei der Implementierung der obigen Funktion `reci_1mz2(...)` realisiert. Beachten Sie, dass  $1/(1 - z^2)$  die 1. Ableitung der `acoth`-Funktion ist.

#### C.2.5.1. Numerische Beispiele

Im **1. Beispiel** wählen wir  $h(z) := 1/(2 + z - 3z^2)$  und zeigen, dass wir mit  $X = [1/2, 7/4]$ ,  $Y = [2, 4]$  und der Funktion `reci_1mz2(...)` ebenfalls eine fast optimale Einschließung der Wertemenge

$$W_{Z,h} = \{1/(2 + z - 3z^2) \mid x \in X \wedge y \in Y\}$$

berechnen können. Mit Hilfe der Umformung

$$\frac{1}{2 + z - 3z^2} = 0.48 \cdot \frac{1}{1 - w^2}, \quad \text{und } w := \frac{6z - 1}{5}$$

wird zunächst in der komplexen Ebene das Rechteck  $Z$  mit nur minimalen Überschätzungen in ein achsenparalleles Rechteck  $W := (6 \diamond Z \diamond 1) \diamond 5$  abgebildet, so dass dann mit mit der Anweisung `T = 0.48  $\diamond$  reci_1mz2(W)` eine nahezu optimale Einschließung  $W_{Z,h} \subset T$  berechnet werden kann. Mit der Präzision `prec = 53` erhalten wir

$$T = ([1.30737618861925e - 2, 6.70527278268821e - 2], [3.15076923076922e - 3, 4.81179721970281e - 2]).$$

Im **2. Beispiel** wählen wir die gleiche Funktion  $h(z)$  und die gleichen Intervalle  $X, Y$  wie im 1. Beispiel und zeigen, dass die naive Intervallauswertung von  $h(z)$  im Vergleich zum 1. Beispiel zu starken Überschätzungen von  $W_{Z,h}$  führt:

$$([1.30737618861925e - 2, 6.70527278268821e - 2], [3.15076923076922e - 3, 4.81179721970281e - 2]), \\ ([3.26276015114256e - 3, 1.64868468905322e - 1], [7.67754318618042e - 4, 9.41176470588236e - 2]).$$

Im **3. Beispiel** wählen wir  $h(z)$  und  $X, Y$  wie im 1. Beispiel und zeigen mit der Umformung

$$h(z) = \frac{1}{2 + z - 3z^2} = -\frac{1}{3} \cdot \frac{1}{z^2 - z/3 - 2/3}, \quad \text{und } A = -1 \diamond 3, B = -2 \diamond 3,$$

dass auch die optimale Einschließung von  $z^2 - z/3 - 2/3$  mit Hilfe der Funktion `poly2(z, A, B)` und anschließender Division ebenfalls zu starken Überschätzungen von  $W_{Z,h}$  führt:

$$([1.30737618861925e - 2, 6.70527278268821e - 2], [3.15076923076922e - 3, 4.81179721970281e - 2]), \\ ([4.41305107318572e - 3, 1.11103762978640e - 1], [1.60574052236746e - 3, 7.61904761904763e - 2]).$$

Die im 3. Beispiel berechnete Einschließung ist jedoch erwartungsgemäß besser als die aus dem 2. Beispiel, da `poly2(z, A, B)` schon eine optimale Einschließung des Nenners berechnet.

### C.2.6. $z^2 + a \cdot z + b$

Mit den gegebenen komplexen Intervallen  $Z = X + i \cdot Y$ ,  $A = A_1 + i \cdot A_2$ ,  $B = B_1 + i \cdot B_2$  ist für ein komplexes Intervallpolynom 2. Grades

$$(C.61) \quad W = Z^2 + A \cdot Z + B := \{w = z^2 + a \cdot z + b \in \mathbb{C} \mid z \in Z, a \in A, b \in B\}$$

eine optimale Einschließung des komplexen Intervalls  $W$  zu berechnen. Mit den reellen Intervallen  $X, Y, A_1, A_2, B_1, B_2$  erhält man eine Einschließung von  $W$  durch

$$(C.62) \quad W \subset P := \left( X + \frac{A_1}{2} \right)^2 - \left( \frac{A_1}{2} \right)^2 + \left( Y + \frac{A_2}{2} \right)^2 - \left( \frac{A_2}{2} \right)^2 + B_1 \\ + i \cdot \left\{ 2 \cdot \left( X + \frac{A_1}{2} \right) \left( Y + \frac{A_2}{2} \right) - 2 \cdot \frac{A_1}{2} \cdot \frac{A_2}{2} + B_2 \right\}.$$

Die obige Einschließung ist sogar **optimal**, d.h. es gilt  $W = P$ , wenn der komplexe Intervallausdruck  $P$  zur Vermeidung von Auslöschung in z.B. doppelter Präzision ausgewertet wird und wenn nur  $A_1, A_2$  als Punktintervalle gewählt werden, d.h. die Intervalle  $X, Y, B_1, B_2$  können beliebige echte Intervalle sein, da sie im Real- und Imaginärteil jeweils nur einmal vorkommen, vgl. Seite 74. Die Einschließung  $W \subseteq P$  bleibt auch dann noch nahezu optimal, wenn die Intervalle  $A_1, A_2$  z.B. eine nicht darstellbare komplexe Zahl  $a = \sqrt{2} - i \cdot \sqrt{5} \in A_1 + i \cdot A_2$  in der Current-Precision selbst optimal einschließen und damit einen nur minimalen Durchmesser besitzen.

Bei der Auswertung der Intervallfunktionen in C.62 besteht noch ein weiteres Problem. Beispielsweise kann die Berechnung von  $(A_1/2)^2$  zu einem vorzeitigen Überlauf führen, während der ganze Realteil in C.62 ein noch darstellbares Intervall ist. Dies kann z.B. eintreten, wenn man im Fall  $B = 0$  für  $Z$  die Nullstelle  $Z_1 = -A$  des Polynoms  $Z^2 + A \cdot Z + 0$  wählt und dabei das komplexe Punktintervall  $A := \text{MaxFloat}() + i \cdot \text{MaxFloat}()/2$  definiert. Wegen  $A_1 = \text{MaxFloat}()$  würde dann die Auswertung von C.62 mit Sicherheit zu einem vorzeitigen Überlauf führen. Um diesen Überlauf zu vermeiden, wird das Intervall  $A_1$  wie folgt skaliert<sup>5</sup>  $A_1 \subseteq 2^{k_1} \cdot \hat{A}_1$ , wobei  $\hat{A}_1$  möglichst groß ausfällt, aber  $2 \cdot \hat{A}_1 \cdot \hat{A}_1$  noch keinen Überlauf erzeugen darf. In diesem Sinne wird in C.62 der ganze Realteil berechnet, und erst am Schluss wird mit dem entsprechenden Faktor  $2^k$  multipliziert, so dass erst dann ein möglicher und nicht zu vermeidender Überlauf eintreten kann. Ganz analog wird der Imaginärteil berechnet. Die Berechnung der Polynomeinschließung  $P$  erfolgt mit Hilfe der Funktion

```
Mpfciclass poly2(const Mpfciclass& Z, const Mpfciclass& A, const Mpfciclass& B);
```

die in `mpfciclass.cpp` definiert ist.

#### C.2.6.1. Numerische Beispiele

Im **1. Beispiel** wählen wir wie oben

$$Z = [-\text{MaxFloat}(), -\text{MaxFloat}()] + i \cdot [-\text{MaxFloat}()/2, -\text{MaxFloat}()/2] = -A, \\ B = [0, 0] + i \cdot [0, 0],$$

so dass das Polynom  $Z^2 + A \cdot Z + 0$  an seiner Nullstelle  $Z_1 = -A$  auszuwerten ist, und tatsächlich erhält man mit dem Funktionsaufruf `P = poly2(Z,A,B);` **ohne** einen vorzeitigen Überlauf die **optimale** Einschließung

$$P = [0, 0] + i \cdot [0, 0] = W.$$

Wertet man das Polynom naiv aus, d.h. durch `sqr(Z) \diamond A \diamond Z \diamond B`, so erhält man die völlig unbrauchbare Einschließung

$$P = ([-\text{Inf}(), \text{Inf}()], [-\text{Inf}(), \text{Inf}()]).$$

<sup>5</sup>Die Skalierung erfolgt mit Hilfe der Funktion `Interval_Scaling(...)`, vgl. Seite 72.

Im **2. Beispiel** wählen wir  $Z = [1, 2] + i \cdot [-2, -1]$ ,  $A = [-2, -2] + i \cdot [2, 2]$ ,  $B = [2, 2] + i \cdot [-1, -1]$  und erhalten mit dem Funktionsaufruf  $P = \text{poly2}(Z, A, B)$ ; die **optimale** Einschließung

$$P = [1, 3] + i \cdot [-1, +1] = W.$$

Wertet man das Polynom jedoch naiv aus, d.h. durch  $\text{sqr}(Z) \diamond A \diamond Z \diamond B$ , so erhält man die viel zu grobe Einschließung

$$P = [-3, +7] + i \cdot [-5, +5].$$

Das Beispiel zeigt, wie wichtig es ist, wenigstens für einfache komplexe Intervallausdrücke, wie hier beim Polynom 2. Grades, optimale Einschließungen zu realisieren.

### C.2.7. $\log(z)$

Der Hauptwert der komplexen Logarithmusfunktion zur Basis  $e$  ist für  $z = x + i \cdot y \in \mathbb{C}$  definiert durch

$$\log(z) := \ln(|z|) + i \cdot \arg(z), \quad |z| := \sqrt{x^2 + y^2}, \quad -\pi < \arg(z) \leq +\pi.$$

Der Verzweigungsschnitt ist wie üblich die negative reelle Achse, wobei der Ursprung  $(0, 0)$  selbst der Verzweigungspunkt ist.  $\ln(|z|)$  bedeutet dabei den reellen Logarithmus zur Basis  $e$  mit dem Argument  $|z|$ . Ist nun  $Z$  ein komplexes Rechteckintervall, das die Null nicht enthält<sup>6</sup> und die negative reelle Achse nur von oben berühren darf, so enthält die komplexe Zahlenmenge

$$\log(Z) := \{y \in \mathbb{C} \mid y = \log(z) \wedge z \in Z\}$$

nur Funktionswerte des Hauptwertes der Logarithmusfunktion. Beachten Sie, dass  $\log(Z)$  i.a. kein Rechteckintervall ist. Beispielsweise ist das Bild  $\log(Z)$  der Parallelen  $Z = [1, 1] + i \cdot [1, 2]$  zur  $y$ -Achse in der komplexen Ebenen eine glatte Kurve, deren Einschließung durch ein achsenparalleles Rechteck mit der üblichen Überschätzung verbunden ist, vgl. dazu auch Seite 243.

Bezeichnen wir mit  $\text{Ln}(Z)$  die Einschließung von  $\log(Z)$  durch ein achsenparalleles Rechteck, so enthält auch dieses Rechteck nur Funktionswerte des Hauptwertes  $\log(z)$  und es gilt

$$\log(Z) := \{y \in \mathbb{C} \mid y = \log(z) \wedge z \in Z\} \subseteq \text{Ln}(Z).$$

Wir bezeichnen daher  $\text{Ln}(Z)$  als **Analytische Logarithmusfunktion**, wobei das Rechteckintervall  $Z$  die Null nicht enthalten darf und die negative reelle Achse nur von oben berühren darf. Man sollte noch beachten, dass z.B. im Fall  $Z = [-2, -1] + i \cdot [0, 1]$  die Einschließung

$$\text{Ln}(Z) = ([0, 8.04719\text{e-}1], [2.35619, 3.14160])$$

mit dem Supremum ihres Imaginärteils einen Wert liefert, der nicht mehr zum Hauptwert des Logarithmus gehört, da das maximale Argument  $\pi$  auf der Maschine durch einen etwas größeren Zahlenwert eingeschlossen werden muss!

Durch analytische Fortsetzung können Funktionswerte der Logarithmusfunktion auch für Rechteckintervalle  $Z$  definiert werden, die den Verzweigungsschnitt in ihrem Innern enthalten, wobei aber die Null weder in  $Z$  noch auf seinem Rand liegen darf. Durch analytische Fortsetzung ist es jetzt auch möglich, dass  $Z$  die negative reelle Achse von unten berührt, wobei dann aber die Funktionswerte auf dem Verzweigungsschnitt den Imaginärteil  $-\pi$  erhalten. Bezeichnen wir mit  $\log(Z)$  auch jetzt wieder die Menge aller dieser Funktionswerte, mit  $z \in Z$ , die also auch durch analytische Fortsetzung definiert sein können und wird deren Einschließung durch Rechteckintervalle mit  $\ln(Z)$  bezeichnet, so gilt

$$\log(Z) \subseteq \ln(Z), \quad z \in Z.$$

Da  $\ln(Z)$ , je nach Lage von  $Z$ , Funktionswerte einschließen kann, die nicht zum Hauptwert der Logarithmusfunktion gehören, wird  $\ln(Z)$  auch als **Nicht-analytische Logarithmusfunktion** bezeichnet. Weitere Einzelheiten zu den erlaubten Intervallen  $Z$  findet man auf Seite 295.

---

<sup>6</sup>Die Null darf auch nicht auf dem Rand von  $Z$  liegen!

### C.2.7.1. Analytische Funktion

Bedeutet  $\log(z)$  den auf Seite 293 definierten Hauptwert der komplexen Logarithmusfunktion und ist  $Z$  ein achsenparalleles Rechteckintervall, wobei die Null nicht in seinem Innern oder auf dem Rand von  $Z$  liegen darf und  $Z$  die negative reelle Achse auch nur von oben berühren darf, so liefert die analytische Logarithmusfunktion  $\text{Ln}(Z)$  eine nahezu optimale achsenparallele Rechteck-Einschließung der komplexen Zahlenmenge

$$\log(Z) := \{y \in \mathbb{C} \mid y = \log(z) \wedge z \in Z\} \subseteq \text{Ln}(Z).$$

In Abbildung C.17 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  angegeben.

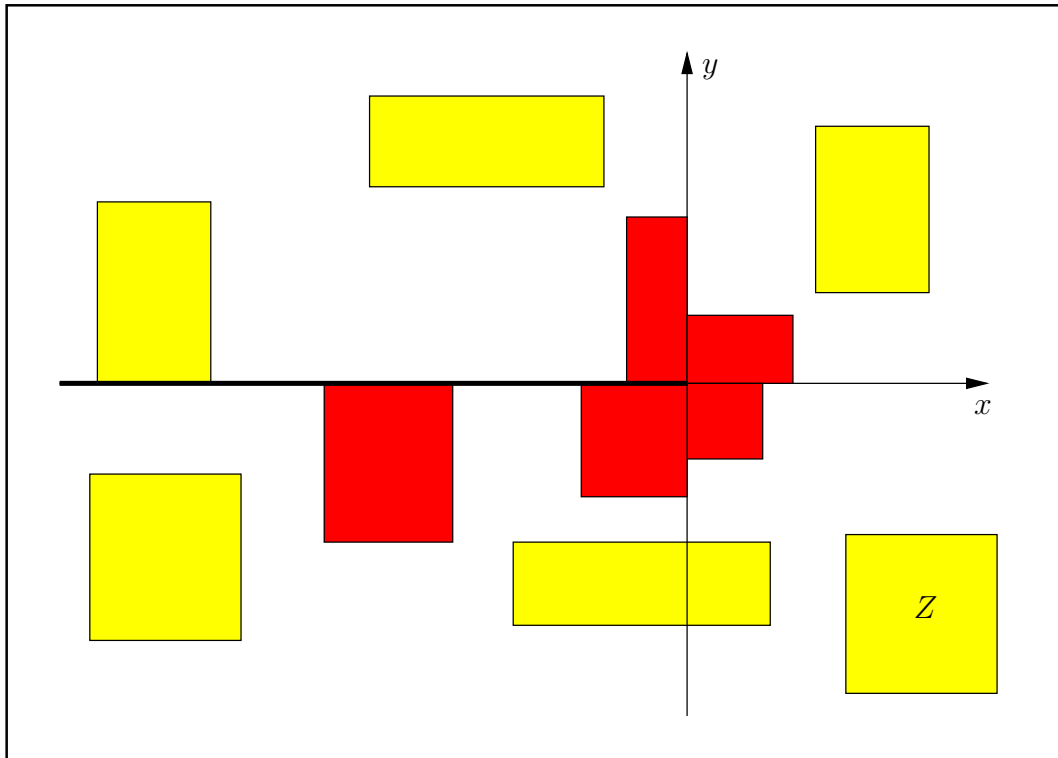


Abbildung C.17.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{Ln}(Z)$ .

Hier einige numerische Beispiele in der niedrigen Präzision  $\text{prec} = 30$ , womit etwa  $30/\log_2(10) \approx 30/3.321928095 \approx 9$  Dezimalstellen berechnet werden:

$$\begin{aligned} \text{Ln}([+2, +4] + i \cdot [+0, +0]) &= ([6.93147180e - 1, 1.38629437], [0.00000000, 0.00000000]), \\ \text{Ln}([-4, -2] + i \cdot [+0, +0]) &= ([6.93147180e - 1, 1.38629437], [3.14159265, 3.14159266]), \\ \text{Ln}([-2, -1] + i \cdot [+0, +1]) &= ([0.00000000, 8.04718957e - 1], [2.35619449, 3.14159266]), \\ \text{Ln}([-2, +1] + i \cdot [-2, -1]) &= ([0.00000000, 1.03972078], [-2.67794505, -7.85398163e - 1]), \\ \text{Ln}([+2, +3] + i \cdot [-3, -2]) &= ([1.03972077, 1.44518588], [-9.82793724e - 1, -5.88002603e - 1]), \\ \text{Ln}([+2, +3] + i \cdot [+2, +3]) &= ([1.03972077, 1.44518588], [5.88002603e - 1, 9.82793724e - 1]), \\ \text{Ln}([-2, +2] + i \cdot [+2, +3]) &= ([6.93147180e - 1, 1.28247468], [7.85398163e - 1, 2.35619450]), \\ \text{Ln}([-2, +2] + i \cdot [+0, +1]) &\rightsquigarrow \text{MpfciClass LN( const MpfciClass\& z ); z contains 0.} \end{aligned}$$

### C.2.7.2. Nicht-analytische Funktion

Ausgangspunkt ist zunächst wieder der auf Seite 293 definierte Hauptwert  $\log(z)$  der komplexen Logarithmusfunktion. Ist  $Z$  dann ein achsenparalleles komplexes Rechteckintervall, das die Null nicht enthält und die negative reelle Achse nur von oben berühren darf, so enthält die komplexe Zahlenmenge

$$\log(Z) := \{y \in \mathbb{C} \mid y = \log(z) \wedge z \in Z\}$$

nur Funktionswerte des Hauptwertes der Logarithmusfunktion. Zusätzlich wird aber jetzt noch vorausgesetzt, dass  $Z$  den Verzweigungsschnitt in seinem Innern enthalten darf, wobei die Funktionswerte in der unteren Halbebene durch analytische Fortsetzung aus der oberen Halbebene definiert werden. Ist dann  $\ln(Z)$  eine Rechteck-Einschließung dieser Zahlenmenge  $\log(Z)$ , so gilt

$$\log(Z) \subseteq \ln(Z), \quad z \in Z.$$

$\ln(Z)$  wird dabei als nicht-analytische Logarithmusfunktion bezeichnet, da jetzt auch Funktionswerte eingeschlossen werden können, die nicht zum Hauptwert der Logarithmusfunktion gehören. In Abbildung C.18 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  der Funktion  $\ln(Z)$  angegeben.

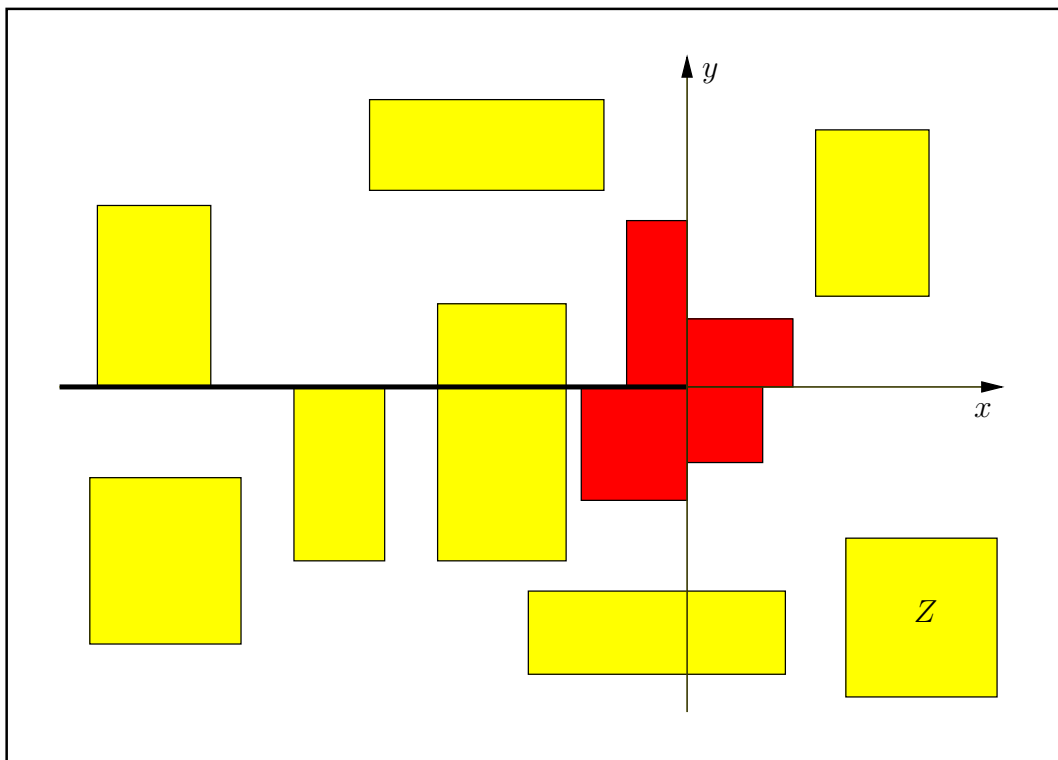


Abbildung C.18.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\ln(Z)$ .

Hier einige numerische Beispiele in der niedrigen Präzision  $\text{prec} = 30$ , womit etwa  $30/\log_2(10) \approx 30/3.321928095 \approx 9$  Dezimalstellen berechnet werden:

$$\begin{aligned} \ln([+2, +4] + i \cdot [+0, +0]) &= ([6.93147180e - 1, 1.38629437], [0.00000000, 0.00000000]), \\ \ln([-4, -2] + i \cdot [+0, +0]) &= ([6.93147180e - 1, 1.38629437], [3.14159265, 3.14159266]), \\ \ln([-4, -2] + i \cdot [+0, +1]) &= ([6.93147180e - 1, 1.41660668], [2.67794504, 3.14159266]), \\ \ln([-4, -2] + i \cdot [-1, +0]) &= ([6.93147180e - 1, 1.41660668], [-3.14159266, -2.67794504]), \\ \ln([-4, -2] + i \cdot [-1, +1]) &= ([6.93147180e - 1, 1.41660668], [2.67794504, 3.60524027]), \\ \ln([-2, +2] + i \cdot [-3, -2]) &= ([6.93147180e - 1, 1.28247468], [-2.35619450, -7.85398163e - 1]). \end{aligned}$$

### C.2.8. $\log(1+z)$

Der Hauptwert der komplexen Logarithmusfunktion zur Basis  $e$  mit dem Argument  $1+z$  ist für  $z = x + i \cdot y \in \mathbb{C}$  definiert durch

$$\log(1+z) := \ln(|1+z|) + i \cdot \arg(1+z), \quad |1+z| := \sqrt{(1+x)^2 + y^2}, \quad -\pi < \arg(1+z) \leq +\pi.$$

Der Verzweigungsschnitt ist die negative reelle Achse von  $-\infty$  bis  $-1$ , wobei der Punkt  $(-1, 0)$  selbst der Verzweigungspunkt ist.  $\ln(|1+z|)$  bedeutet dabei den reellen Logarithmus zur Basis  $e$  mit dem Argument  $|1+z|$ . Ist nun  $Z$  ein komplexes Rechteckintervall, das  $(-1, 0)$  nicht enthält<sup>7</sup> und den Verzweigungsschnitt nur von oben berühren darf, so enthält die komplexe Zahlenmenge

$$\log(1+Z) := \{y \in \mathbb{C} \mid y = \log(1+z) \wedge z \in Z\}$$

nur Funktionswerte des Hauptwertes der Logarithmusfunktion  $\log(1+z)$ . Beachten Sie, dass  $\log(1+Z)$  i.a. kein Rechteckintervall ist. Beispielsweise ist das Bild  $\log(1+Z)$  der Parallelen  $Z = [1, 1] + i \cdot [1, 2]$  zur  $y$ -Achse in der komplexen Ebenen eine glatte Kurve, deren Einschließung durch ein achsenparalleles Rechteck mit der üblichen Überschätzung verbunden ist, vgl. dazu auch Seite 243.

Bezeichnen wir mit  $\text{Lnp1}(Z)$  die Einschließung von  $\log(1+Z)$  durch ein achsenparalleles Rechteck, so enthält auch dieses Rechteck nur Funktionswerte des Hauptwertes  $\log(1+z)$  und es gilt

$$\log(1+Z) := \{y \in \mathbb{C} \mid y = \log(1+z) \wedge z \in Z\} \subseteq \text{Lnp1}(Z).$$

Wir bezeichnen daher  $\text{Lnp1}(Z)$  als **Analytische Logarithmusfunktion**, wobei das Rechteckintervall  $Z$  die Singularität  $z_S = -1 + 0 \cdot i$  nicht enthalten darf und den Verzweigungsschnitt nur von oben berühren darf. Man sollte noch beachten, dass z.B. im Fall  $Z = [-3, -2] + i \cdot [0, 1]$  die Einschließung

$$\text{Lnp1}(Z) = ([0, 8.04719\text{e-}1], [2.35619, 3.14160])$$

mit dem Supremum ihres Imaginärteils einen Wert liefert, der nicht mehr zum Hauptwert des Logarithmus gehört, da das maximale Argument  $\pi$  auf der Maschine durch einen etwas größeren Zahlenwert eingeschlossen werden muss!

Durch analytische Fortsetzung können Funktionswerte der Logarithmusfunktion auch für Rechteckintervalle  $Z$  definiert werden, die den Verzweigungsschnitt in ihrem Innern enthalten, wobei aber die Singularität  $z_S = -1 + 0 \cdot i$  weder in  $Z$  noch auf seinem Rand liegen darf. Durch analytische Fortsetzung ist es jetzt auch möglich, dass  $Z$  den Verzweigungsschnitt von unten berührt, wobei dann aber die Funktionswerte auf dem Verzweigungsschnitt den Imaginärteil  $-\pi$  erhalten. Bezeichnen wir mit  $\log(1+Z)$  auch jetzt wieder die Menge aller dieser Funktionswerte, mit  $z \in Z$ , die also auch durch analytische Fortsetzung definiert sein können und wird deren Einschließung durch Rechteckintervalle mit  $\text{lnp1}(Z)$  bezeichnet, so gilt

$$\log(1+Z) \subseteq \text{lnp1}(Z), \quad z \in Z.$$

Da  $\text{lnp1}(Z)$ , je nach Lage von  $Z$ , Funktionswerte einschließen kann, die nicht zum Hauptwert der Logarithmusfunktion gehören, wird  $\text{lnp1}(Z)$  auch als **Nicht-analytische Logarithmusfunktion** bezeichnet. Weitere Einzelheiten zu den erlaubten Intervallen  $Z$  findet man auf Seite 299.

---

<sup>7</sup>  $z_S = -1 + 0 \cdot i$  darf auch nicht auf dem Rand von  $Z$  liegen!



### C.2.8.1. Analytische Funktion

Bedeutet  $\log(1+z)$  den auf Seite 296 definierten Hauptwert der komplexen Logarithmusfunktion und ist  $Z$  ein achsenparalleles Rechteckintervall, wobei die Singularität  $z_S = -1 + 0 \cdot i$  nicht in seinem Innern oder auf dem Rand von  $Z$  liegen darf und  $Z$  den Verzweigungsschnitt auch nur von oben berühren darf, so liefert die analytische Logarithmusfunktion  $\text{Ln}p1(Z)$  eine nahezu optimale achsenparallele Rechteck-Einschließung der komplexen Zahlenmenge

$$\log(1+Z) := \{y \in \mathbb{C} \mid y = \log(1+z) \wedge z \in Z\} \subseteq \text{Ln}p1(Z).$$

In Abbildung C.19 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  angegeben.

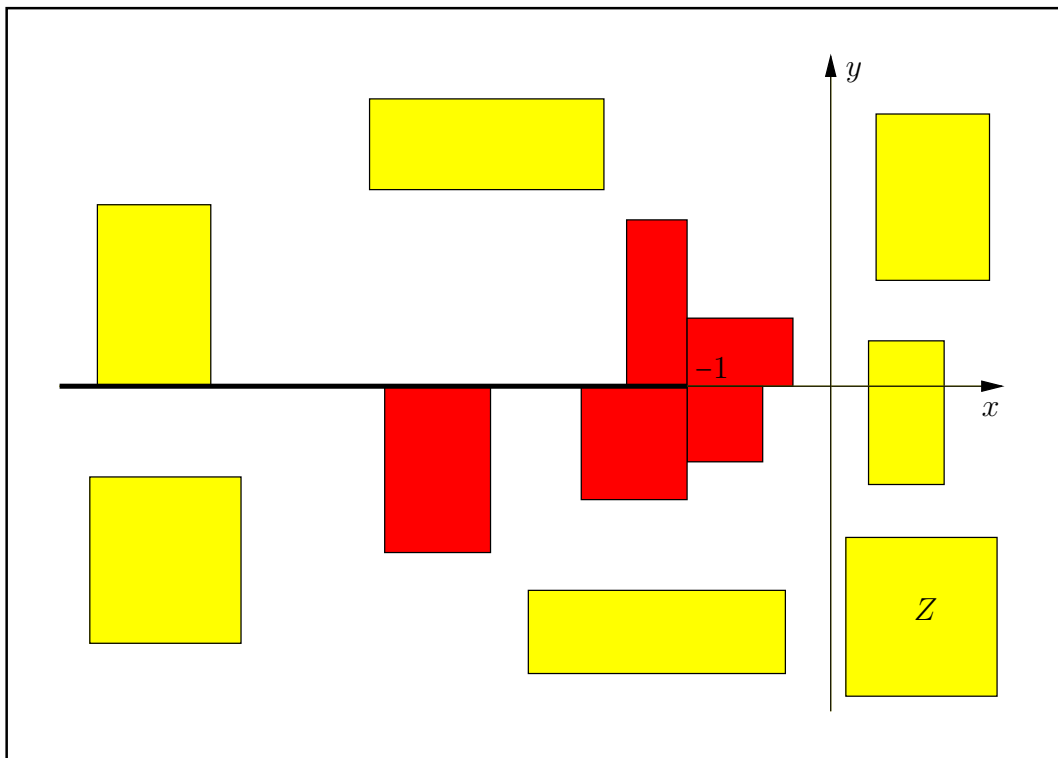


Abbildung C.19.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{Ln}p1(Z)$ .

Die Einschließung  $\text{Ln}p1(Z)$  wird berechnet mit der C-XSC-Funktion

```
MpfiClass Ln p1(const MpfiClass& Z);
```

die in der Datei `mpficlass.cpp` definiert ist. Mit  $Z = X + i \cdot Y$  wird der Realteil realisiert durch die C-XSC-Funktion

```
MpfiClass ln_sqrtxp1_2y2(const MpfiClass& X, const MpfiClass& Y);
```

die in der Datei `mpficlass.cpp` definiert ist, vgl. Seite 205. Der Imaginärteil wird berechnet mit der C-XSC-Funktion

```
MpfiClass argp1(const MpfiClass& Z, PrecisionType prec);
```

die in der Datei `mpficlass.cpp` definiert ist, vgl. Seite 116.

## Numerische Ergebnisse:

Bezeichnungen:  $Z = X + i \cdot Y$ ,  $X = [x_1, x_2]$ ,  $Y = [y_1, y_2]$ ,  $t = \log(1 + Z)$ ;

1. Mit  $x_1 = x_2 = \text{minfloat}() = 2^{-1073741824}$ ;  $y_1 = y_2 = \text{minfloat}()$  und  $\text{prec}=30\,000$  erhält man für  $t$  die notwendigerweise grobe Einschließung:

$$t \subset ([0, 2.38256490488...41164e - 323228497], [0, 2.38256490488...41164e - 323228497]).$$

2. Mit  $x_1 = x_2 = 2 \cdot \text{minfloat}() = 2^{-1073741823}$ ;  $y_1 = y_2 = 2 \cdot \text{minfloat}()$  und  $\text{prec}=30\,000$  erhält man für  $t$  die fast optimale Einschließung, wobei Real- und Imaginärteil mit jeweils 9030 gemeinsamen Dezimalziffern berechnet werden:

$$t \subset ([4.765129809775...82325e - 323228497, 4.765129809775...82328e - 323228497], \\ [4.765129809775...82324e - 323228497, 4.765129809775...82327e - 323228497]).$$

3. Mit  $x_1 = y_1 = 2 \cdot \text{minfloat}() = 2^{-1073741823}$ ,  $x_2 = y_2 = 4 \cdot \text{minfloat}()$ , und  $\text{prec}=30\,000$  erhält man für  $t$  die fast optimale Einschließung, wobei Infimum und Supremum von Real- und Imaginärteil mit jeweils 9030 korrekten Dezimalziffern berechnet werden:

$$t \subset ([4.765129809775...82325e - 323228497, 9.530259619551...64655e - 323228497], \\ [4.765129809775...82324e - 323228497, 9.530259619551...64653e - 323228497]).$$

4. Mit  $x_1 = x_2 = -1$ ,  $y_1 = y_2 = \text{minfloat}() = 2^{-1073741824}$  und  $\text{prec}=30\,000$  erhält man für  $t$  die fast optimale Einschließung, wobei Real- und Imaginärteil mit jeweils 9030 gemeinsamen Dezimalziffern berechnet werden:

$$t \subset ([-7.44261117954...03830425613365e8, -7.44261117954...03830425613363e8], \\ [1.570796326794...67470301701083272, 1.570796326794...7470301701083273]).$$

Beachten Sie, dass das Argument  $Z = [-1, -1] + i \cdot [\text{minfloat}(), \text{minfloat}()]$  jetzt in der unmittelbaren Nähe der Singularität  $z_S = -1 + i \cdot 0$  liegt und damit einen betragsmäßig sehr großen Realteil verursacht, der nahezu optimal eingeschlossen wird.

5. Mit  $x_1 = x_2 = -2^{3\,000\,001}$ ,  $y_1 = y_2 = 2^{-1\,500\,000}$  und  $\text{prec}=800\,000$  erhält man für  $t$  die fast optimale Einschließung, wobei Real- und Imaginärteil mit jeweils **240823** gemeinsamen Dezimalziffern berechnet werden:

$$t \subset ([1.327168685081681...30691e - 1806181, 1.327168685081681...30692e - 1806181], \\ [1.0150887817716533...57311e - 451545, 1.0150887817716533...57314e - 451545]).$$

Die sehr große Präzision  $\text{prec}=800\,000$  ist hier notwendig, um den Realteil von  $t$  nahezu optimal einschließen zu können, vgl. dazu die Bemerkungen auf Seite 269.

6. Mit  $x_1 = -3$ ,  $x_2 = -2$ ,  $y_1 = y_2 = 0$  und  $\text{prec}=30\,000$  erhält man für  $t$  die fast optimale Einschließung, wobei jetzt nur der Imaginärteil mit 9030 gemeinsamen Dezimalziffern eingeschlossen wird:

$$t \subset ([0, 6.931471805599453094172321214581765680755001343602...6258516803e - 1], \\ [3.141592653589793...6034021665442, 3.141592653589793...6034021665448]).$$

Beachten Sie, dass  $Z = [-3, -2] + i \cdot [0, 0]$  den Verzweigungsschnitt von oben berührt.

7.  $Z = [-3, -2] + i \cdot [-1, 0]$  und  $Z = [-1, 0] + i \cdot [0, 1]$  sind nicht erlaubte komplexe Intervalle und erzeugen entsprechende Fehlermeldungen.

### C.2.8.2. Nicht-analytische Funktion

Ausgangspunkt ist zunächst wieder der schon auf Seite 296 definierte Hauptwert  $\log(1+z)$  der komplexen Logarithmusfunktion. Ist  $Z$  dann ein achsenparalleles komplexes Rechteckintervall, das die Singularität  $z_S = -1 + i \cdot 0$  nicht enthält und den Verzweigungsschnitt nur von oben berührt, so enthält die komplexe Zahlenmenge

$$\log(1+Z) := \{y \in \mathbb{C} \mid y = \log(1+z) \wedge z \in Z\}$$

nur Funktionswerte des Hauptwertes  $\log(1+z)$ . Zusätzlich wird aber jetzt noch vorausgesetzt, dass  $Z$  den Verzweigungsschnitt in seinem Innern enthalten darf, wobei die Funktionswerte in der unteren Halbebene durch analytische Fortsetzung aus der oberen Halbebene definiert werden. Ist dann  $\text{lnp1}(Z)$  eine Rechteck-Einschließung dieser Zahlenmenge  $\log(1+Z)$ , so gilt

$$\log(1+Z) \subseteq \text{lnp1}(Z), \quad z \in Z.$$

$\text{lnp1}(Z)$  wird dabei als nicht-analytische Logarithmusfunktion bezeichnet, da jetzt auch Funktionswerte eingeschlossen werden können, die nicht zum Hauptwert der Logarithmusfunktion gehören. In Abbildung C.20 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  der Funktion  $\text{lnp1}(Z)$  angegeben.

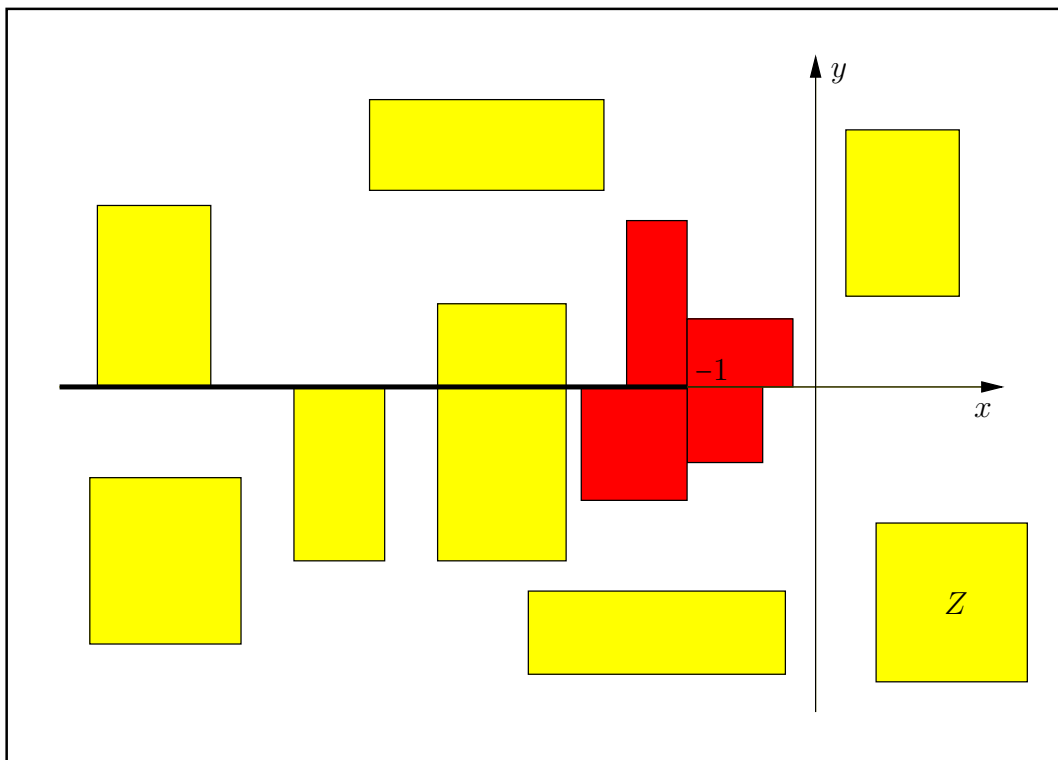


Abbildung C.20.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{lnp1}(Z)$ .

#### Numerische Ergebnisse:

Bezeichnungen:  $Z = X + i \cdot Y$ ,  $X = [x_1, x_2]$ ,  $Y = [y_1, y_2]$ ,  $t = \log(1+Z)$ ;

1. Mit Ausnahme des letzten Beispiels gelten alle numerischen Ergebnisse von Seite 298 auch für die nicht-analytische Funktion  $\text{lnp1}(Z)$ .
2. Mit  $Z = [-3, -2] + i \cdot [-1, 0]$  und  $\text{prec} = 30\,000$  erhält man die Einschließung

$$t \subset ([0, 8.04718956217050187300379666613093819762800677134258860\dots 765881e - 1], [2.35619449019234492884698\dots 249077, 3.92699081698724154807830\dots 081817]).$$

### C.2.9. $\sqrt{z}$

Der Verzweigungsschnitt der komplexen Quadratwurzel ist in der komplexen Ebene wie üblich die negative reelle Achse. Ein achsenparalleles Argumentintervall  $Z$  darf diesen Verzweigungsschnitt nicht im Innern enthalten und auch nicht von unten berühren.  $\sqrt{Z}$  liefert eine achsenparallele Einschließung der komplexen Zahlenmenge

$$\{y \in \mathbb{C} \mid y = \sqrt{z}, z \in Z\} \subseteq \sqrt{Z},$$

wobei  $\sqrt{z}$  als Hauptwert der komplexen Quadratwurzel definiert ist.

In Abbildung C.21 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  der Funktion  $\sqrt{Z}$  angegeben.

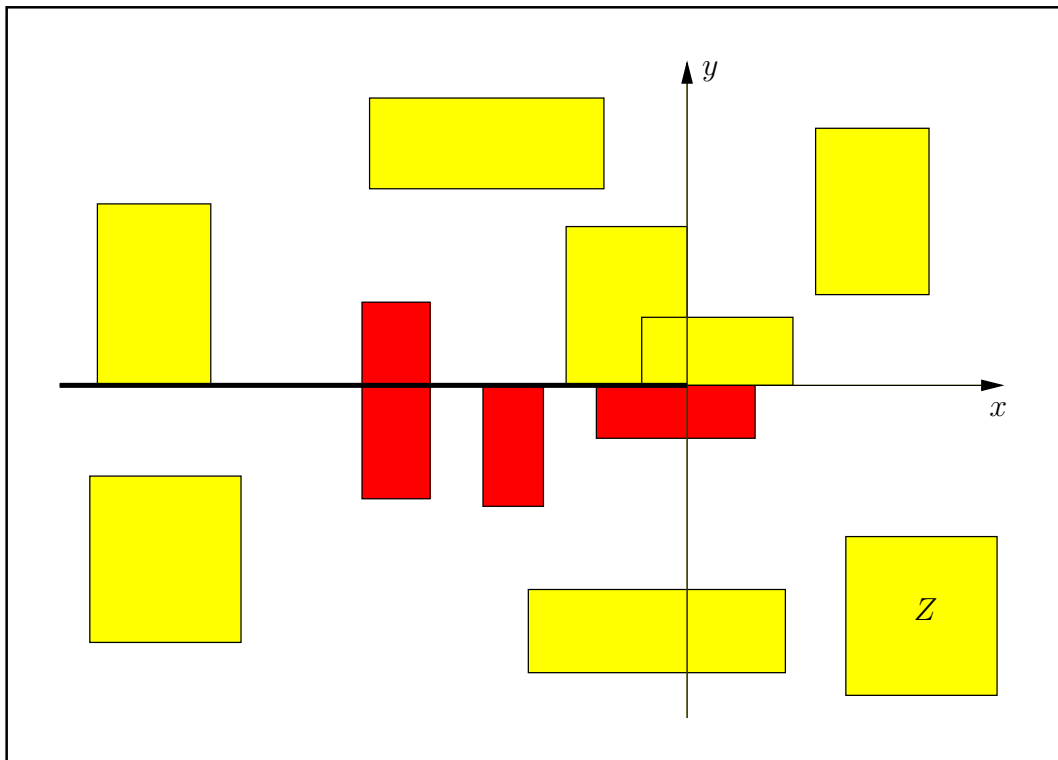


Abbildung C.21.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\sqrt{Z}$ .

Hier einige numerische Beispiele in der niedrigen Präzision  $\text{prec} = 30$ , womit etwa  $30/\log_2(10) \approx 30/3.321928095 \approx 9$  Dezimalstellen berechnet werden:

```

sqrt([+0,+0] + i * [+0,+0]) = ([0.00000000, 0.00000000], [0.00000000, 0.00000000]),
sqrt([+0,+0] + i * [+1,+2]) = ([7.07106781e - 1, 1.00000001], [7.07106781e - 1, 1.00000001]),
sqrt([+0,+0] + i * [-2,-1]) = ([7.07106781e - 1, 1.00000001], [-1.00000001, -7.07106781e - 1]),
sqrt([-4,-2] + i * [+0,+2]) = ([0.00000000, 6.43594253e - 1], [1.41421356, 2.05817103]),
sqrt([-1,+0] + i * [+0,+1]) = ([0.00000000, 7.07106782e - 1], [0.00000000, 1.09868412]),
sqrt([-1,+2] + i * [+0,+1]) = ([0.00000000, 1.45534670], [0.00000000, 1.09868412]),
sqrt([-1,+2] + i * [-3,-2]) = ([7.86151377e - 1, 1.67414923], [-1.44261528, -6.43594252e - 1]),
sqrt([+2,+3] + i * [-2,+1]) = ([1.41421356, 1.81735403], [-6.43594253e - 1, 3.43560750e - 1]),
sqrt([-3,-2] + i * [-2,+2]) ~
MpfciClass sqrt(const MpfciClass& z); z not in the principal branch.

```

## C.2.10. $\sqrt{z}$ , Beide Quadratwurzeln

Während mit `sqrt(Z)` nach Seite 300 eine achsenparallele Einschließung nur der Hauptwerte  $\sqrt{z}$ ,  $z \in Z$ , geliefert wird, werden mit `sqrt_all(Z)` für **beliebige** achsenparallele Argumentintervalle  $Z$  in einer Liste die **beiden** möglichen Einschließungen für  $\pm\sqrt{z}$ ,  $z \in Z$ , berechnet. Mit dem folgenden Programm

```
1 // MPFR-07.cpp
2 #include "mpfciclass.hpp"
3
4 using namespace MPFR;
5 using namespace MPFI;
6 using namespace cxsc;
7 using namespace std;
8
9 int main(void)
10 {
11     interval re(-4,4), im(0,0);
12     Mpfciclass Z(re, im, 20);
13     cout.precision(30/3.321928095); // Ausgabe mit 9 Dez.-Stellen
14     cout << "Z = " << Z << endl;
15     cout << "Z.GetPrecision() = " << Z.GetPrecision() << endl;
16
17     list<Mpfciclass> res;
18     // Berechnung beider Einschliessungen in einer Liste:
19     res = sqrt_all(Z);
20
21     list<Mpfciclass>::iterator pos;
22     // Ausgabe der beiden Einschliessungen:
23     for (pos = res.begin(); pos != res.end(); ++pos )
24     {
25         cout << *pos << endl; // Jede Einschliessung in neue Zeile
26         cout << "Praezision = " << (*pos).GetPrecision() << endl;
27     }
28
29     return 0;
30 }
```

erhält man die Ausgabe

```
Z = ([-4.00000000,4.00000000], [0,-0])
Z.GetPrecision() = 20
([0,2.00000000], [0,2.00000000])
Praezision = 53
([-2.00000000,-0], [-2.00000000,-0])
Praezision = 53
```

wobei nach jeder Einschließung, die mit jeweils 9 Dezimalstellen erfolgt, auch noch die Präzision dieser berechneten Einschließung ausgegeben wird. Hier wird mit der Default-Präzision `prec = 53` gerechnet, da keine Current-Präzision vereinbart wurde.

Es folgen noch weitere Beispiele in der niedrigen Präzision `prec = 30`, womit etwa  $30/\log_2(10) \approx 30/3.321928095 \approx 9$  Dezimalstellen berechnet werden:

```
sqrt_all([-2,-1] + i * [-1,+1]) ~
  ([-4.55089861e-1,4.55089861e-1], [1.00000000,1.45534670])
  ([-4.55089861e-1,4.55089861e-1], [-1.45534670,-1.00000000]).
sqrt_all([-4,-2] + i * [-4,+0]) ~
  ([0.00000000,1.11178595], [-2.19736823,-1.41421356])
  ([-1.11178595,-0.00000000], [1.41421356,2.19736823]).
```

### C.2.11. $\sqrt[n]{z}$

Schreibt man eine komplexe Zahl  $z = x + i \cdot y \in \mathbb{C}$  in Polarkoordinaten, so gilt

$$z = r \cdot e^{i\varphi}, \quad r = \sqrt{x^2 + y^2} \geq 0, \quad -\pi < \varphi \leq +\pi.$$

Die komplexen Lösungen  $w_k \in \mathbb{C}$  der Gleichung

$$w^n = z, \quad n \in \mathbb{N}$$

sind gegeben durch die  $n$  komplexen Zahlen

$$(C.63) \quad w_k := \sqrt[n]{r} \cdot \left\{ \cos\left(\frac{\varphi + 2\pi k}{n}\right) + i \cdot \sin\left(\frac{\varphi + 2\pi k}{n}\right) \right\} = \sqrt[n]{r} \cdot e^{i\frac{\varphi + 2\pi k}{n}}, \quad k = 0, 1, 2, \dots, n-1,$$

wobei die  $w_k$  auch allgemein als  $\sqrt[n]{z}$  oder  $z^{1/n}$  bezeichnet werden. In diesem Abschnitt definieren wir den Hauptzweig der  $n$ -ten Wurzel aus  $z = r \cdot e^{i\varphi}$  durch:

$$\sqrt[n]{z} := w_0 = \sqrt[n]{r} \cdot \left\{ \cos\left(\frac{\varphi}{n}\right) + i \cdot \sin\left(\frac{\varphi}{n}\right) \right\} = \sqrt[n]{r} \cdot e^{i\frac{\varphi}{n}}, \quad -\pi < \varphi < +\pi.$$

Im **ersten Beispiel** betrachten wir  $z = -1 + i = \sqrt{2} \cdot \{\cos(3\pi/4) + i \cdot \sin(3\pi/4)\} = \sqrt{2} \cdot e^{i(3\pi/4)}$ . Dann gilt nach (C.63) für die dritten Wurzeln aus  $z$ , [62]

$$\begin{aligned} \sqrt[3]{-1+i} &:= w_0 = 2^{1/6} \cdot \{\cos(\pi/4) + i \cdot \sin(\pi/4)\}, \\ w_1 &= 2^{1/6} \cdot \{\cos(11\pi/12) + i \cdot \sin(11\pi/12)\}, \\ w_2 &= 2^{1/6} \cdot \{\cos(19\pi/12) + i \cdot \sin(19\pi/12)\}. \end{aligned}$$

Im **zweiten Beispiel** wählen wir  $z = -2\sqrt{3} - 2i = 4 \cdot \{\cos(-5\pi/6) + i \cdot \sin(-5\pi/6)\} = 4 \cdot e^{-5\pi/6}$ . Beachten Sie, dass jetzt  $\varphi = -5\pi/6$  negativ sein muss, da  $z$  in der unteren Halbebene liegt. Dann gilt wieder nach (C.63) für die vierten Wurzeln aus  $z$

$$\begin{aligned} \sqrt[4]{-2\sqrt{3}-2i} &:= w_0 = \sqrt{2} \cdot \{\cos(-5\pi/24) + i \cdot \sin(-5\pi/24)\}, \\ w_1 &= \sqrt{2} \cdot \{\cos(7\pi/24) + i \cdot \sin(7\pi/24)\}, \\ w_2 &= \sqrt{2} \cdot \{\cos(19\pi/24) + i \cdot \sin(19\pi/24)\}, \\ w_3 &= \sqrt{2} \cdot \{\cos(-17\pi/24) + i \cdot \sin(-17\pi/24)\}, \end{aligned}$$

wobei auch jetzt wieder zu beachten ist, dass für die 4 komplexen Wurzelwerte die Bedingung  $-\pi < \varphi \leq +\pi$  erfüllt sein muss.

Für ein vorgegebenes achsenparalleles Rechteckintervall  $Z \subset \mathbb{C}$  sollen jetzt für alle komplexen Zahlen  $z \in Z$  die Funktionswerte des Hauptzweiges der  $n$ -ten Wurzel durch ein achsenparalleles Rechteckintervall möglichst optimal eingeschlossen werden. Dazu definieren wir die komplexe Zahlenmenge  $\sqrt{Z, n}$  durch:

$$\sqrt{Z, n} := \begin{cases} [1, 1] + i \cdot [0, 0], & \text{für } n = 0 \text{ und für alle } Z \subset \mathbb{C}, \\ Z \subset \mathbb{C}, & \text{für } n = 1, \\ \square \{ \sqrt[n]{z} \mid z \in Z \}, & Z \subset \mathbb{C}_0^-, n = 2, 3, 4, \dots, \\ \text{undefiniert,} & Z \cap \mathbb{R}^- \neq \emptyset, n = 2, 3, 4, \dots \end{cases}$$

$\mathbb{C}_0^-$  ist die längs der negativen reellen Achse aufgeschnittene komplexe Ebene  $\mathbb{C}$ , wobei die Null selbst zu  $\mathbb{C}_0^-$  gehört.  $\mathbb{R}^-$  bedeutet die Menge der negativen reellen Zahlen, und  $\square \{ \sqrt[n]{z} \mid z \in Z \}$  symbolisiert die achsenparallele Rechteckeinschließung der komplexen Menge  $\{ \sqrt[n]{z} \mid z \in Z \}$ . Die in `mpfciclass.cpp` definierte Funktion

```
Mpfciclass sqrt( const Mpfciclass& z, int n )
```

liefert eine fast optimale Einschließung von  $\sqrt{Z, n}$ , d.h. es gilt

$$\sqrt{Z, n} \subseteq \text{sqrt}(Z, n).$$

In der folgenden Abbildung C.22 sind für  $n = 2, 3, 4, \dots$  einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben:

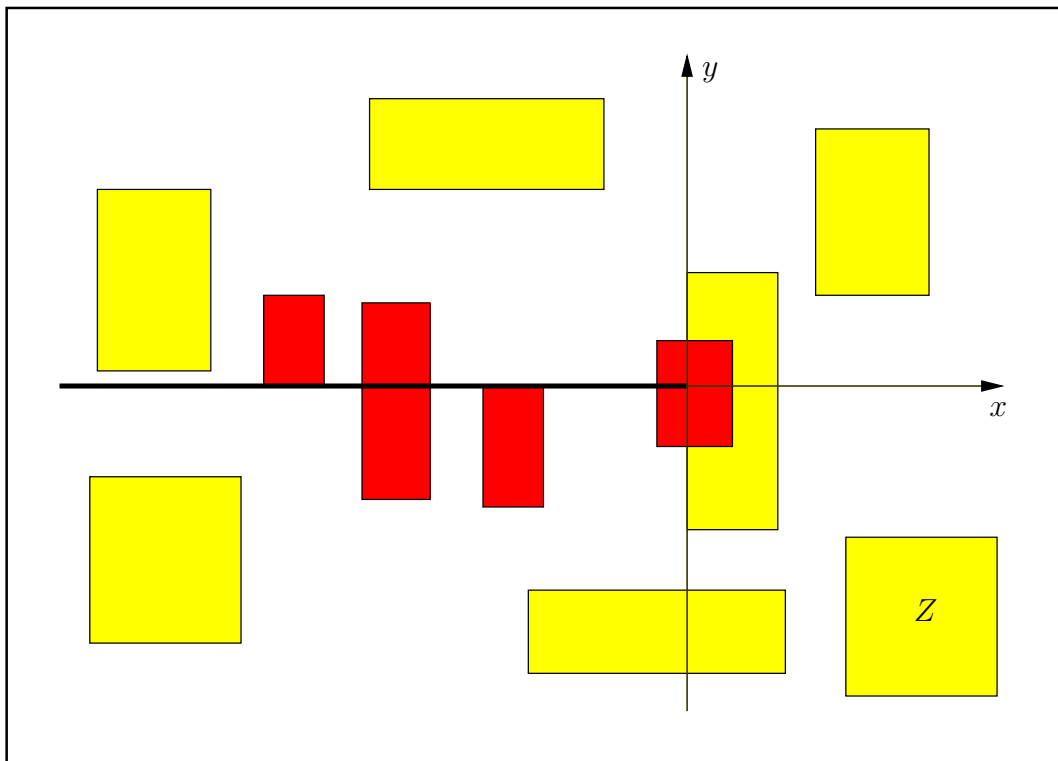


Abbildung C.22.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{sqrt}(Z, n)$ .

Es folgen noch einige numerische Beispiele in der niedrigen Präzision  $\text{prec} = 30$ , womit etwa  $30/\log_2(10) \approx 30/3.321928095 \approx 9$  Dezimalstellen berechnet werden:

```

sqrt([-2,+1] + i·[-1,+1],0) = ([1.00000000,1.00000000],[0.00000000,0.00000000])
sqrt([-2,+1] + i·[-1,+1],1) = ([-2.00000000,1.00000000],[-1.00000000,1.00000000]),
sqrt([+1,+2] + i·[-1,+0],2) = ([1.00000000,1.45534670],[-4.55089861e-1,-0.00000000]),
sqrt([-1,-1] + i·[+1,+1],3) = ([0.793700525,0.793700526],[0.793700525,0.793700526]),
sqrt([-2√3,-2√3] + i·[-2,-2],4) =
    ([1.12197105,1.12197106],[-8.60918670e-1,-8.60918669e-1]),
sqrt([-5,-5] + i·[-10,-10],9) = ([1.27439289,1.27439290],[-0.293084789,-0.293084788]),
sqrt([-2,+1] + i·[-1,+1],6) ~
MpfciClass sqrt(const MpfciClass& z, const int n ); z contains negative real values.

```

**Anmerkung:**

Der Parameter  $n$  sollte nicht zu groß gewählt werden.  $n = 100000$  erfordert beispielsweise schon große Laufzeiten!

## C.2.12. $\sqrt[n]{z}$ Alle Wurzeln

Für ein vorgegebenes achsenparalleles Rechteckintervall  $Z \subset \mathbb{C}$  berechnet die Funktion, [13], [55]

```
std::list<Mpfciclass> sqrt_all(const Mpfciclass& Z, int n);
```

eine Liste von  $n$  achsenparallelen Rechteckintervallen, die jeweils alle Lösungen von  $w^n = z$  für alle  $z \in Z \subset \mathbb{C}$  einschließen. Mit dem folgenden Programm MPFR-08.cpp berechnen wir zunächst für das Punktintervall  $Z = [-1, -1] + i \cdot [1, 1]$  Einschließungen der drei Wurzeln  $w_0, w_1, w_2$  der Gleichung  $w^3 = Z$  aus dem ersten Beispiels von Seite 302.

```
1 // MPFR-08.cpp
2 #include "mpfciclass.hpp"
3
4 using namespace MPFR;
5 using namespace MPFI;
6 using namespace cxsc;
7 using namespace std;
8
9 int main(void)
10 {
11     MPFI::MpfiClass::SetCurrPrecision(30);
12     cout << "GetCurrPrecision() = " << MPFI::MpfiClass::GetCurrPrecision() << endl;
13     int n = 3; // n-te Wurzel berechnen
14     interval re(-1,-1), im(1,1);
15     Mpfciclass Z(re, im, 53);
16     cout.precision(30/3.321928095); // Ausgabe mit 9 Dez.-Stellen
17     cout << "Z = " << Z << endl;
18     cout << "Z.GetPrecision() = " << Z.GetPrecision() << endl;
19     cout << "Berechnung aller " << n << "-ten Wurzeln aus Z" << endl;
20
21     list<Mpfciclass> res;
22     res = sqrt_all(Z, n);
23
24     list<Mpfciclass>::iterator pos;
25     // Ausgabe der n n-ten Wurzeln:
26     for (pos = res.begin(); pos != res.end(); ++pos)
27     {
28         cout << *pos << endl; // Jede Einschliessung in neue Zeile
29         cout << "Praezision = " << (*pos).GetPrecision() << endl;
30     }
31
32     return 0;
33 }
```

Das Programm liefert die Ausgabe

```
GetCurrPrecision() = 30
Z = ([-1.00000000, -1.00000000], [1.00000000, 1.00000000])
Z.GetPrecision() = 53
Berechnung aller 3-ten Wurzeln aus Z
([7.93700525e-1, 7.93700527e-1], [7.93700525e-1, 7.93700527e-1])
Praezision = 30
([-1.08421509, -1.08421508], [2.90514555e-1, 2.90514556e-1])
Praezision = 30
([2.90514555e-1, 2.90514556e-1], [-1.08421509, -1.08421508])
Praezision = 30
```

mit den Einschließungen der drei Wurzeln  $w_0, w_1, w_2$ .



Mit  $Z = [-2\sqrt{3}, -2\sqrt{3}] + i \cdot [-2, -2]$  liefert `sqrt_all(Z, 4)` Einschließungen der vier Wurzeln  $w_0, w_1, w_2, w_3$  der Gleichung  $w^4 = Z$  aus dem zweiten Beispiel von Seite 302. Die Einschließungen sind:

$$\begin{aligned} w_0 &\in ([1.12197105, 1.12197106], [-8.60918670e - 1, -8.60918668e - 1]), \\ w_1 &\in ([8.60918668e - 1, 8.60918670e - 1], [1.12197105, 1.12197106]), \\ w_2 &\in ([-1.12197106, -1.12197105], [8.60918668e - 1, 8.60918670e - 1]), \\ w_3 &\in ([-8.60918670e - 1, -8.60918668e - 1], [-1.12197106, -1.12197105]). \end{aligned}$$

Beachten Sie, dass die Intervalle  $Z \subset \mathbb{C}$  jetzt die negative reelle Achse nicht nur berühren sondern auch ganz in ihrem Innern enthalten dürfen. Mit  $Z = [-1, -1] + i \cdot [0, 0]$  erhält man daher für die beiden Quadratwurzeln  $w_0 = i, w_1 = -i$  die Einschließungen

$$\begin{aligned} w_0 &\in ([0.00000000, 0.00000000], [+1.00000000, +1.00000000]), \\ w_1 &\in ([0.00000000, 0.00000000], [-1.00000000, -1.00000000]). \end{aligned}$$

Mit  $Z = [-1, -1] + i \cdot [-1, 1]$  erhält man die folgenden sechs Einschließungen der sechs Wurzeln  $w_0, w_1, w_2, w_3, w_4, w_5$  der Gleichung  $w^6 = Z$ .

$$\begin{aligned} w_0 &\in ([7.93353340e - 1, 9.78816269e - 1], [3.82683432e - 1, 6.44960268e - 1]), \\ w_1 &\in ([-1.38287684e - 1, 1.38287684e - 1], [9.91444861e - 1, 1.05946310]), \\ w_2 &\in ([-9.78816269e - 1, -7.93353340e - 1], [3.82683432e - 1, 6.44960268e - 1]), \\ w_3 &\in ([-9.78816269e - 1, -7.93353340e - 1], [-6.44960268e - 1, -3.82683432e - 1]), \\ w_4 &\in ([-1.38287684e - 1, 1.38287684e - 1], [-1.05946310, -9.91444861e - 1]), \\ w_5 &\in ([7.93353340e - 1, 9.78816269e - 1], [-6.44960268e - 1, -3.82683432e - 1]). \end{aligned}$$

Mit  $Z = [-1, +1] + i \cdot [4, 4]$  erhält man die folgenden acht Einschließungen der acht Wurzeln  $w_0, w_1, w_2, w_3, w_4, w_5, w_6, w_7$  der Gleichung  $w^8 = Z$ .

$$\begin{aligned} w_0 &\in ([1.15870665, 1.17736602], [1.96183044e - 1, 2.68620905e - 1]), \\ w_1 &\in ([6.30104013e - 1, 6.93274511e - 1], [9.68097359e - 1, 1.01238336]), \\ w_2 &\in ([-2.68620905e - 1, -1.96183044e - 1], [1.15870665, 1.17736602]), \\ w_3 &\in ([-1.01238336, -9.68097359e - 1], [6.30104013e - 1, 6.93274511e - 1]), \\ w_4 &\in ([-1.17736602, -1.15870665], [-2.68620905e - 1, -1.96183044e - 1]), \\ w_5 &\in ([-6.93274511e - 1, -6.30104013e - 1], [-1.01238336, -9.68097359e - 1]), \\ w_6 &\in ([1.96183044e - 1, 2.68620905e - 1], [-1.17736602, -1.15870665]), \\ w_7 &\in ([9.68097359e - 1, 1.01238336], [-6.93274511e - 1, -6.30104013e - 1]). \end{aligned}$$

Mit  $Z = [+1, +1] + i \cdot [-0, +0]$  erhält man die folgenden fünf Einschließungen der fünf Einheitswurzeln  $w_0, w_1, w_2, w_3, w_4$  der Gleichung  $w^5 = 1$ .

$$\begin{aligned} w_0 &\in ([1.00000000, 1.00000000], [0.00000000, 0.00000000]), \\ w_1 &\in ([3.09016994e - 1, 3.09016995e - 1], [9.51056516e - 1, 9.51056517e - 1]), \\ w_2 &\in ([-8.09016995e - 1, -8.09016994e - 1], [5.87785252e - 1, 5.87785253e - 1]), \\ w_3 &\in ([-8.09016995e - 1, -8.09016994e - 1], [-5.87785253e - 1, -5.87785252e - 1]), \\ w_4 &\in ([3.09016994e - 1, 3.09016995e - 1], [-9.51056517e - 1, -9.51056516e - 1]). \end{aligned}$$

### C.2.13. $\sqrt{z+1}-1$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = \sqrt{z+1}-1, \quad z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind. Für  $f(z) = \sqrt{z+1}-1$  ist  $z_V = -1 + 0 \cdot i$  der Verzweigungspunkt, und der Verzweigungsschnitt ist die Menge  $\{r \in \mathbb{R} \mid -\infty < r < -1\}$ . Für den Hauptwert von  $f(z)$  darf  $Z$  den Verzweigungsschnitt nicht in seinem Inneren enthalten und auch nicht von unten berühren. In Abbildung C.21 sind einige erlaubte und nicht erlaubte achsenparallele Argumentintervalle  $Z$  der Funktion  $F(Z) = \sqrt{Z+1}-1 \supseteq W$  angegeben.

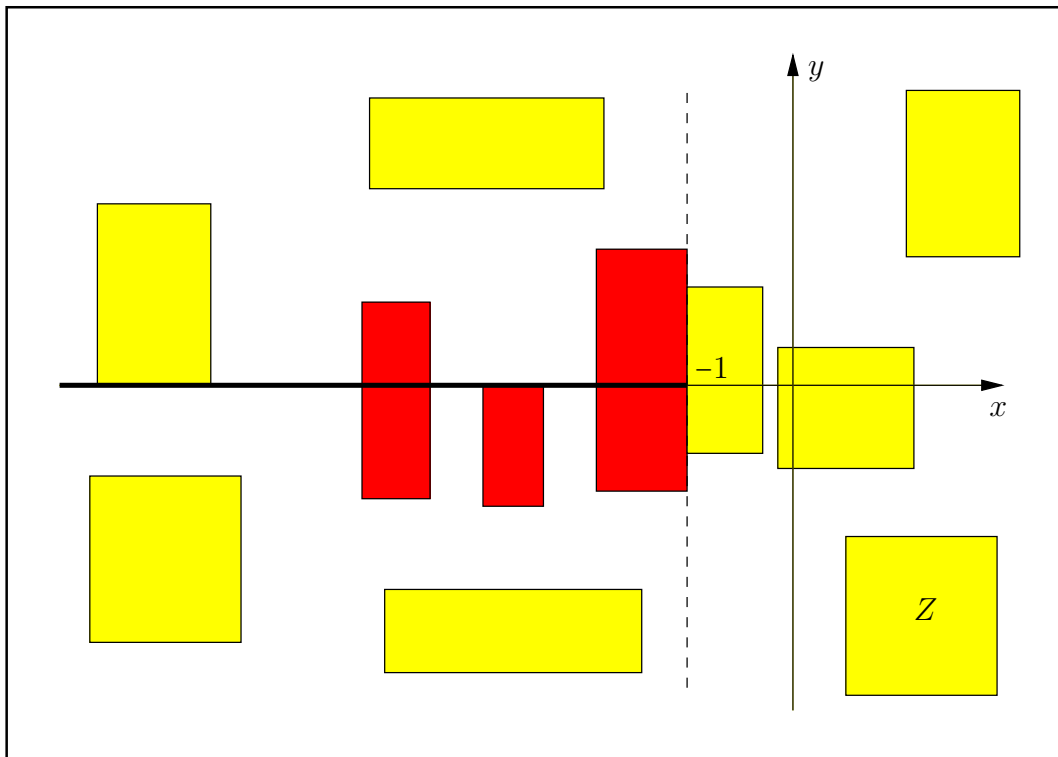


Abbildung C.23.: Erlaubte und nicht erlaubte Intervalle  $Z$  von  $\sqrt{Z+1}-1$ .

Zum vorgegebenen, erlaubten Eingangsintervall  $Z$  wird die nahezu optimale Einschließung  $F$  berechnet mithilfe der Funktion

```
Mpfciclass sqrtp1m1 (const Mpfciclass& Z);
```

die in der Datei `mpfciclass.cpp` definiert ist.

Um für  $F(Z) := \sqrt{Z+1}-1 = U + i \cdot V$  die reellen Intervalle  $U, V$  für die Einschließungen des Real- und Imaginärteils optimal berechnen zu können, benötigen wir jetzt für die Real- und Imaginärteilmfunktionen  $u(x, y), v(x, y)$  von  $f(z) = \sqrt{z+1}-1 = u + i \cdot v$  jeweils die Minimum- und Maximumpunkte  $m$  und  $M$  auf dem **Rand** des erlaubten Eingangsintervall  $Z$ , vgl. Seite 244. Um die Lage der Punkte  $m$  und  $M$  auf dem Rand von  $Z$  bestimmen zu können, müssen für  $u(x, y)$  und  $v(x, y)$  jeweils die Partiellen Ableitungen nach  $x$  und  $y$  berechnet und auf Minima und Maxima der Funktionen  $u(x, y)$  und  $v(x, y)$  jeweils auf dem Rand von  $Z$  untersucht werden.

### C.2.13.1. Realteil

Nach (C.40) gilt für die Realteilstfunktion  $u(x, y)$  von  $f(z) = \sqrt{z+1} - 1 = u(x, y) + i \cdot v(x, y)$  mit einem positiven Faktor  $\delta(x, y) > 0$

$$u(x, y) = \sqrt{\frac{|z+1| + x + 1}{2}} - 1, \quad z = x + i \cdot y \in \mathbb{C},$$

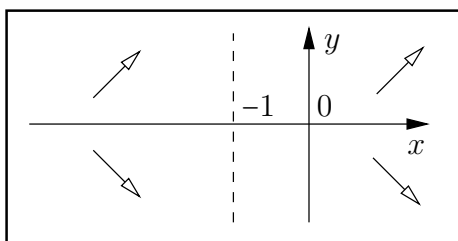
$$\frac{\partial u(x, y)}{\partial x} = \delta \cdot \frac{(x+1) + \sqrt{(x+1)^2 + y^2}}{\sqrt{(x+1)^2 + y^2}}, \quad x \neq -1,$$

$$\frac{\partial u(x, y)}{\partial y} = \delta \cdot \frac{y}{\sqrt{(x+1)^2 + y^2}}, \quad x \neq -1.$$

Im Fall  $z \neq -1 + 0 \cdot i$  gilt also  $\partial u / \partial x \geq 0$ , und  $\partial u / \partial x = 0$  ist nur möglich auf dem Verzweigungsschnitt, d.h für  $y = 0 \wedge -\infty < x < -1$ . Außerhalb des Verzweigungsschnitts gilt damit:

$$(C.64) \quad \text{sign}\left(\frac{\partial u}{\partial x}\right) = 1, \quad x \neq -1,$$

$$(C.65) \quad \text{sign}\left(\frac{\partial u}{\partial y}\right) = \text{sign}(y), \quad x \neq -1.$$



Nach (C.64) und (C.65) geben im linken Bild die vier Pfeile die Richtung wachsender Funktionswerte von  $u(x, y)$  an.

In der unteren Abbildung sind zu einigen Eingangsintervallen  $Z$  die Randpunkte  $m, M$  angegeben, in denen das Minimum bzw. Maximum von  $u(x, y)$  angenommen wird.

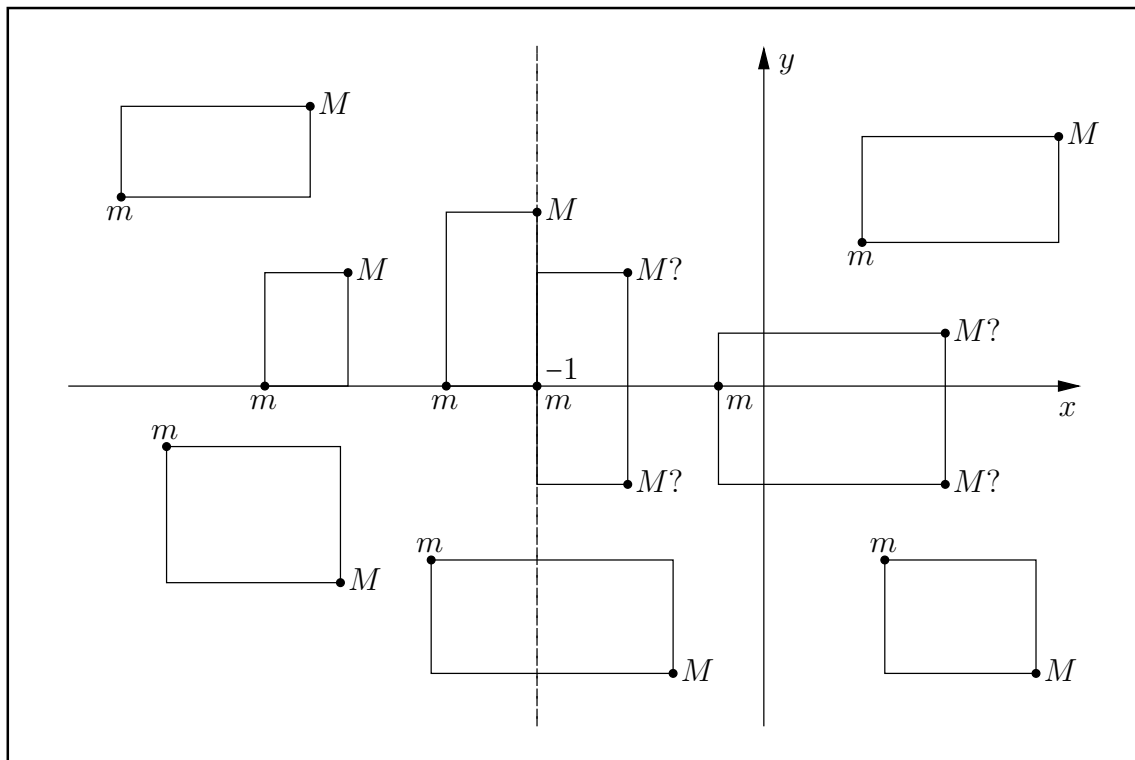


Abbildung C.24.: Minimum- und Maximumpunkte  $m, M$  der Realteilstfunktion  $u(x, y)$ .

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]; \quad M? = (x_2, y_1), \text{ falls } |y_1| > |y_2|, \text{ sonst } M? = (x_2, y_2).$$

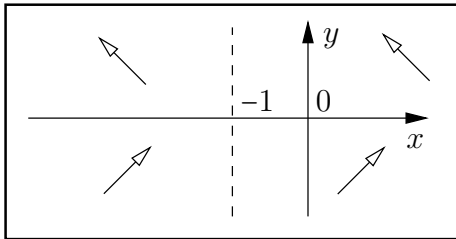
### C.2.13.2. Imaginärteil

Nach (C.42) auf Seite 274 gilt für die Imaginärteilstfunktion  $v(x, y)$  von  $f(z) = u(x, y) + i \cdot v(x, y)$  die allgemeine Darstellung

$$v(x, y) = \frac{y}{\sqrt{2 \cdot (|z+1| + x+1)}}, \quad y \neq 0, \quad \text{und einfache Rechnungen ergeben:}$$

$$(C.66) \quad \text{sign}\left(\frac{\partial v}{\partial x}\right) = -\text{sign}(y), \quad y \neq 0,$$

$$(C.67) \quad \text{sign}\left(\frac{\partial v}{\partial y}\right) = +1, \quad y \neq 0.$$



Nach (C.66) und (C.67) geben im linken Bild die vier Pfeile die Richtung wachsender Funktionswerte von  $v(x, y)$  an.

In der unteren Abbildung sind zu einigen Eingangsintervallen  $Z$  die Randpunkte  $m, M$  angegeben, in denen das Minimum bzw. Maximum des Imaginärteils  $v(x, y)$  angenommen wird.

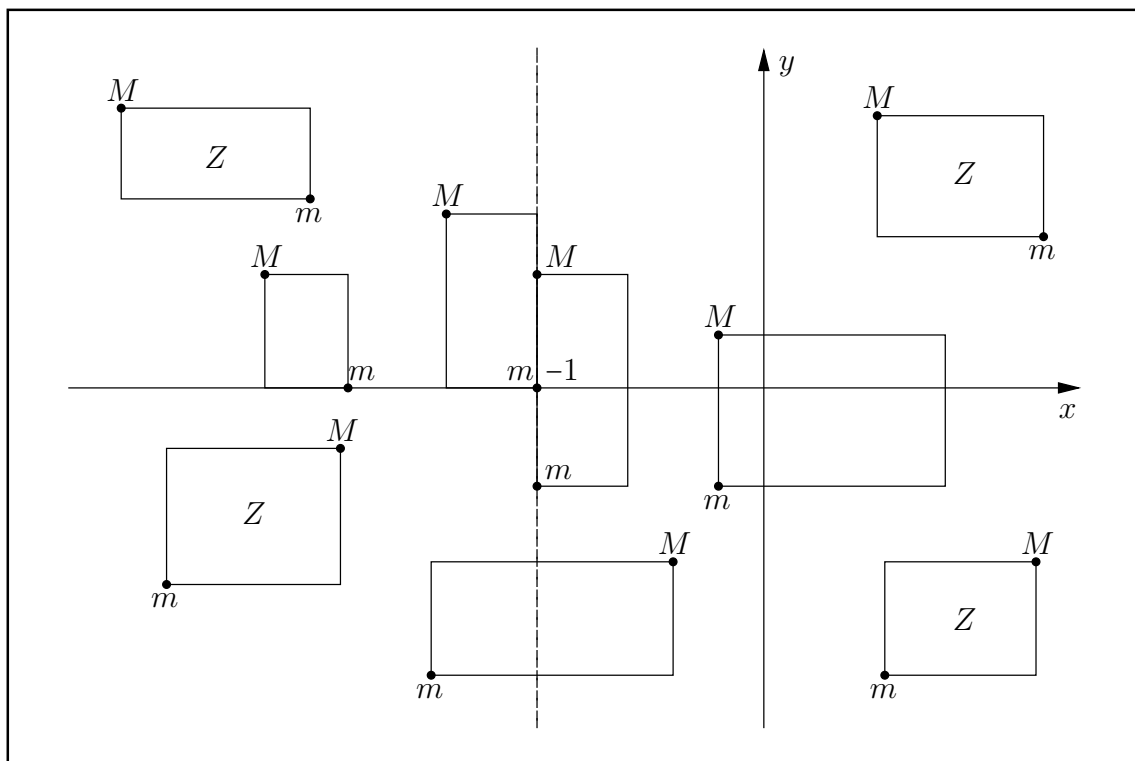


Abbildung C.25.: Minimum- und Maximumpunkte  $m, M$  der Imaginärteilstfunktion  $v(x, y)$ .

### C.2.13.3. Numerische Ergebnisse

Im **1. Beispiel** wählen wir mit  $Z = X + i \cdot Y = [x_1, x_2] + i \cdot [y_1, y_2]$  und mit der Current-Präzision  $\text{prec} = 800$  für  $Z$  das Punktintervall  $x_1 = x_2 = y_1 = y_2 = \text{MaxFloat}() = 2.0985 \dots \cdot 10^{+323228496}$ . Für den Wertebereich  $W$  von  $f(z)$ , mit  $z \in Z$ , erhalten wir die nahezu optimale Einschließung

$$W \subseteq F(Z) = ([1.5916056670864 \dots 8630 \cdot 10^{+161614248}, 1.5916056670864 \dots 8631 \cdot 10^{+161614248}], \\ [6.5926465325706 \dots 2455 \cdot 10^{+161614247}, 6.5926465325706 \dots 2457 \cdot 10^{+161614247}]).$$

Man erkennt die fast optimale Einschließung an den jeweils 240 gemeinsamen Dezimalstellen.

Im **2. Beispiel** wählen wir mit  $Z = X + i \cdot Y = [x_1, x_2] + i \cdot [y_1, y_2]$  und mit der Current-Präzision  $\text{prec} = 800$  für  $Z$  das Punktintervall  $x_1 = x_2 = 0.5$ ,  $y_1 = y_2 = 4$ . Nach Abbildung C.14 auf Seite 275 liegt dann  $Z$  im oberen rechten Eckpunkt des roten Bereiches. Für den Wertebereich  $W$  von  $f(z) = \sqrt{z+1} - 1$ , mit  $z \in Z$ , erhalten wir die nahezu optimale Einschließung

$$W \subseteq F(Z) = ([6.98823397628306387333...5497 \cdot 10^{-1}, 6.98823397628306387333...5498 \cdot 10^{-1}], \\ [1.1772854098855480133550677...3088, 1.1772854098855480133550677...3089]).$$

Man erkennt die fast optimale Einschließung an den jeweils 240 gemeinsamen Dezimalstellen.

Im **3. Beispiel** wählen wir mit  $Z = X + i \cdot Y = [x_1, x_2] + i \cdot [y_1, y_2]$  und mit der Current-Präzision  $\text{prec} = 800$  für  $Z$  das Punktintervall  $x_1 = x_2 = 1$ ,  $y_1 = y_2 = +6$ . Für den Wertebereich  $W$  von  $f(z) = \sqrt{z+1} - 1$ , mit  $z \in Z$ , erhalten wir die nahezu optimale Einschließung

$$W \subseteq F(Z) = ([1.0401660864175689291956325...8935, 1.0401660864175689291956325...8936], \\ [1.4704685172312868433025417...4107, 1.4704685172312868433025417...4108]).$$

Man erkennt die fast optimale Einschließung an den jeweils 240 gemeinsamen Dezimalstellen.

Im **4. Beispiel** wählen wir mit  $Z = X + i \cdot Y = [x_1, x_2] + i \cdot [y_1, y_2]$  und mit der Current-Präzision  $\text{prec} = 800$  für  $Z$  das echte Intervall  $Z = [0.5, 1] + i \cdot [4, 6]$ . Für den Wertebereich  $W$  von  $f(z) = \sqrt{z+1} - 1$ , mit  $z \in Z$ , erhalten wir die nahezu optimale Einschließung

$$W \subseteq F(Z) = ([6.98823397628306387333...5497 \cdot 10^{-1}, 1.040166086417568929195632...8936], \\ [1.1117859405028423439840960...8368, 1.530466993833334980217900...2684]).$$

Vergleichen Sie bitte die Realteileinschließung dieses Beispiels mit den Realteileinschließungen der Beispiele 2 und 3. Beachten Sie dabei die Punkte  $m, M$  des Eingangsintervalls oben rechts in Abb. C.24 auf Seite 307 und wählen Sie  $m(0.5|4)$  und  $M(1|6)$ . Die entsprechenden Intervallgrenzen werden dabei im 4. Beispiel genau bestätigt.

Im **5. Beispiel** soll für  $\alpha := \sqrt{2^{-200000}(1+i) + (1+2i)} - \sqrt{1+2i}$  mit der Current-Präzision  $\text{prec} = 10000$  eine möglichst optimale Einschließung berechnet werden. Mit der Umformung

$$\alpha = \sqrt{(1+2i) \cdot \left( \frac{2^{-200000}(1+i)}{1+2i} + 1 \right)} - \sqrt{1+2i} \\ = \sqrt{1+2i} \cdot \left( \sqrt{\frac{2^{-200000}(1+i)}{1+2i} + 1} - 1 \right)$$

kann die Funktion  $f(z) = \sqrt{z+1} - 1$  direkt benutzt werden, wenn für  $z := 2^{-200000}(1+i)/(1+2i)$  vorher eine optimale Einschließung berechnet wird. Damit erhält man für  $\alpha$  die nahezu optimale Einschließung:

$$\alpha \in ([+3.005996416218562137...5710 \cdot 10^{-60207}, +3.005996416218562137...5712 \cdot 10^{-60207}], \\ [-1.001998805406187379...1904 \cdot 10^{-60207}, -1.001998805406187379...1903 \cdot 10^{-60207}])$$

mit 3009 von 3010 möglichen Dezimalstellen. Noch eine Anmerkung zur obigen Umformung: Mit  $u = |u| \cdot e^{i\varphi_u}$ ,  $v = |v| \cdot e^{i\varphi_v}$  und  $-\pi < \varphi_u, \varphi_v \leq +\pi$  gilt  $\sqrt{u \cdot v} = \sqrt{|u|} \cdot \sqrt{|v|}$  nur, falls  $-\pi < \varphi_u + \varphi_v \leq +\pi$ , wobei diese Voraussetzung bei der obigen Umformung von  $\alpha$  erfüllt ist.

Als Gegenbeispiel betrachten wir  $u = v = -1 = e^{i\varphi_u}$  mit  $\varphi_u = \varphi_v = \pi$  und erhalten wegen  $\varphi_u + \varphi_v = 2\pi > \pi$  die Ungleichung  $1 = \sqrt{(-1) \cdot (-1)} \neq \sqrt{-1} \cdot \sqrt{-1} = i \cdot i = -1$ .

### C.2.14. $\sqrt{1+z^2}$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = \sqrt{1+z^2}, \quad z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind. Für  $f(z) = \sqrt{1+z^2}$  sind  $z_{\pm 1} = 0 \pm i$  die Verzweigungspunkte, und die Verzweigungsschnitte sind die Mengen  $\{\pm r \cdot i \mid 1 \leq r < +\infty\}$ . Für den Hauptwert von  $f(z)$  darf  $Z$  einen Verzweigungsschnitt nicht in seinem Inneren enthalten. In der oberen Halbebene darf  $Z$  den Verzweigungsschnitt nur von rechts und in der unteren Halbebene nur von links berühren. In Abbildung C.26 sind einige erlaubte und nicht erlaubte achsenparallele Argumentintervalle  $Z$  der Funktion  $F(Z) = \sqrt{1+Z^2} \supseteq W$  angegeben.

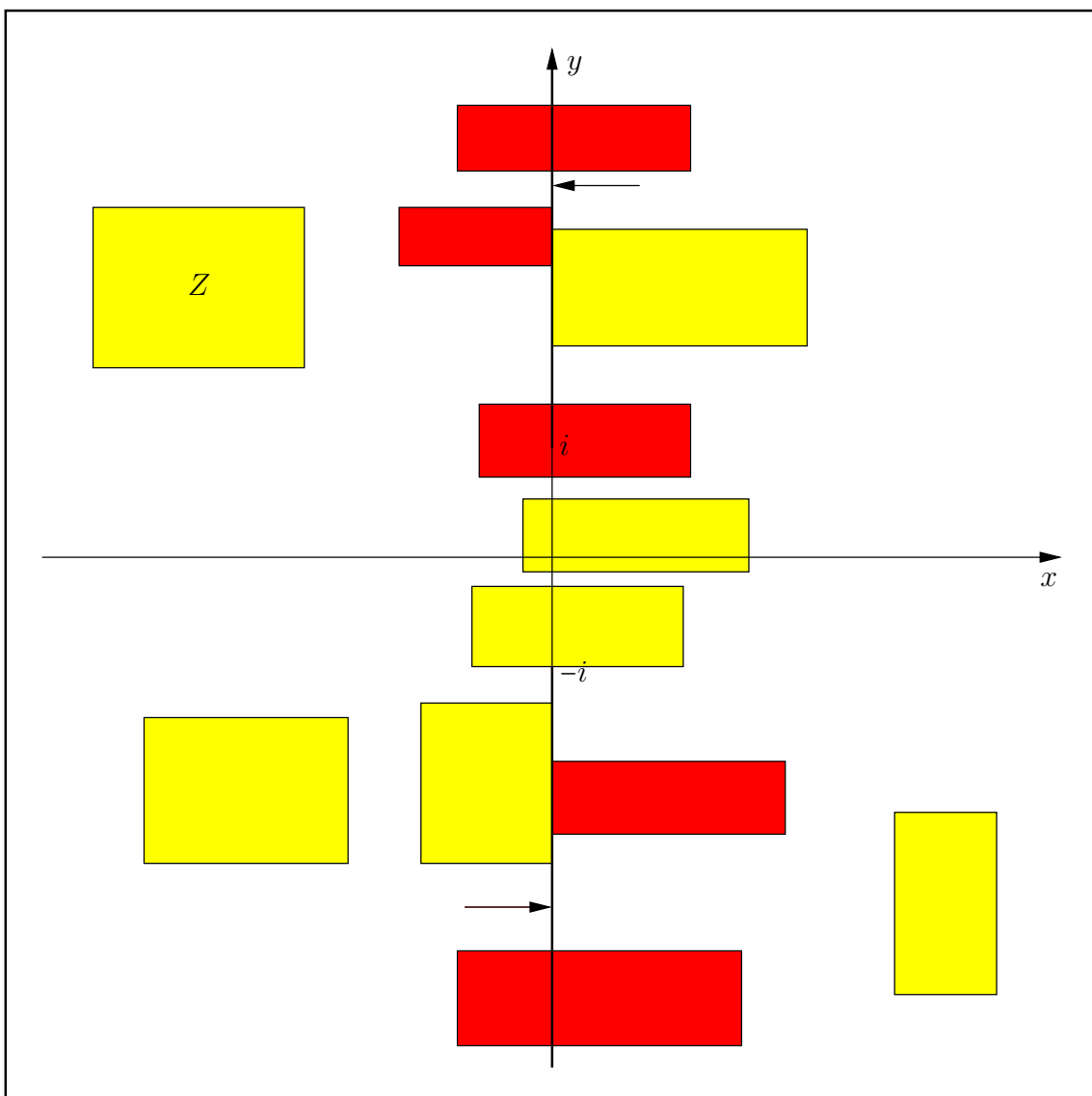


Abbildung C.26.: Erlaubte und nicht erlaubte Intervalle  $Z$  für den Hauptzweig von  $f(z)$

Zu einem vorgegebenen Intervall  $Z$  wird die optimale Einschließung  $F$  berechnet mit Hilfe der Funktion

```
Mpfciclass sqrt1px2(const Mpfciclass& Z);
```

die in `mpfciclass.cpp` definiert ist. Da die Funktion  $f(z) = \sqrt{1+z^2} = u(x,y) + i \cdot v(x,y)$  für  $z \in Z$  holomorph ist, nehmen die Real- und Imaginärteildfunktionen  $u(x,y)$ ,  $v(x,y)$  als harmonische Funktionen ihre Extrema jeweils auf dem Rand von  $Z = U + i \cdot V$  an, vgl. Seite 244.

### C.2.14.1. Realteil

Nach (C.45) auf Seite 279 gilt für die Realteildfunktion  $u(\pm x, \pm y) = u(x,y)$ , so dass man für die Berechnung der Extrema  $m, M$  nur die Intervalle  $X^* := \text{abs}(X)$  und  $Y^* := \text{abs}(Y)$  der jeweiligen Absolutbeträge zu betrachten hat. Das komplexe Intervall  $Z^* = X^* + i \cdot Y^*$  liegt damit nur im ersten Quadranten und besitzt für  $u(|x|, |y|)$  die gleichen Extrema  $m, M$  wie das ursprüngliche Intervall  $Z$  für  $u(x,y)$ . Im ersten Quadranten erhält man für die beiden partiellen Ableitungen der Realteildfunktion  $u(x,y)$  die Ergebnisse:

$$\frac{\partial u(x,y)}{\partial x} = \begin{cases} 0, & x = 0 \wedge 0 \leq y < 1, \\ > 0 & \text{sonst.} \end{cases} \quad \frac{\partial u(x,y)}{\partial y} = \begin{cases} 0, & y = 0 \vee (x = 0 \wedge y \geq +1), \\ < 0 & \text{sonst.} \end{cases}$$

Damit ergibt sich in Abb. C.27 für drei typische komplexe Intervalle  $Z^*$  aus dem 1. Quadranten die Lage der Punkte mit den gesuchte Extrema  $m, M$ .

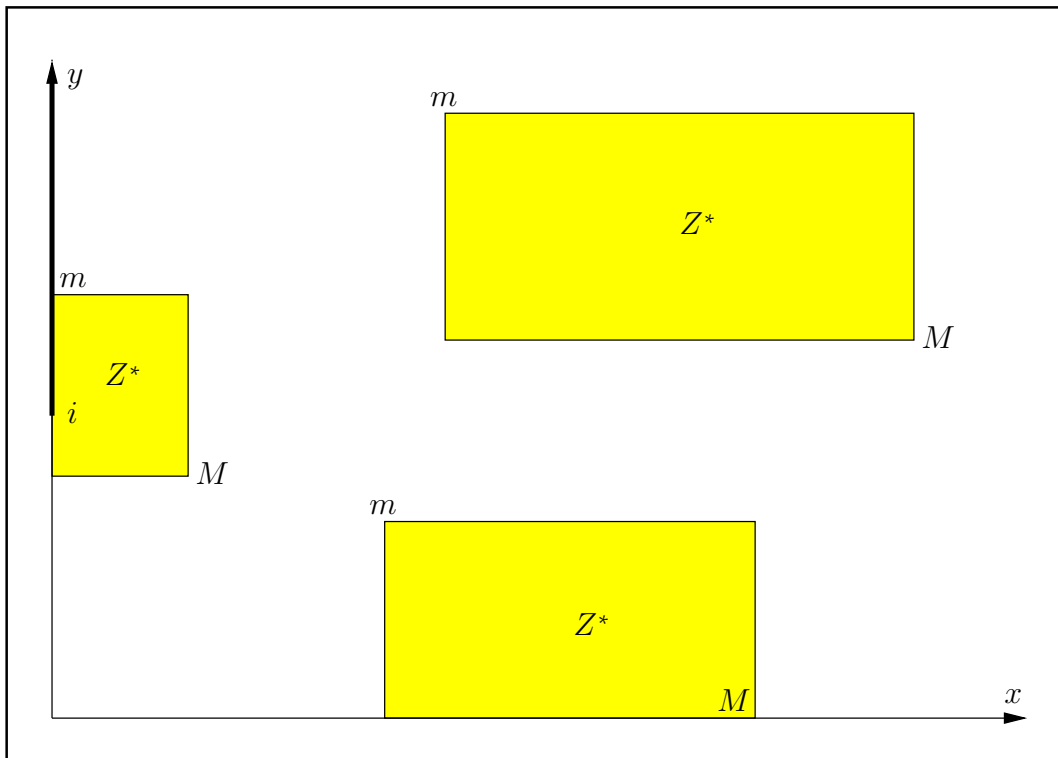


Abbildung C.27.: Lage der Extrema  $m, M$  der Realteildfunktion  $u(x,y)$  auf Intervallen  $Z^*$

Mit Hilfe der Koordinaten  $x_m, y_m$  bzw.  $x_M, y_M$  lassen sich dann das Minimum  $m$  bzw. das Maximum  $M$  von  $u(x,y)$  für ein Intervall  $Z^*$  aus dem 1. Quadranten berechnen mit der Funktion

```
MpfrClass Re_sqrt1px2(const MpfrClass& x, const MpfrClass& y, RoundingMode rnd)
```

Eine Unterschranke  $m$  des Minimums  $m \geq m$  und eine Oberschranke  $M$  des Maximums  $M \leq M$  erhält man also mit den Funktionsaufrufen

```
m = Re_sqrt1px2(xm, ym, RoundDown);    M = Re_sqrt1px2(xM, yM, RoundUp);
```

### C.2.14.2. Imaginärteil

Nach (C.57) von Seite 282 gilt für die Imaginärteilmfunktion

$$(C.68) \quad v(-x, y) = -v(x, y) \quad \text{und} \quad v(x, -y) = -v(x, y)$$

so dass über einem gegebenen Rechteck  $Z = X + i \cdot Y$ , mit  $\text{Sup}(X) \leq 0$ , die Imaginärteilmfunktion  $v(x, y)$  die gleichen Funktionswerte liefert wie über dem Rechteck  $Z^* := -X + i \cdot (-Y) = -Z$ , das dann jedoch nur noch in der rechten Halbebene liegt. Um die gesuchten Extrema  $m, M$  über  $Z$  zu bestimmen, müssen neben den Intervallen  $Z^*$  aus der rechten Halbebene wegen der auf der imaginären Achse liegenden Verzweigungsschnitte nur noch komplexe Intervalle  $Z = X + i \cdot Y$  betrachtet werden, für die mit  $X = [x_1, x_2]$  gilt:  $0 \in (x_1, x_2)$  und  $|y| < 1 \forall y \in Y$ .

Zur Bestimmung der Koordinaten der Extremstellen benötigen wir die partiellen Ableitungen

$$\text{sign}\left(\frac{\partial v(x, y)}{\partial x}\right) = \text{sign}(y), \quad \frac{\partial v(x, y)}{\partial y} = \begin{cases} 0, & x = 0 \wedge |y| \leq 1, \\ > 0, & x > 0 \vee (x = 0 \wedge |y| > 1), \\ < 0, & x < 0 \vee (x = 0 \wedge |y| > 1). \end{cases}$$

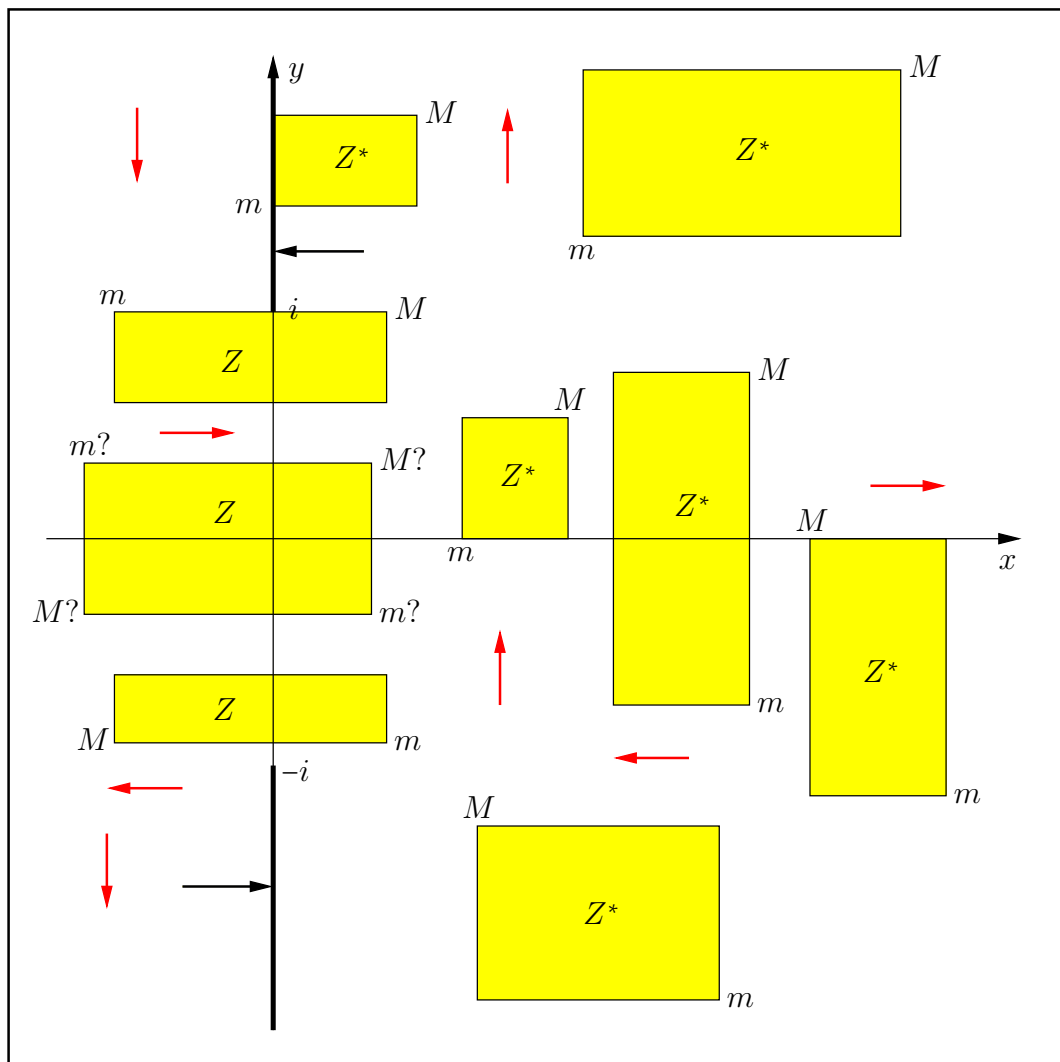


Abbildung C.28.: Extremstellen auf den Intervallrändern von  $Z$  bzw.  $Z^*$

Die beiden schwarzen Pfeile geben die Richtungen an, aus denen die Funktionswerte  $v(x, y)$  auf den jeweiligen Verzweigungsschnitt stetig ergänzt werden. Die roten Pfeile geben die Richtung wachsender Funktionswerte von  $v(x, y)$  an.



Beachten Sie bitte, dass in Abb. C.28 mit  $Z^*$  neben den am Ursprung gespiegelten Intervallen aus der linken Halbebene natürlich auch solche Intervalle symbolisiert werden, die ursprünglich in der rechten Halbebene vorgegeben waren.

### C.2.14.3. Numerische Ergebnisse

Im **1. Beispiel** wählen wir  $Z = [0, \text{minfloat}()) + i \cdot [1, 1] = [1, 2^{-1073741824}] + i \cdot [1, 1]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$ , mit  $z \in Z$ , die nahezu optimale Einschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 4.881152430408...027 \cdot 10^{-161614249}], [0, 4.881152430408...027 \cdot 10^{-161614249}]).$$

Im **2. Beispiel** wählen wir  $Z = [-2, +1] + i \cdot [-1, +0.5]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 2.23606797749978969640917366873127623544061835961152572427...8351], [-7.86151377757423286069...1307 \cdot 10^{-1}, 9.10179721124454682608...9583 \cdot 10^{-1}]).$$

Im **3. Beispiel** wählen wir  $Z = [-2, -1] + i \cdot [-4, +2]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([+1.03065889036991751569220...6366, +2.236067977499789696409173...8351] [-1.879129818333282376255771...9882, 3.899774854972245632918720...7245]).$$

Zeigen Sie mithilfe der Symmetrieeigenschaften von  $u(x, y)$  und  $v(x, y)$ , dass mit dem komplexen Eingangsintervall  $Z_1 = [1, 2] + i \cdot [-2, 4]$  die gleiche obige Einschließung  $F(Z)$  berechnet wird.

Im **4. Beispiel** wählen wir  $Z = [2^{+1073741000}, 2^{+1073741820}] + i \cdot [2^{+1073741000}, 2^{+1073741820}]$ . Mit  $\text{prec} = 400$  und

$$Z = [9.379481717334163189...964 \cdot 10^{323228247}, 6.558058488960586538...127 \cdot 10^{323228494}] [9.379481717334163189...964 \cdot 10^{323228247}, 6.558058488960586538...127 \cdot 10^{323228494}]$$

erhalten wir für den Wertebereich  $W$  von  $f(z) = \sqrt{1+z^2}$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = [9.3794817173341631...963 \cdot 10^{323228247}, 6.5580584889605865...128 \cdot 10^{323228494}] [9.3794817173341631...963 \cdot 10^{323228247}, 6.5580584889605865...127 \cdot 10^{323228494}].$$

Beachten Sie, dass bei der naiven Intervallauswertung von  $Z^2$  ein vorzeitiger Überlauf eintreten würde und dass wegen  $\text{Inf}(\Re(Z)) \gg 1$  und  $\text{Inf}(\Im(Z)) \gg 1$  die Intervalle  $Z$  und  $F(Z)$  nahezu identisch sind.

#### Anmerkung:

Es sei nochmals betont, dass die Funktion  $\text{sqrt1px2}(Z)$  eine **optimale** Einschließung  $F(Z)$  der Wertemenge  $W$  von  $f(z) = \sqrt{1+z^2}$  für alle  $z \in Z$  berechnet. Wenn man jedoch, z.B. mit dem Intervall  $Z = [-2, +1] + i \cdot [-1, 0.5]$  aus dem 2. Beispiel, den Funktionsterm  $\sqrt{1+z^2}$  in naiverweise intervallmäßig auswertet, so erhält man mit dem Ergebnis

$$\sqrt{1 \oplus Z^2} = ([0, 2.38779440461619817894179314811304091311811195514305331557...6078], [-1.000...001, 1.41421356237309504880168872420969807856967187537...9703]).$$

eine deutliche Überschätzung der optimalen Einschließung  $F(Z)$  aus dem 2. Beispiel. Der obige linke Ausdruck  $\sqrt{1 \oplus Z^2}$  soll die naive intervallmäßige Auswertung symbolisieren.

### C.2.15. $\sqrt{1-z^2}$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = \sqrt{1-z^2}, z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind. Für  $f(z) = \sqrt{1-z^2}$  sind  $z_{\pm 1} = \pm 1$  die Verzweigungspunkte, und die Verzweigungsschnitte sind die Mengen  $\{\pm r \mid 1 \leq r < +\infty\}$ . Für den Hauptwert von  $f(z)$  darf  $Z$  einen Verzweigungsschnitt nicht in seinem Inneren enthalten. In der rechten Halbebene darf  $Z$  den Verzweigungsschnitt nur von unten und in der linken Halbebene nur von oben berühren. In Abbildung C.29 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  der Funktion  $F(Z) = \sqrt{1-Z^2} \supseteq W$  angegeben.

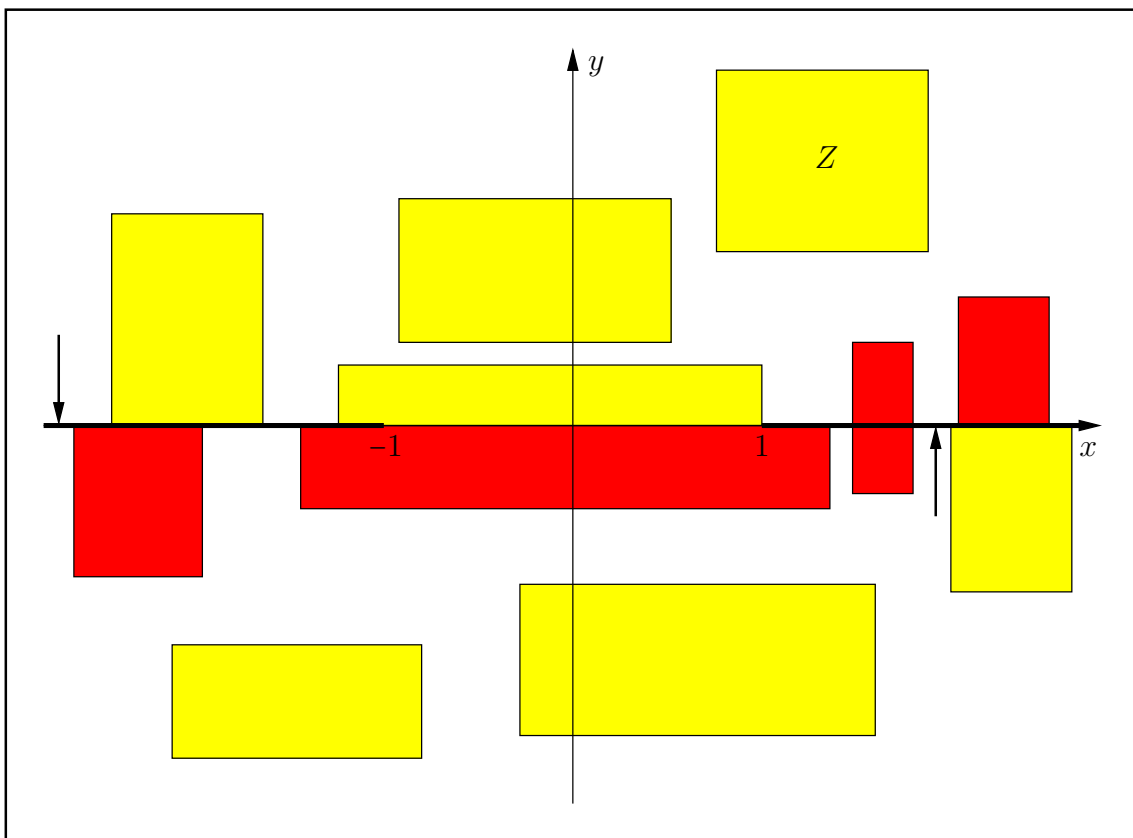


Abbildung C.29.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  für den Hauptzweig von  $f(z)$

Ganz analog zu Seite 283 erhält man mit der Transformation  $Z = i \cdot Z_T$ , bzw. mit  $Z_T = Y + i \cdot (-X)$  die Beziehung  $\sqrt{1-Z^2} = \sqrt{1+Z_T^2}$ , so dass die Auswertung von  $\sqrt{1-Z^2}$  zurückgeführt werden kann auf die Auswertung von  $\sqrt{1+Z_T^2}$ .

Die **optimale** Einschließung des Wertebereichs  $W$  durch das komplexe Rechteck  $F(Z) \supseteq W$  erfolgt mithilfe der Funktion

```
Mpfciclass sqrt1mx2(const Mpfciclass& Z);
```

die in der Datei `pmfciclass.cpp` definiert ist.

### C.2.16. $\sqrt{z^2 - 1}$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = \sqrt{z^2 - 1}, \quad z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind. Für  $f(z) = \sqrt{z^2 - 1}$  sind  $z_{\pm 1} = \pm 1$  die Verzweigungspunkte. Die beiden Verzweigungsschnitte sind in der folgenden Abbildung C.30 schwarz bzw. grün dargestellt, vgl. Seite 284. Typische erlaubte und nicht erlaubte komplexe Intervalle  $Z$  sind entsprechend gelb bzw. rot dargestellt.

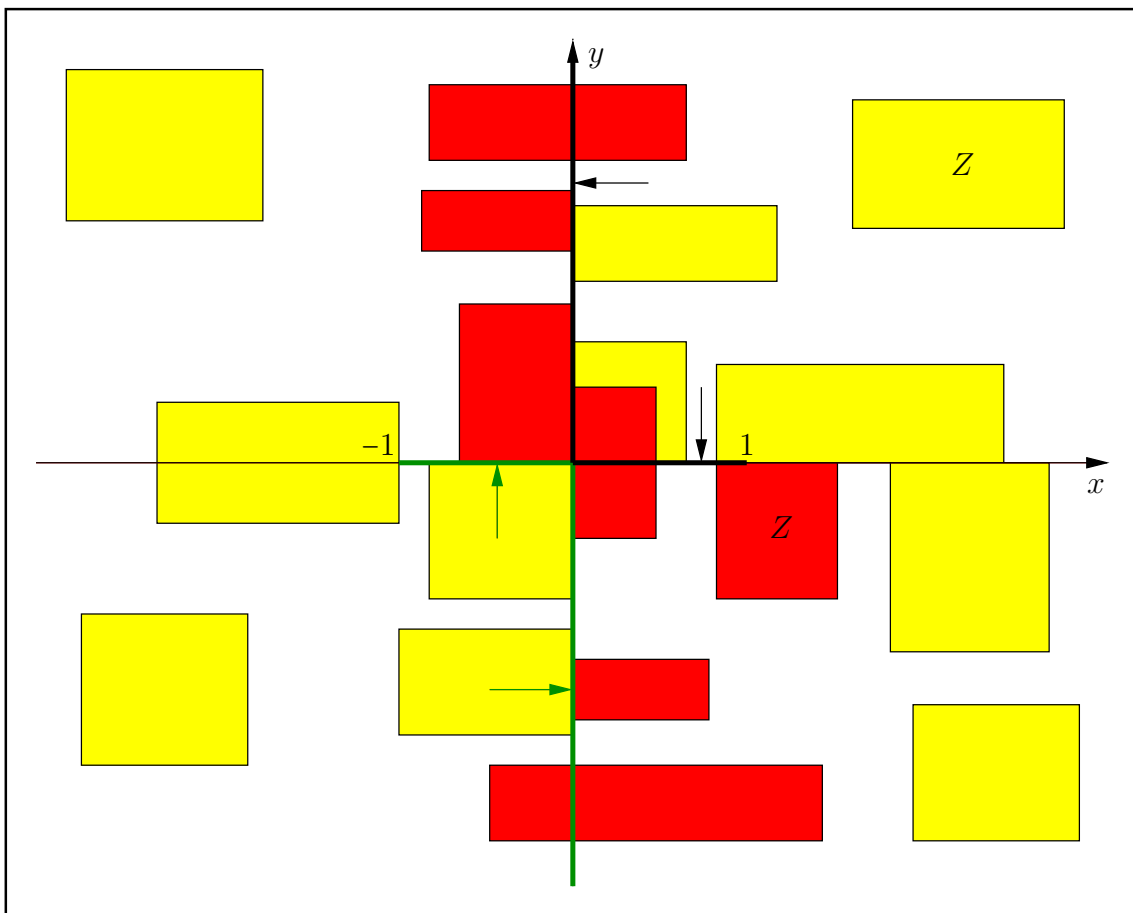


Abbildung C.30.: Erlaubte und nicht erlaubte Intervalle  $Z$  für den Hauptzweig von  $f(z)$

Beachten Sie, dass auch jedes Intervall der Form  $Z = [0, 0] + i \cdot Y$  oder  $Z = X + i \cdot [0, 0]$ , d.h. jedes Intervall  $Z$ , das auf einer Koordinatenachse liegt, ein erlaubtes Intervall ist. Die schwarzen und die grünen Pfeile geben die Richtungen an, aus denen die Funktionswerte von  $f(z) = \sqrt{z^2 - 1}$ , mit  $z \in \mathbb{C}$ , auf die jeweiligen Verzweigungsschnitte analytisch fortgesetzt werden.

Die **optimale** Einschließung des Wertebereichs  $W$  durch das komplexe Rechteck  $F(Z) \supseteq W$  erfolgt mithilfe der Funktion

```
Mpfciclass sqrtx2m1(const Mpfciclass& Z);
```

die in der Datei `pmfciclass.cpp` definiert ist.

### C.2.16.1. Realteil

Nach (C.59) auf Seite 284 gilt für die Realteilfunktion  $\hat{u}(\pm x, \pm y) = \hat{u}(x, y)$ , so dass man für die Berechnung der Extrema  $m, M$  nur die Intervalle  $X^* := \text{abs}(X)$  und  $Y^* := \text{abs}(Y)$  der jeweiligen Absolutbeträge zu betrachten hat. Das komplexe Intervall  $Z^* = X^* + i \cdot Y^*$  liegt damit nur im ersten Quadranten und besitzt für  $\hat{u}(|x|, |y|)$  die gleichen Extrema  $m, M$  wie das ursprüngliche Intervall  $Z$  für  $\hat{u}(x, y)$ . Für die partiellen Ableitungen erhält man die Ergebnisse

$$\begin{aligned} \frac{\partial \hat{u}(x, y)}{\partial x} &> 0, \quad \text{falls } y > 0, \\ \frac{\partial \hat{u}(x, y)}{\partial x} &= 0, \quad \text{falls } y = 0 \wedge 0 \leq x < 1, \\ \frac{\partial \hat{u}(x, y)}{\partial x} &= \frac{x}{\sqrt{x^2 - 1}} > 0, \quad \text{falls } y = 0 \wedge x > 1; \end{aligned}$$

$$\begin{aligned} \frac{\partial \hat{u}(x, y)}{\partial y} &= 0, \quad \text{falls } x = 0, \\ \frac{\partial \hat{u}(x, y)}{\partial y} &> 0, \quad \text{falls } x > 0. \end{aligned}$$

Damit erhält man in Abb. C.31 für einige typische Intervalle  $Z^*$  aus dem 1. Quadranten die Lage der Punkte mit den gesuchte Extrema  $m, M$  der Realteilfunktion  $\hat{u}(x, y)$ .

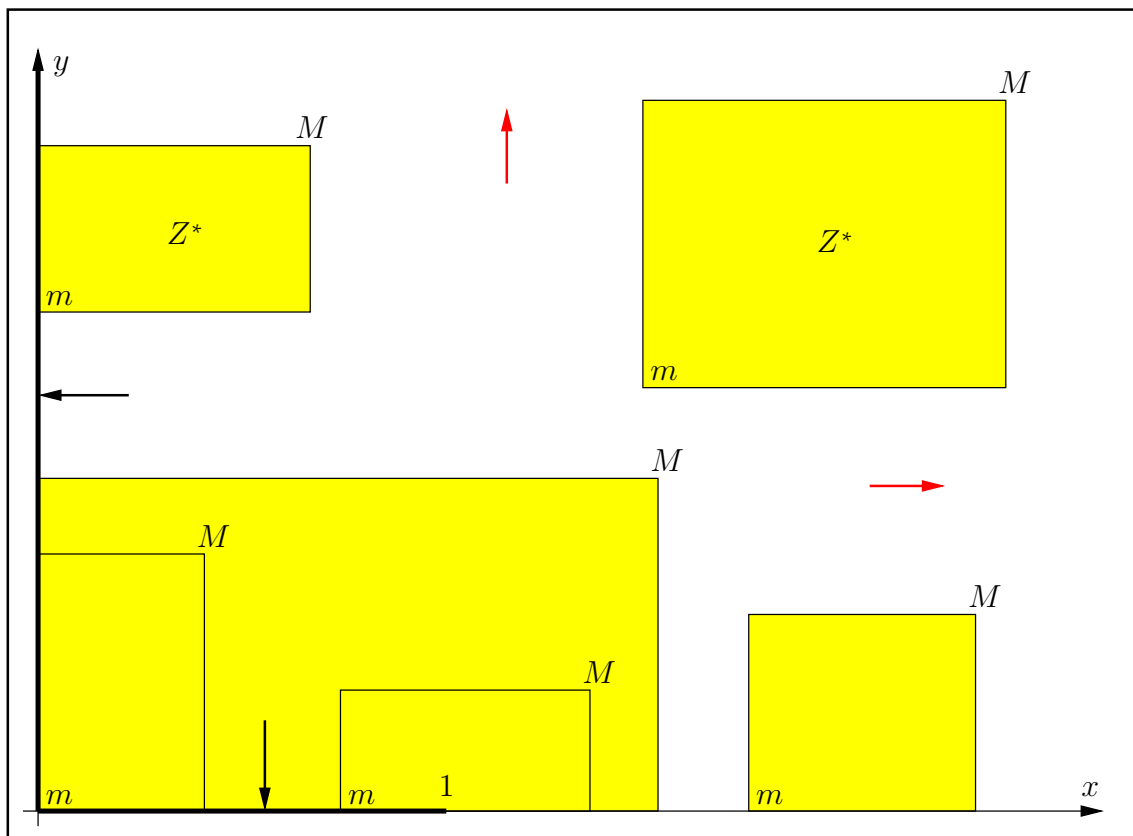


Abbildung C.31.: Extremstellen  $m, M$  von  $\hat{u}(x, y)$  auf den Intervallrändern von  $Z$  bzw.  $Z^*$

Die beiden schwarzen Pfeile geben die Richtungen an, aus der die Funktionswerte  $\hat{u}(x, y)$  auf den Verzweigungsschnitt stetig ergänzt werden. Die roten Pfeile geben die Richtung wachsender Funktionswerte von  $\hat{u}(x, y)$  an.

### C.2.16.2. Imaginärteil

Mit Ausnahme der Intervalle  $Z$ , die auf der reellen Achse liegen, befinden sich nach Abb. C.30 alle restlichen erlaubten Intervalle  $Z$  entweder in der rechten oder linken Halbebene. Im Falle  $y_1 = y_2 = 0$  liegen dann nach der Transformation  $Z^* := \text{abs}(X) + i \cdot [0, 0]$  **alle** erlaubten Intervalle  $Z$  bzw.  $Z^*$  entweder in der rechten oder linken Halbebene. Im Fall  $x_1 = x_2 = 0$  wählen wir die Transformation  $Z^* := [0, 0] + i \cdot \text{abs}(Y)$ . Wählt man dann noch im Fall  $x_1 < 0$  die Transformation  $Z^* := -X - i \cdot Y = -Z$ , so befinden sich anschließend wegen  $\hat{v}(-x, y) = -\hat{v}(x, y)$  und  $\hat{v}(x, -y) = -\hat{v}(x, y)$  alle erlaubten Intervalle  $Z, Z^*$  nur noch in der rechten Halbebene, wobei die Wertemengen der Funktion  $\hat{v}(x, y)$  über den Intervallen  $Z^*$  und deren Urbildern  $Z$  jeweils gleich sind, vgl. (C.60) auf Seite 284. Für die partiellen Ableitungen erhält man in der rechten Halbebene

$$\begin{aligned} \frac{\partial \hat{v}(x, y)}{\partial x} &= 0, & y = 0 \wedge x > 1, \\ \frac{\partial \hat{v}(x, y)}{\partial x} &< 0, & y = 0 \wedge 0 < x < 1, \\ \text{sign}\left(\frac{\partial \hat{v}(x, y)}{\partial x}\right) &= -\text{sign}(y), & y \neq 0 \\ \frac{\partial \hat{v}(x, y)}{\partial y} &= 0, & y = 0 \wedge 0 \leq x < 1, \\ \frac{\partial \hat{v}(x, y)}{\partial y} &> 0, & \text{sonst.} \end{aligned}$$

Damit erhält man in Abb. C.32 für einige typische Intervalle  $Z^*$  aus der rechten Halbebene die Lage der Punkte mit den gesuchten Extrema  $m, M$  der Imaginärteilmfunktion  $\hat{v}(x, y)$ .

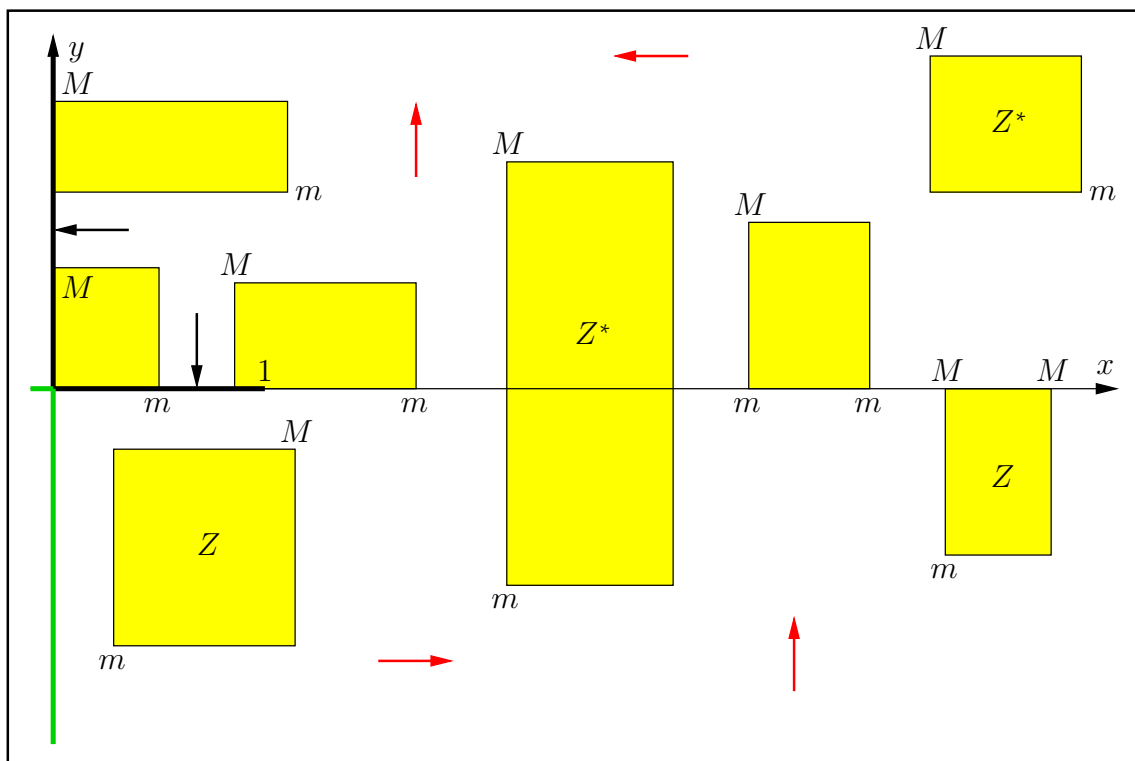


Abbildung C.32.: Extremstellen  $m, M$  von  $\hat{v}(x, y)$  auf den Intervallrändern von  $Z$  bzw.  $Z^*$

Die roten Pfeile geben die Richtung wachsender Funktionswerte von  $\hat{v}(x, y)$  an. Auf der reellen Achse ist  $\hat{v}(x, 0)$  für  $x \geq 1$  konstant, und für  $0 \leq x < 1$  wächst  $\hat{v}(x, 0)$  in Richtung zum Ursprung.

### C.2.16.3. Numerische Ergebnisse

Im **1. Beispiel** liegt  $Z = [-3, 2] + i \cdot [0, 0]$  auf der reellen Achse. Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$ , mit  $z \in Z$ , die nahezu optimale Einschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 2.8284271247461900976033774484193961571393437507538\dots9406], [0, 1]).$$

Beachten Sie, dass man nach (C.59) und (C.60) auf Seite 284 z.B. für  $Z_1 = [0, 3] + i \cdot [0, 0]$  die gleiche obige Einschließung  $F(Z) = F(Z_1)$  erhält.

Im **2. Beispiel** liegt  $Z = [0, 0] + i \cdot [-3, 2]$  auf der imaginären Achse. Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$ , mit  $z \in Z$ , die nahezu optimale Einschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 0], [1, 3.1622776601683793319988935444327185337195551393252168268]).$$

Beachten Sie, dass man nach (C.59) und (C.60) auf Seite 284 z.B. für  $Z_1 = [0, 0] + i \cdot [0, 3]$  die gleiche obige Einschließung  $F(Z) = F(Z_1)$  erhält.

Im **3. Beispiel** wählen wir  $Z = [-2, -1] + i \cdot [-4, +2]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 1.9513492931928285028764303881123826661836623638977857854\dots1558] \\ [-2.197368226935619932079602\dots4282, 4.116342054542984500643962\dots7563]).$$

Zeigen Sie mithilfe der Symmetrieeigenschaften von  $\hat{u}(x, y)$  und  $\hat{v}(x, y)$ , dass mit dem komplexen Eingangsintervall  $Z_1 = [1, 2] + i \cdot [-2, 4]$  die gleiche obige Einschließung  $F(Z)$  berechnet wird.

Im **4. Beispiel** wählen wir  $Z = [-1, 0] + i \cdot [-2, 0]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 9.10179721124454682608\dots9583 \cdot 10^{-1}], [0, 2.23606797749978969640\dots8351]).$$

Zeigen Sie mithilfe der Symmetrieeigenschaften von  $\hat{u}(x, y)$  und  $\hat{v}(x, y)$ , dass mit dem komplexen Eingangsintervall  $Z_1 = [0, 1] + i \cdot [0, 2]$  die gleiche obige Einschließung  $F(Z)$  berechnet wird.

Im **5. Beispiel** wählen wir  $Z = [0.5, 1.5] + i \cdot [0, 1]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0, 1.276792591750530535249\dots3706], [0, 1.37214511569925383987424\dots0792]).$$

Zeigen Sie mithilfe der Symmetrieeigenschaften von  $\hat{u}(x, y)$  und  $\hat{v}(x, y)$ , dass mit dem komplexen Eingangsintervall  $Z_1 = [-1.5, -0.5] + i \cdot [-1, 0]$  die gleiche obige Einschließung  $F(Z)$  berechnet wird.

Im **6. Beispiel** wählen wir das Punktintervall  $Z = [2, 2] + i \cdot [3, 3]$ . Mit  $\text{prec} = 400$  erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([1.92566973609167184450654\dots2310, 1.92566973609167184450654\dots2311] \\ [3.11579908410336514259287\dots6742, 3.11579908410336514259287\dots6743]).$$

Zeigen Sie mithilfe der Symmetrieeigenschaften von  $\hat{u}(x, y)$  und  $\hat{v}(x, y)$ , dass mit dem komplexen Eingangsintervall  $Z_1 = [-2, -2] + i \cdot [-3, -3]$  die gleiche obige Einschließung  $F(Z)$  berechnet wird.

### C.2.17. $\sqrt{-z^2 - 1}$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = \sqrt{-z^2 - 1}, \quad z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind. Für  $f(z) = \sqrt{-z^2 - 1}$  sind  $z_{\pm 1} = \pm i$  die Verzweigungspunkte. Die beiden Verzweigungsschnitte sind in der folgenden Abbildung C.33 schwarz bzw. grün dargestellt, vgl. Seite 286. Typische erlaubte und nicht erlaubte komplexe Intervalle  $Z$  sind entsprechend gelb bzw. rot dargestellt.

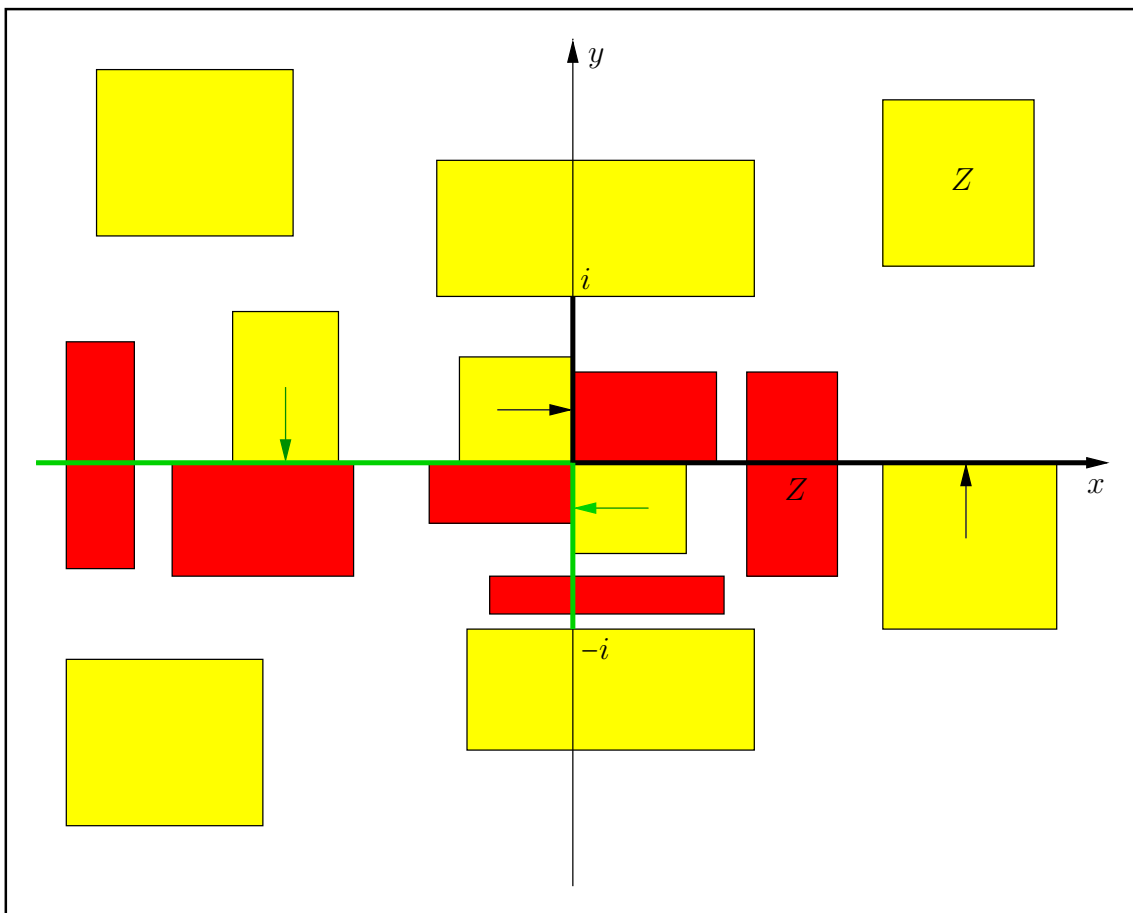


Abbildung C.33.: Erlaubte und nicht erlaubte Intervalle  $Z$  für den Hauptzweig von  $f(z)$

Beachten Sie, dass auch jedes Intervall der Form  $Z = [0, 0] + i \cdot Y$  oder  $Z = X + i \cdot [0, 0]$ , d.h. jedes Intervall  $Z$ , das auf einer Koordinatenachse liegt, ein erlaubtes Intervall ist. Die schwarzen und die grünen Pfeile geben die Richtungen an, aus denen die Funktionswerte von  $f(z) = \sqrt{-z^2 - 1}$ , mit  $z \in \mathbb{C}$ , auf die jeweiligen Verzweigungsschnitte analytisch fortgesetzt werden.

Die **optimale** Einschließung des Wertebereichs  $W$  durch das komplexe Rechteck  $F(Z) \supseteq W$  erfolgt mithilfe der Funktion

```
Mpfciclass sqrtmx2m1(const Mpfciclass& Z);
```

die in der Datei `pmfciclass.cpp` definiert ist.

### C.2.17.1. Numerische Ergebnisse

Im **1. Beispiel** wählen wir das Intervall  $Z = [2, 3] + i \cdot [1, 2]$ . Mit `prec = 400` erhalten wir für den Wertebereich  $W$  von  $f(z)$  für alle  $z \in Z$  die nahezu optimale Rechteckeinschließung  $F(Z)$  mit 120 korrekten Dezimalstellen:

$$W \subseteq F(Z) = ([0.910179721124454682608715...9580, 1.92566973609167184450654...2311] \\ [-3.14774949975311471570686...0773, -2.1286448445312042768138...1189]).$$

Im **2. Beispiel** wählen wir das gleiche Intervall  $Z = [2, 3] + i \cdot [1, 2]$  und berechnen jetzt aber eine Einschließung  $F_1(Z)$  des Wertebereichs  $W$  mit dem naiven Ansatz

$$\text{MpfciClass } z(\text{interval}(2,3), \text{interval}(1,2)), \text{ res}; \quad \text{res} = \text{sqrt}(-\text{sqr}(z)-1);$$

und erhalten die viel gröbere Einschließung

$$W \subseteq F_1(Z) = ([0.651482080258584212175291...2081, 2.34963769321913708102402...5287] \\ [-3.4641016151377545870548...0607, -1.60048518044024083831781...4802]).$$

Es gilt also:

$$W \subseteq F(Z) \subset F_1(Z).$$

Wird ein mit holomorphen Funktionen verschachtelter, komplexer Intervallausdruck, in dem die Intervallvariable  $Z$  nur einmal vorkommt, naiv ausgewertet, so erhält man nach obigem Beispiel mit  $F_1(Z)$  i.a. **keine optimale** Einschließung des Wertebereichs  $W$ .

Wertet man im Gegensatz dazu einen **reellen** Intervallausdruck aus, in dem die Intervallvariable  $x$  nur einmal vorkommt und der aus nur stetigen Funktionen zusammengesetzt ist, so erhält man stets **optimale** Einschließungen für den entsprechenden Wertebereich  $W$ . So erhält man z.B. für den Ausdruck  $\sqrt{x^2+1}$  mit `sqrt(sqr(x) ⊕ 1)` eine optimale Einschließung, während man mit `sqrt(x ⊕ x ⊕ 1)` i.a. eine starke Überschätzung erhält, weil im letzten Ausdruck die Variable  $x$  zweimal vorkommt und daher gilt: `sqr(x) ⊆ x ⊕ x`, vgl. auch Seite 74.

Um den Wertebereich  $W$  eines aus holomorphen Funktionen zusammengesetzten Intervallausdrucks möglichst optimal einzuschließen, sollte man auf möglichst viele Teilfunktionen zurückgreifen können, deren Wertebereiche selbst jeweils optimal eingeschlossen werden, vgl. dazu die Tabelle auf Seite 115.

Die Einschließung des Wertebereichs  $W_f$  von z.B.  $f(z) := \sqrt{-z^2-1} \cdot \log(1+z)$  wird man daher realisieren durch

$$\text{MpfciClass } z(\dots), F; \quad F = \text{sqr}(\text{tmx}2\text{m}1(z) \oplus \text{lnp}1(z));$$

Beachten Sie aber, dass die Einschließung von  $W_f$  durch  $F$  noch keinesfalls optimal sein muss, da uns nach der Tabelle auf Seite 115 nur optimale Einschließungen der beiden Teilausdrücke  $\sqrt{-z^2-1}$  und  $\log(1+z)$  zur Verfügung stehen, nicht aber eine optimale Einschließung des ganzen Produkts  $f(z) = \sqrt{-z^2-1} \cdot \log(1+z)$ . Da in der Realität nur wenige Intervallausdrücke wirklich optimal eingeschlossen werden können, kommen wir zum Ergebnis:

Die Einschließung des Wertebereichs  $W$  eines komplexen Intervallausdrucks ist i.a. ein eigenständiges Problem, das nur in Ausnahmefällen optimal gelöst werden kann, vgl. dazu auch Seite 120.



### C.2.18. $1/\sqrt{z}$

Der Verzweigungsschnitt von  $f(z) = 1/\sqrt{z}$ ,  $z = x + i \cdot y \in \mathbb{C} \setminus \{0\}$ , ist in der komplexen Ebene die negative reelle Achse. Ein achsenparalleles Argumentintervall  $Z$  darf diesen Verzweigungsschnitt nicht in seinem Innern enthalten und auch nicht von unten berühren. Die in `mpfciclass.cpp` definierte Funktion

```
Mpfciclass sqrt_r(const Mpfciclass& Z)
```

liefert mit  $W = \text{sqrt\_r}(Z)$  eine nahezu optimale, achsenparallele Rechteck-Einschließung der komplexen Wertemenge

$$\{w \in \mathbb{C} \mid w = 1/\sqrt{z}, z \in Z\} \subseteq W,$$

wobei  $\sqrt{z}$  der Hauptwert der komplexen Quadratwurzel ist.

In Abbildung C.34 sind einige erlaubte und **nicht erlaubte** achsenparallele Argumentintervalle  $Z$  der Funktion `sqrt_r(Z)` angegeben.

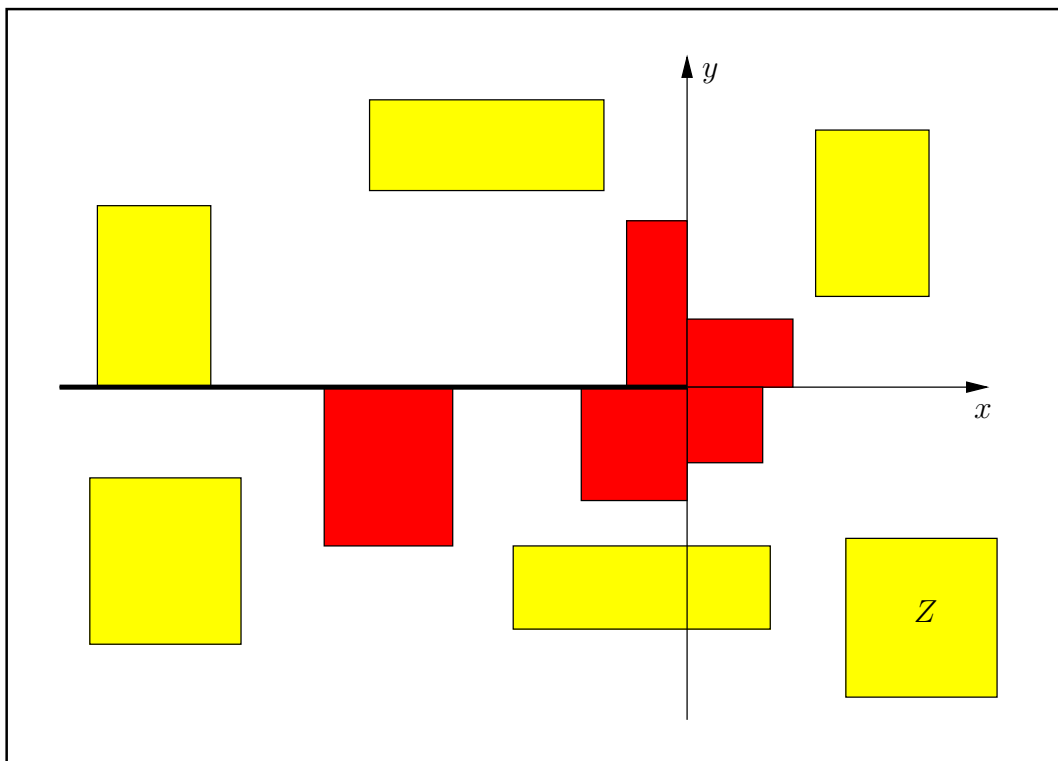


Abbildung C.34.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von `sqrt_r(Z)`.

#### C.2.18.1. Realteil

Der Realteil von  $f(z) = u(x, y) + i \cdot v(x, y)$  ist nach (C.19) von Seite 258 gegeben durch

$$u(x, y) := \frac{\sqrt{|z| + x}}{\sqrt{2} \cdot |z|}, \quad |z| = \sqrt{x^2 + y^2} \neq 0.$$

Für ein erlaubtes Intervall  $Z = X + i \cdot Y$  liegen dann die Maximum- und Minimumpunkte  $M, m$  der harmonischen Funktion  $u(x, y)$  auf dem Rand von  $Z(X, Y)$ , und wegen  $u(x, -y) \equiv u(x, y)$  kann man sich bei der Berechnung von  $M, m$  auf die obere Halbebene beschränken, d.h.  $u(x, y)$  nimmt auf den Rechtecken  $Z(X, Y)$  und  $Z^*(X, \text{abs}(Y))$  die gleichen Funktionswerte an, wobei  $\text{abs}(Y)$  das Intervall der Absolutbeträge von  $Y$  ist. Zur Vereinfachung werden wir in diesem Teilabschnitt das Rechteck  $Z^*(X, \text{abs}(Y))$  aus der oberen Halbebene wieder mit  $Z$  bezeichnen.

Zur Berechnung der Extrema  $m, M$  benötigen wir von  $u(x, y)$  die partiellen Ableitungen, die mit  $\alpha(x, y) > 0$  gegeben sind durch

$$\frac{\partial u(x, y)}{\partial x} = (y^2 - x \cdot (x + \sqrt{x^2 + y^2})) \cdot \alpha,$$

$$\frac{\partial u(x, y)}{\partial y} = -y \cdot (2x + \sqrt{x^2 + y^2}) \cdot \alpha.$$

Die Extremalkurven bez.  $\partial u(x, y)/\partial x = 0$  ergeben sich aus:  $y^2 - x \cdot (x + \sqrt{x^2 + y^2}) = 0$  zu:

1.  $y = 0$  für  $x < 0$ , d.h. auf der negativen reellen Achse gilt  $u(x, y) \equiv 0$ ,
2.  $y = \sqrt{3} \cdot x$  für  $x > 0$ , d.h. die Extremalkurve liegt nur im 1. Quadranten.

Die Extremalkurven bez.  $\partial u(x, y)/\partial y = 0$  ergeben sich aus:  $y \cdot (2x + \sqrt{x^2 + y^2}) = 0$  zu:

1.  $y = 0$  für  $x > 0$ , d.h. die Extremalkurve ist die positiven reellen Achse,
2.  $y = -\sqrt{3} \cdot x$  für  $x < 0$ , d.h. die Extremalkurve liegt nur im 2. Quadranten.

In folgender Abbildung ist in der oberen Halbebene das Monotonieverhalten von  $u(x, y)$  angegeben, und für typische Lagen von  $Z$  sind die Punkte  $m, M$  der gesuchten Extrema von  $u(x, y)$  auf dem Rand von  $Z$  dargestellt.

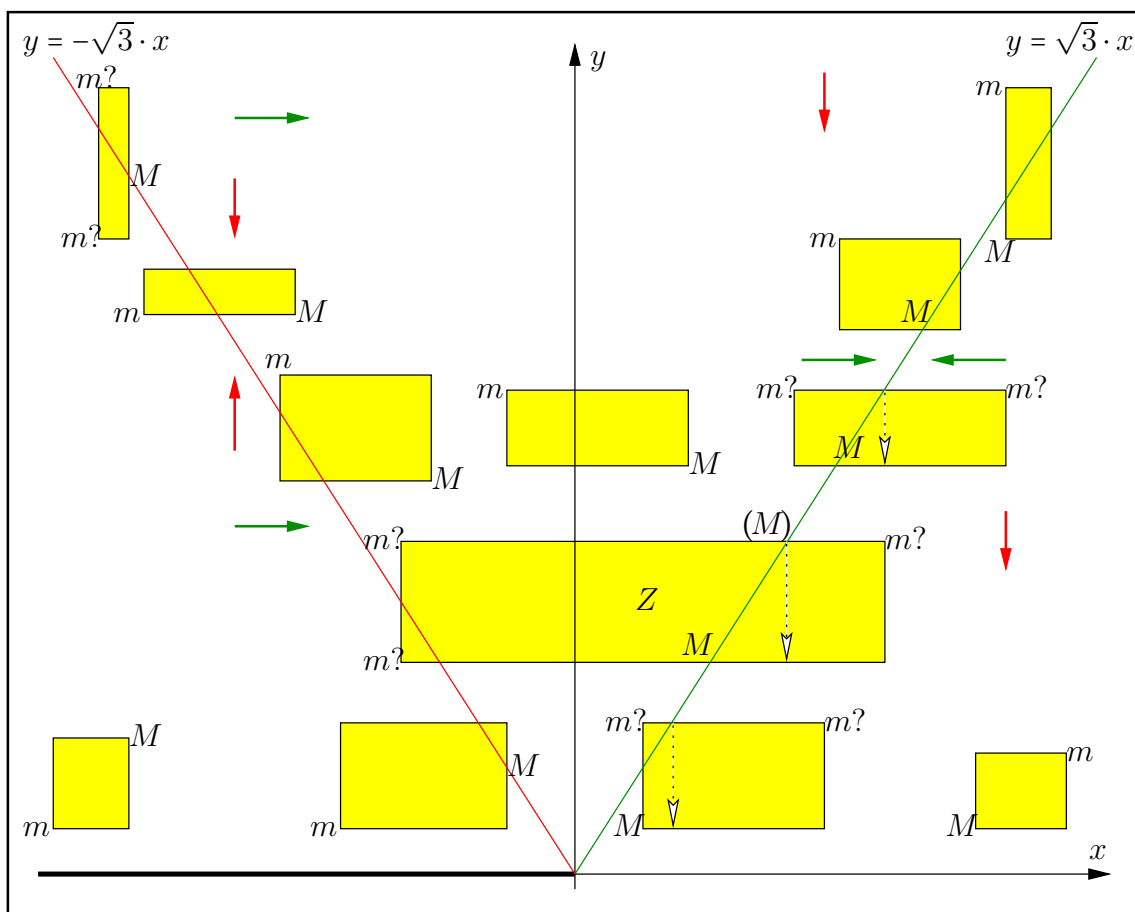


Abbildung C.35.: Typische Lagen von  $Z$  mit den Punkten  $m, M$  bez.  $u(x, y)$

Im Rechteck  $Z$  kann  $(M)$  nicht der Maximumpunkt sein, denn längs des punktierten Pfeils kann man nach unten und anschließend bis  $M$  nur in Richtung wachsender Funktionswerte von  $u(x, y)$  laufen. Analog argumentiert man bei den zwei anderen Rechtecken mit den punktierten Pfeilen.

Mit Abb. C.35 und mit den Abkürzungen:

ULP: Unterer linker Eckpunkt von  $Z$ ;    URP: Unterer rechter Eckpunkt von  $Z$ ;  
OLP: Oberer linker Eckpunkt von  $Z$ ;    ORP: Oberer rechter Eckpunkt von  $Z$ ;

ergibt sich der folgende Grobalgorithmus zur Bestimmung der Punkte  $m, M$  für das Minimum bzw. Maximum von  $u(x, y)$  über dem Rand von  $Z$ :

```
if (ULP unterhalb von  $y=\sqrt{3}\cdot x$ )
  M = ULP;    m = ORP oder m = OLP;
else
  if (Untere Parallele schneidet  $y=\sqrt{3}\cdot x$ )
    M ist dieser Schnittpunkt;    m = ULP oder m = OLP oder m = ORP;
  else
    if (URP liegt rechts von  $y=-\sqrt{3}\cdot x$ )
      M = LRP;    m = OLP oder m = ULP;
    else
      if (Rechte Parallele schneidet  $y=-\sqrt{3}\cdot x$ )
        M ist dieser Schnittpunkt;    m = ULP oder m = OLP;
      else
        M = ORP;    m = ULP;
```

#### Anmerkungen:

1. Weitere Einzelheiten findet man mit dem komplexen Intervall  $z = x + i \cdot y$  in der Funktion

```
MpfiClass sqrt_rRe(const MpfiClass& x, const MpfiClass& y)
```

die in `mpficlass.cpp` definiert ist.

2. Der maximale Funktionswert  $r$  wird z.B. im unteren linken Eckpunkt von  $z$  berechnet mit dem Aufruf: `r = sqrt_rRe(x1, y1, RoundUp)`;
3. Schneidet die rechte Parallele vom komplexen Rechteckintervall  $z$  die Extremalkurve  $y = -\sqrt{3} \cdot x$  im Punkt  $S(x_2, y_S)$ , so wird mit  $y_S = -\sqrt{3} \cdot x_2$  der maximale Funktionswert  $r$  bei geeigneter Rundung berechnet durch:  $r = \sqrt{2} \cdot (1/\sqrt{-x_2})/4$ .

### C.2.18.2. Imaginärteil

Der Imaginärteil von  $f(z) = u(x, y) + i \cdot v(x, y)$  ist für  $|z| = \sqrt{x^2 + y^2} \neq 0$  nach (C.19) von Seite 258 gegeben durch

$$v(x, y) := \begin{cases} \frac{-1}{\sqrt{-x}}, & y = 0 \wedge x < 0, \\ \frac{-y}{\sqrt{2} \cdot |z| \cdot \sqrt{|z| + x}}, & \text{sonst, d.h. falls } y \neq 0 \vee x \geq 0. \end{cases}$$

Wegen der für  $x < 0$  gültigen Grenzwerte

$$\lim_{y \rightarrow 0^+} v(x, y) = -1/\sqrt{-x} \quad \text{und} \quad \lim_{y \rightarrow 0^-} v(x, y) = +1/\sqrt{-x}$$

besitzt  $v(x, y)$  auf der negativen reellen Achse eine Sprungstelle und im Ursprung eine Singularität. Für ein erlaubtes Intervall  $Z = X + i \cdot Y$  liegen die Maximum- und Minimumpunkte  $M, m$  der harmonische Funktion  $v(x, y)$  auf dem Rand von  $Z(X, Y)$ , und wegen  $v(x, -y) \equiv -v(x, y)$  kann man sich bei der Berechnung der Funktionswerte  $l = v(m)$  und  $r = v(M)$  auf die obere Halbebene beschränken. Mit  $X = [x_1, x_2]$ ,  $Y = [y_1, y_2]$  erfolgt dies mit folgendem Algorithmus:

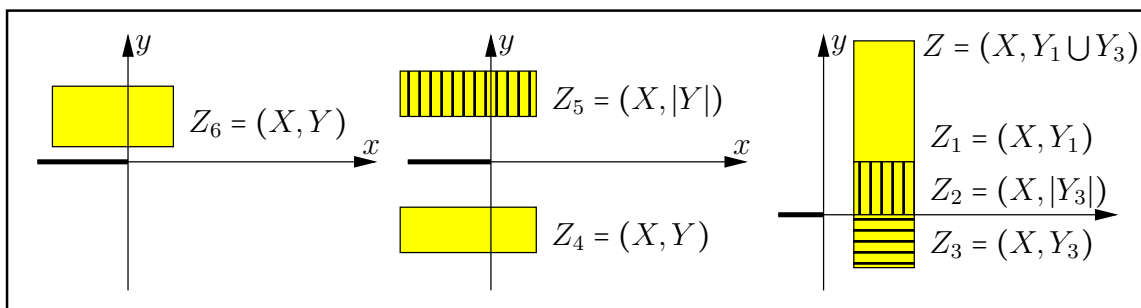


Abbildung C.36.: Zur Bestimmung von  $l, r$  bei Auswertung in der oberen Halbebene

**Algorithmus** zur Berechnung der Extrema  $l, r$  nur in der oberen Halbebene:

```

if (y1<0)
  if (x1>0 && y2>0)
    l in Z1 berechnen; // l ist Minimum in Z;
    l2 in Z2 = (X, |Y3|) berechnen;
    r = -l2; // r ist Maximum in Z;
  else // x1<=0 oder y2<=0:
    l in Z5 berechnen;
    r in Z5 berechnen;
    l4 = -r; // l4: Minimum in Z4
    r4 = -l; // r4: Maximum in Z4
  else // y1>=0
    l in Z6 berechnen;
    r in Z6 berechnen;

```

Wir benötigen daher die folgenden Funktionen:

1. Zur Berechnung von  $l$ , wenn  $Z_1$  oder  $Z_2$  im 1. Quadranten liegen:  
`MpfrClass sqrt_rIm1Q(const MpfiClass& x, const MpfiClass& y)`
2. Zur Berechnung von  $l, r$ , wenn  $Z_5$  oder  $Z_6$  in der oberen Halbebene liegen:  
`MpfiClass sqrt_rImOH(const MpfiClass& x, const MpfiClass& y)`

Um in der oberen Halbebene die Extrema  $l = v(m)$  und  $r = v(M)$  nach obigem Algorithmus berechnen zu können, benötigen wir von  $v(x, y)$  noch die partiellen Ableitungen, die für  $y \geq 0$  mit den positiven Faktoren  $\beta(x, y) > 0$  und  $\delta(x, y) > 0$  gegeben sind durch:

$$\frac{\partial v(x, y)}{\partial x} = y \cdot \left( y^2 + 3x \cdot (x + \sqrt{x^2 + y^2}) \right) \cdot \beta(x, y),$$

$$\frac{\partial v(x, y)}{\partial y} = (-2x + \sqrt{x^2 + y^2}) \cdot \sqrt{x + \sqrt{x^2 + y^2}} \cdot \delta(x, y).$$

Die Extremalkurven bez.  $\partial v(x, y)/\partial x = 0$  ergeben sich aus:  $y \left( y^2 + 3x \cdot (x + \sqrt{x^2 + y^2}) \right) = 0$  zu:

1.  $y = 0$  für  $x > 0$ , d.h. auf der positiven reellen Achse gilt:  $v(x, y) \equiv 0$ .
2.  $y = -\sqrt{3} \cdot x$  für  $x < 0$ , d.h. die Extremalkurve liegt nur im 2. Quadranten.

Die Extremalkurven bez.  $\partial v(x, y)/\partial y = 0$  ergeben sich aus:  $(-2x + \sqrt{x^2 + y^2}) \cdot \sqrt{x + \sqrt{x^2 + y^2}} = 0$  zu:

1.  $y = 0$  für  $x < 0$ , d.h. auf der negativen reellen Achse gilt:  $\partial v(x, y)/\partial y = 0$ .
2.  $y = \sqrt{3} \cdot x$  für  $x > 0$ , d.h. die Extremalkurve liegt nur im 1. Quadranten.

Um den Minimumpunkt  $m$  für ein Rechteck  $Z = (X, Y)$  aus dem 1. Quadranten mit Hilfe der Funktion `MpfrClass sqrt_rIm1Q(const MpfiClass& x, const MpfiClass& y)` bestimmen zu können, betrachten wir in folgender Abbildung typische Lagen dieser Rechtecke  $Z = (X, Y)$ :

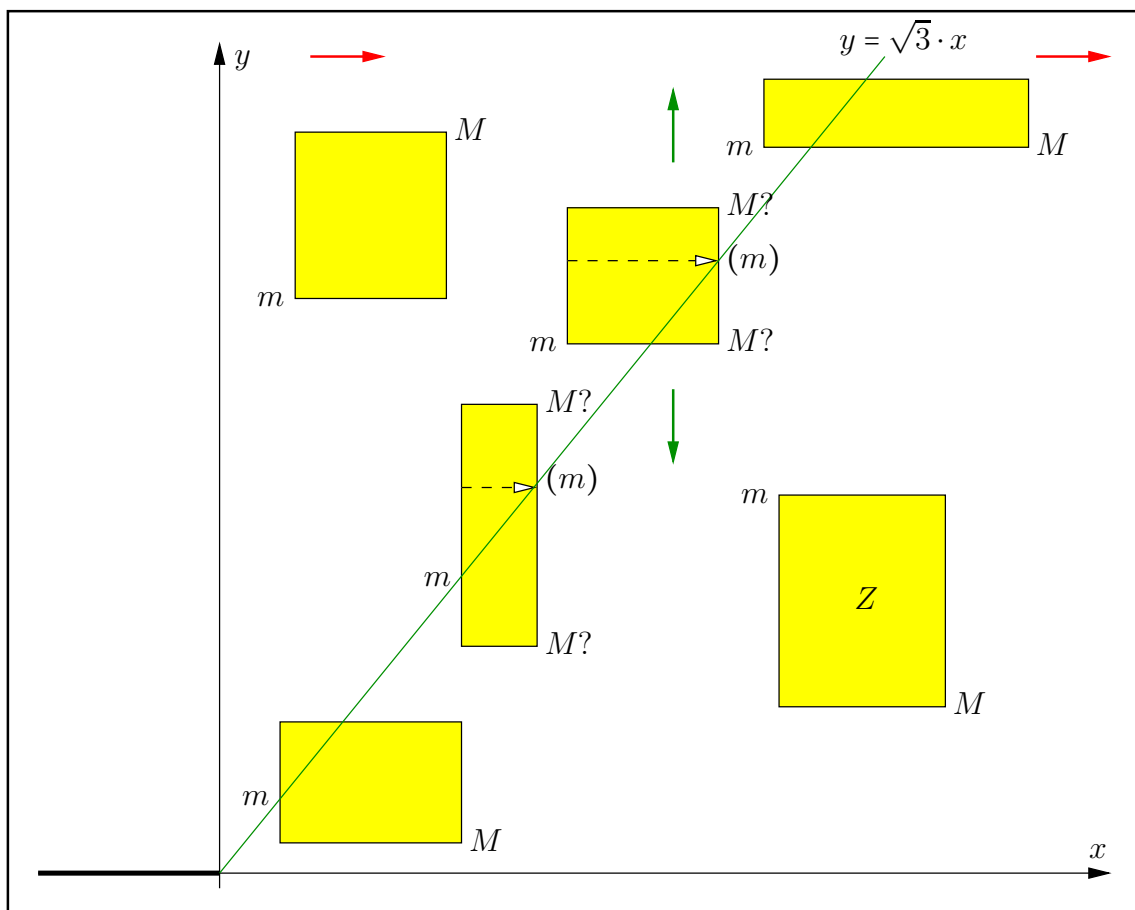


Abbildung C.37.: Zur Bestimmung von  $m, M$  bei Auswertung von  $v(x, y)$  im 1. Quadranten

Die roten und grünen Pfeile geben die Richtung wachsender Funktionswerte an. An den Punkten  $(m)$  kann kein Minimum angenommen werden und bei  $M?$  ist jeweils ein Maximum möglich.

Der Minimumpunkt  $m$  des Imaginärteils  $v(x, y)$  über einem Rechteck  $Z$  aus dem 1. Quadranten wird mit folgendem Algorithmus berechnet:

```

if (ULP liegt links von  $y = \sqrt{3} \cdot x$ )
   $m = \text{ULP}$ ;
else
  if (Linke Parallele schneidet  $y = \sqrt{3} \cdot x$  im Punkt S)
     $m = S$ ;
  else
     $m = \text{OLP}$ ; //  $m = \text{Oberer linker Eckpunkt von } Z$ ;

```

Um die Punkte  $m, M$  der Extrema für ein Rechteck  $Z = (X, Y)$  aus der oberen Halbebene mit Hilfe der Funktion `MpfiClass sqrt_rImOH( const MpfiClass& x, const MpfiClass& y )` bestimmen zu können, betrachten wir in folgender Abbildung typische Lagen dieser Rechtecke  $Z = (X, Y)$ :

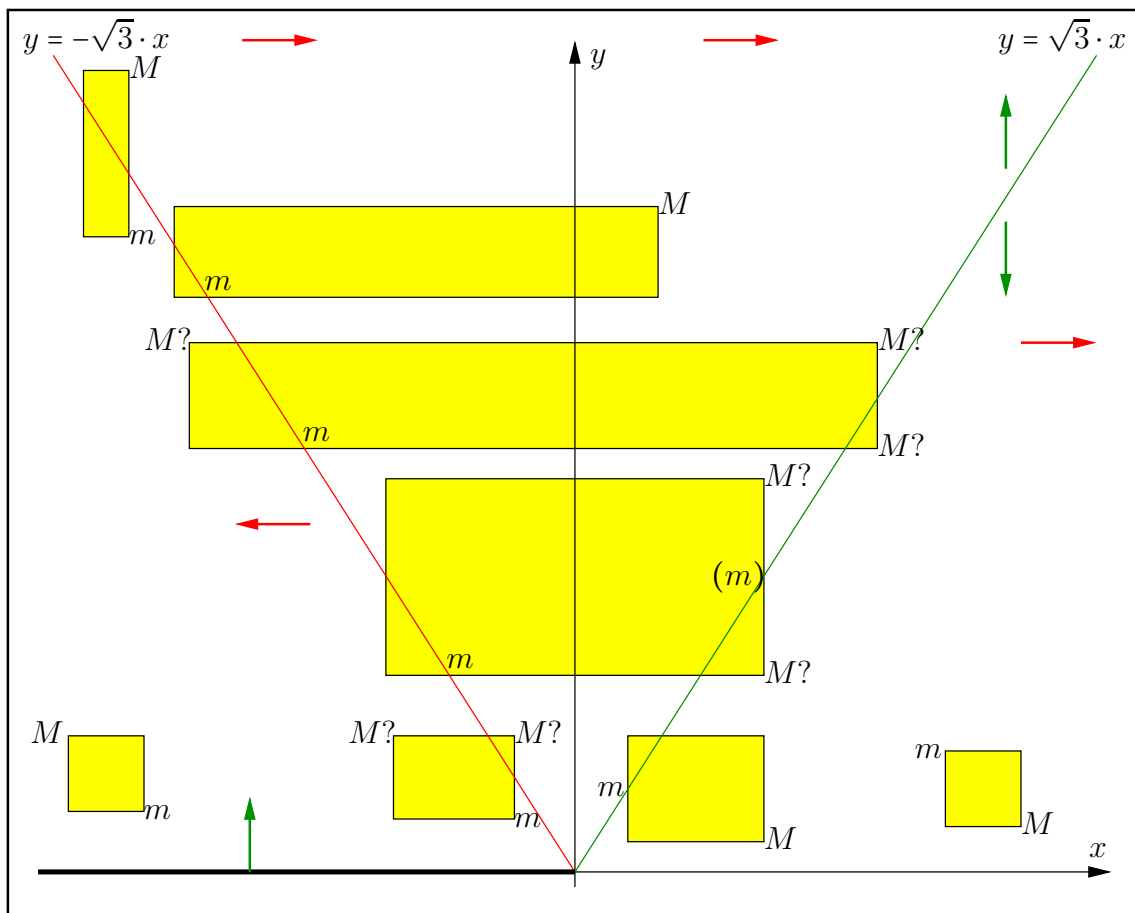


Abbildung C.38.: Lage von  $m, M$  bei Auswertung von  $v(x, y)$  in der oberen Halbebene

Die Lage der Extremalpunkte  $m, M$  für Rechtecke  $Z = (X, Y)$ , die mit ihren unteren Parallelen die Extremalkurve  $y = -\sqrt{3} \cdot x$  nicht schneiden und rechts von dieser Geraden liegen, findet man in Abb. C.37. Mit den Abbildungen C.37, C.38 und mit den Abkürzungen:

ULP: Unterer linker Eckpunkt von  $Z$ ;    URP: Unterer rechter Eckpunkt von  $Z$ ;  
 OLP: Oberer linker Eckpunkt von  $Z$ ;    ORP: Oberer rechter Eckpunkt von  $Z$ ;

ergibt sich der folgende Grobalgorithmus zur Bestimmung der Punkte  $m, M$  für das Minimum bzw. Maximum von  $v(x, y)$  über dem Rand von  $Z$ , wenn  $Z$  in der oberen Halbebene liegt:

```

if (URP liegt unterhalb von  $y = -\sqrt{3} \cdot x$ )
    m = URP; M = ORP oder M = OLP;
else
    if (Untere Parallele schneidet  $y = -\sqrt{3} \cdot x$  im Punkt S)
        m = S; M = OLP oder M = ORP oder M = URP;
    else // Z liegt rechts von  $y = -\sqrt{3} \cdot x$ ;
        if (ULP liegt links von  $y = +\sqrt{3} \cdot x$ )
            m = ULP; M = ORP oder M = URP;
        else
            if (Linke Parallele schneidet  $y = +\sqrt{3} \cdot x$  im Punkt S)
                m = S; M = ORP oder M = URP;
            else
                m = OLP; M = URP;

```

Der obige Algorithmus wird realisiert mit der Funktion

```
MpfiClass sqrt_rImOH(const MpfiClass& x, const MpfiClass& y)
```

wobei das Intervall  $y$  in der oberen Halbebene liegen muss.

### C.2.18.3. Numerische Beispiele

Im **1. Beispiel** wählen wir mit  $Z = [-1/2, -2^{-5000}] + i \cdot [0, 1]$  ein Argumentintervall in der Nähe der Singularität  $z_0 = 0$  und erhalten bei einer Präzision von `prec = 700` mit dem Funktionsaufruf  $W = \text{sqrt\_r}(Z)$  für den Wertebereich  $\{w \in \mathbb{C} \mid w = 1/\sqrt{z}, z \in Z\} \subseteq W$  mit 210 korrekten Dezimalstellen die nahezu optimale Einschließung

$$W = ([0, 1.328752 \dots 574 \cdot 10^{752}], [-3.758280 \dots 802 \cdot 10^{752}, -7.071067 \dots 853 \cdot 10^{-1}]).$$

Zum Vergleich erhalten wir bei der Intervallauswertung von  $W_I = 1 \diamond \text{sqrt}(Z)$  die viel gröbere Einschließung

$$W_I = ([0, 1.879140 \dots 901 \cdot 10^{752}], [-3.758280 \dots 802 \cdot 10^{752}, -5.321582 \dots 764 \cdot 10^{-753}]),$$

wobei die Oberschranke des Imaginärteils um 752 Zehnerpotenzen zu groß ausfällt.

Im **2. Beispiel** wählen wir  $Z = [1/4, 1/2] + i \cdot [-2, +1]$  und erhalten mit `prec = 700` und mit dem Funktionsaufruf  $W = \text{sqrt\_r}(Z)$  für den Wertebereich  $\{w \in \mathbb{C} \mid w = 1/\sqrt{z}, z \in Z\} \subseteq W$  mit 210 korrekten Dezimalstellen die nahezu optimale Einschließung

$$W = ([5.280517 \dots 552 \cdot 10^{-1}, 2.000000 \dots 000], [-7.071067 \dots 861 \cdot 10^{-1}, 7.071067 \dots 861 \cdot 10^{-1}]).$$

Zum Vergleich erhalten wir bei der Intervallauswertung von  $W_I = 1 \diamond \text{sqrt}(Z)$  auch jetzt die viel gröbere Einschließung

$$W_I = ([4.4139110926 \dots 396 \cdot 10^{-1}, 2.00000000 \dots 000], [-1.0000000 \dots 000, +1.0000000 \dots 000]),$$

wobei aber die Größenordnungen der Intervallgrenzen in diesem Beispiel übereinstimmen.

### C.2.19. $\cot(z)$

Mit  $z = x + i \cdot y$  gilt nach Abramowitz (4.3.58)

$$(C.69) \quad \cot(z) = \frac{\sin(2x) - i \cdot \sinh(2y)}{\cosh(2y) - \cos(2x)} = u(x, y) + i \cdot v(x, y)$$

$$(C.70) \quad = \frac{\cos(z)}{\sin(z)}, \quad z \neq k\pi, \quad k \in \mathbb{Z}$$

$$(C.71) \quad = \frac{1}{\tan(z)}$$

$$(C.72) \quad = \tan(\pi/2 - z) = \tan(\pi/2 - z + k\pi), \quad k \in \mathbb{Z}$$

Die Polstellen der Cotangens-Funktion sind gegeben durch  $z_{p,k} = k\pi$ , und ihre Nullstellen liegen bei  $z_{s,k} = \pi \cdot (k + 1/2)$ ,  $k \in \mathbb{Z}$ .

Da der Tangens für komplexe Intervallargumente  $Z$  bereits implementiert ist, soll der ebenfalls  $\pi$ -periodische Cotangens mit Hilfe des Tangens realisiert werden. Die Gleichung (C.71) ist dabei ungeeignet, denn im Falle  $\pi/2 \in Z$  liegt eine Polstelle des Tangens in  $Z$ , so dass mit (C.71) Programmabbruch erfolgt, obwohl der Cotangens in der Umgebung von  $\pi/2$  existiert und bei  $z_{s,0} = \pi/2$  eine Nullstelle besitzt. Die Darstellung (C.70) ist ebenfalls ungeeignet, da zwei Ergebnisintervalle zu dividieren sind, womit eine erhebliche, zusätzliche Überschätzung verbunden ist. Zur Implementierung des Cotangens wählen wir daher (C.72). Ist ein rechteckiges Argumentintervall  $Z = X + i \cdot Y$  vorgegeben, so kann es bei der intervallmäßigen Auswertung von  $\pi/2 - Z$  zu den bekannten Überschätzungen kommen, wenn  $\text{Inf}(Z) \approx \pi/2$  oder  $\text{Sup}(Z) \approx \pi/2$  realisiert werden. Diese Überschätzungen sind prinzipiell unvermeidbar, da  $\pi/2$  keine Maschinenzahl ist und daher durch ein echtes Intervall eingeschlossen werden muss. Insbesondere bei Punktintervallen  $Z \approx \pi/2$  lassen sich diese Überschätzungen jedoch vermeiden, wenn man die Intervalldifferenz  $\pi/2 - Z$  in doppelter Präzision ausführt. Beachten Sie jedoch, dass die Intervalldifferenz  $\pi/2 - Z$  selbst dann von  $[0, 0]$  verschieden ist, wenn  $Z$  eine optimale Einschließung von  $\pi/2$  ist, so dass dann die Auswertung von  $\pi/2 - Z$  zu großen Überschätzungen führen muss.

Im folgenden **1. Beispiel** wählen wir ein Punktintervall  $X = [x, x]$ , wobei im Zahlenformat mit der Current-Präzision `prec = 300`

$$x = \text{Sup}(\text{Pi}()/2)$$

die kleinste Maschinenzahl größer  $\pi/2$  ist. Mit  $Z = ([x, x] + i \cdot [1, 1])$  liefert die Version mit doppelter Präzision `prec = 600` die Einschließung

$$\cot(Z) = ([-3.08764047\text{e-}91, -3.08764046\text{e-}91], [-7.61594156\text{e-}1, -7.61594155\text{e-}1])$$

während man bei nur einfacher Präzision für den Realteil die praktisch unbrauchbare Einschließung erhält:

$$\cot(Z) = ([-4.12338660\text{e-}91, 0.00000000], [-7.61594156\text{e-}1, -7.61594155\text{e-}1]).$$

Im **2. Beispiel** wählen wir mit  $Z = (\text{Pi}()/2 + i \cdot [1, 1])$  ein echtes komplexes Maschinenintervall, das im Zahlenformat mit der Current-Präzision `prec = 300` den Realteil  $\pi/2$  optimal einschließt. Aber selbst mit doppelter Präzision `prec = 600` erhält man jetzt nur die grobe Einschließung

$$\cot(Z) = ([-3.08764047\text{e-}91, 1.03574613\text{e-}91], [-7.61594156\text{e-}1, -7.61594155\text{e-}1]),$$

die auch bei einer drei- oder vierfachen Präzision beim Realteil **grundsätzlich** nicht verbessert werden kann.



### C.2.20. $\arcsin(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - \{(-\infty, -1) \cup (+1, +\infty)\}$  liefert die Funktion

```
MpfciClass asin(const MpfciClass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\text{asin}(Z)$  für die komplexen Funktionswerte  $\arcsin(z)$ , mit  $z \in Z$ .

$$\{\arcsin(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \text{asin}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der reellen Achse von  $-1$  bis  $-\infty$  bzw. von  $+1$  bis  $+\infty$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben.

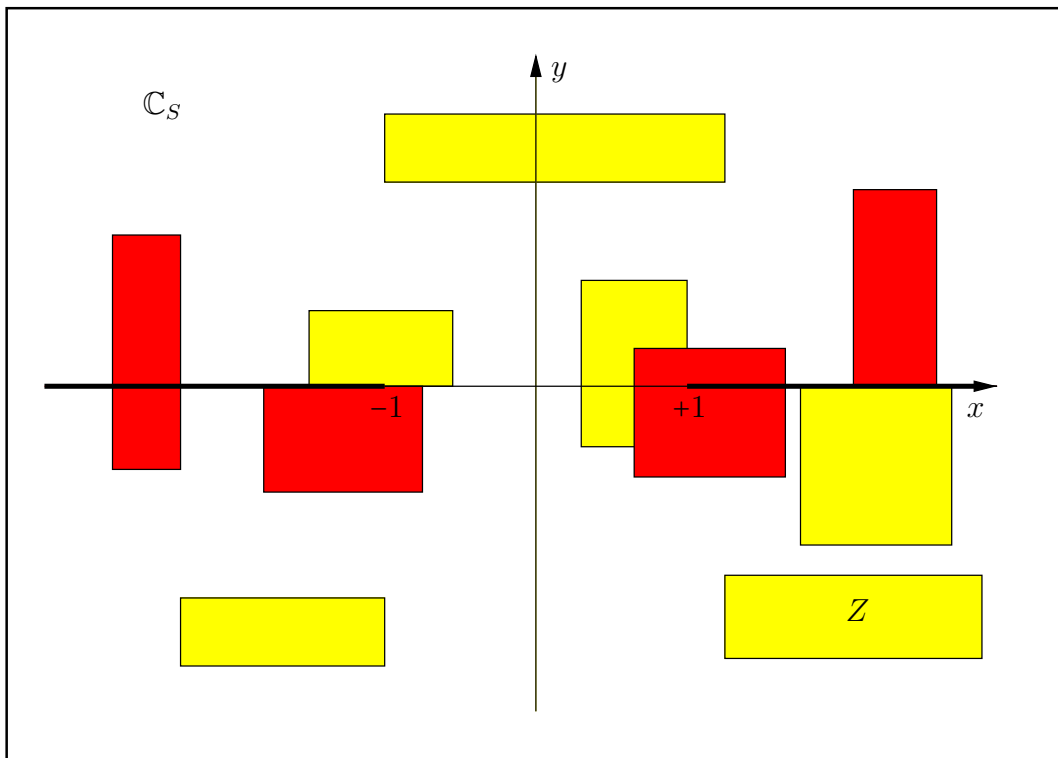


Abbildung C.39.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{asin}(Z)$ .

Mit  $Z = [1, 1] + i \cdot [2^{-1073741824}, 2^{-1073741824}]$  erhält man mit der Current-Präzision  $\text{prec} = 70$  die Einschließung<sup>8</sup>

$$\text{asin}(Z) = ([1.57079632679489661922, 1.57079632679489661924], [4.88115243040816240520e - 161614249, 4.88115243040816240522e - 161614249]).$$

Mit  $Z = [0.5, 1.5] + i \cdot [-1, 0]$  erhält man mit der Current-Präzision  $\text{prec} = 80$  die Einschließung

$$\text{asin}(Z) = ([3.49439062857213293627411825e - 1, 1.57079632679489661923132170], [-1.26047518779845407285290814, 0.0000000000000000000000000000000000]).$$

<sup>8</sup>Die Maschinenzahl  $2^{-1073741824} = \text{minfloat}()$  ist präzisionsunabhängig die kleinste positive Maschinenzahl.

### C.2.20.1. Algorithmus

Der nachfolgende Algorithmus zur Berechnung der Rechteck-Einschließung  $\text{asin}(Z)$  für ein vorgegebenes achsenparalleles Rechteckintervall  $Z$  basiert auf [55], [54], [12], wobei sich die Verbesserungen in [12] auf das dort verwendete IEEE-Format beziehen. Da  $\text{asin}(Z)$  jetzt im MPFCI-Format implementiert wird, müssen die Verbesserungen entsprechend an dieses Format angepasst werden.

Mit  $z = x + i \cdot y \in \mathbb{C}$  gelten nach W. Krämer für den komplexen Funktionswert  $w = \arcsin(z) = \Re(w) + i \cdot \Im(w)$  die folgenden Beziehungen, [38]:

$$(C.73) \quad \alpha := \frac{1}{2} \cdot \left\{ \sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2} \right\}$$

$$(C.74) \quad \beta := \frac{1}{2} \cdot \left\{ \sqrt{(x+1)^2 + y^2} - \sqrt{(x-1)^2 + y^2} \right\}$$

$$(C.75) \quad \beta = \frac{x}{\frac{1}{2} \cdot \left\{ \sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2} \right\}} = \frac{x}{\alpha}$$

$$(C.76) \quad \Re(w) := \arcsin(\beta)$$

$$(C.77) \quad \Im(w) = \begin{cases} +\text{arcosh}(\alpha), & \text{falls } y > 0 \\ +\text{arcosh}(\alpha), & \text{falls } y = 0 \text{ und } x \leq -1 \\ 0, & \text{falls } y = 0 \text{ und } -1 \leq x \leq +1 \\ -\text{arcosh}(\alpha), & \text{falls } y = 0 \text{ und } x \geq +1 \\ -\text{arcosh}(\alpha), & \text{falls } y < 0 \end{cases}$$

Wegen  $\alpha(x, y) \geq \alpha(x, 0) = (|x+1| + |x-1|)/2$  erhält man mit einfachen Fallunterscheidungen ( $|x| \leq 1$ ,  $x < -1$ ,  $x > +1$ ):

$$(C.78) \quad \alpha(x, y) \geq \max(1, |x|) = \begin{cases} 1, & \text{wenn } |x| \leq 1, \\ |x|, & \text{wenn } |x| > 1, \end{cases} \quad \text{oder größer: } \alpha(x, y) \geq 1.$$

Unter den Voraussetzungen  $y = 0$  und  $-1 \leq x \leq +1$  gilt:

$$(C.79) \quad \alpha(x, 0) \equiv 1, \quad \text{d.h. } \Im(w) = \pm \text{arcosh}(1) = 0, \quad \Re(w) = \arcsin(x);$$

Es gilt außerdem

$$|\beta(x, y)| \leq |\beta(x, 0)| \leq 1$$

Zum Beweis gilt nach (C.75):

$$|\beta(x, y)| \leq |\beta(x, 0)| = \frac{2 \cdot |x|}{|x+1| + |x-1|} =: R(x),$$

und wegen  $R(-x) \equiv R(x)$  kann man sich auf  $x \geq 0$  beschränken:

$$\text{Sei } 0 \leq x \leq 1, \quad \text{d.h. } |x-1| = -(x-1) \rightsquigarrow R(x) = \frac{2x}{(x+1)-(x-1)} = x \leq 1;$$

$$\text{Sei } x > 1, \quad \text{d.h. } |x-1| = +(x-1) \rightsquigarrow R(x) = \frac{2x}{(x+1)+(x-1)} = 1 \quad \blacksquare$$

In [54] wird für das komplexe Intervallargument

$$z = \mathbf{x} + i \cdot \mathbf{y} = [x_1, x_2] + i \cdot [y_1, y_2]$$

eine Einschließung für  $\arcsin(z)$  berechnet, wobei die beiden reellen Funktionen  $\arcsin(\beta)$  und  $\text{arcosh}(\alpha)$  für **Intervallargumente**  $\alpha$  und  $\beta$  auszuwerten sind. Dabei sind in (C.73) und (C.75) die reellen Werte  $x, y$  durch entsprechende Punktintervalle  $\mathbf{x} = [x_1, x_2]$ ,  $\mathbf{y} = [y_1, y_2]$ ,  $x_i, y_i \in S(2, 53)$  zu ersetzen.

Bei der Auswertung von  $\alpha$  und  $\beta$  nach (C.73) und (C.75) ist folgender Term  $T(x, y)$  zu berechnen:

$$(C.80) \quad T := \sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2}$$

Um  $T < \text{MaxFloat}()$  für jede Präzision zu gewährleisten, betrachten wir folgende Abschätzungen:

Zunächst gilt:  $T(x, y) = T(|x|, y) = \sqrt{(|x|+1)^2 + y^2} + \sqrt{(|x|-1)^2 + y^2}$ ;

Der Beweis ergibt sich direkt mit den Fallunterscheidungen:  $x \geq 0$  und  $x < 0$ . Man findet zusätzlich:

$$\begin{aligned} T(|x|, y) \equiv T(|x|, |y|) &\leq \sqrt{(|x|+1)^2 + y^2} + \sqrt{(|x|+1)^2 + y^2} \\ &= 2 \cdot \sqrt{(|x|+1)^2 + y^2} < 2 \cdot \sqrt{(|x|+1)^2 + (|y|+1)^2} \\ &\leq 2\sqrt{2} \cdot (M+1), \quad \text{falls } M = \max\{|x|, |y|\}. \end{aligned}$$

Bei der Auswertung von  $T(x, y)$  wird also Overflow verhindert, falls gilt<sup>9</sup>:

$$\begin{aligned} 2\sqrt{2} \cdot (M+1) &< \text{MaxFloat}(2) \\ \Leftrightarrow M &< \frac{\text{MaxFloat}(2)}{2\sqrt{2}} - 1 = 0.530330085889 \dots \cdot 2^{1073741822} \end{aligned}$$

Mit  $gr := 0.530330085889 \cdot 2^{1073741822}$  und der Anweisung

```
if (|x|>gr or |y|>gr) 'Programm-Abbruch'
```

verhindert man daher einen vorzeitigen Overflow bei der Berechnung von  $T(x, y)$ , wobei die obige Einschränkung in der Praxis keine Bedeutung haben wird.  $\alpha := T(x, y)/2$  wird mit der Hilfsfunktion

```
MpfiClass f_aux_asin(const MpfiClass& x, const MpfiClass& y)
```

zusammen mit `sqrtx2y2()` ausgewertet, wobei für  $x, y$  nur Punktintervalle übergeben werden. Die oben beschriebene Abfrage `if (|x|>gr or |y|>gr)` erfolgt nicht in `f_aux_asin(x, y)` selbst, sondern in einer übergeordneten Funktion. In `f_aux_asin(x, y)` wird zusätzlich noch die erste Ungleichung in (C.78) überprüft, die durch Intervallüberschätzungen verletzt sein kann.

Eine nahezu optimale Einschließung des Realteils  $\Re(w)$  nach (C.76) durch die Auswertung von  $\arcsin(\beta)$  scheint damit gewährleistet zu sein, wenn man zusätzlich im entsprechenden Algorithmus aus [54] (`2nd: real part`) den offensichtlichen Tippfehler, `hx1` statt `hxu`, beseitigt:

```
if( srez < 0.0 )
    resxu = Sup(asin(hxu/f_aux_asin(hxu, interval(max(-iimz, simz))));
```

Für  $\beta \rightarrow -1$  oder  $\beta \rightarrow +1$  wird jedoch der Realteil  $\Re(w)$  durch die Auswertung von  $\arcsin(\beta)$  noch nicht optimal eingeschlossen, da wegen der nahezu vertikalen Tangenten an den Intervallrändern von  $\beta \in [-1, +1]$  nur minimale Überschätzungen beim Argument  $\beta$  große Überschätzungen beim reellen  $\arcsin(\beta)$  bewirken.

Diese Überschätzungen lassen sich jedoch im Falle  $\beta \geq 0.75$ , d.h.  $\beta \rightarrow +1$  durch die folgende Transformation  $\beta = 1 - \delta$  bzw.  $\delta = 1 - \beta$  beseitigen:

$$(C.81) \quad \arcsin(\beta) = \arcsin(1 - \delta) = \frac{\pi}{2} - \arcsin(\sqrt{\delta \cdot (2 - \delta)}), \quad 0 \leq \delta \leq 2;$$

Wertet man jetzt in Gleichung (C.81) den letzten Ausdruck intervallmäßig aus, so wird von der optimalen Einschließung von  $\pi/2$  der betragsmäßig kleine Intervallausdruck  $\arcsin(\sqrt{\delta \cdot (2 - \delta)})$  subtrahiert, so dass keine nennenswerten Überschätzungen auftreten können. Die Berechnung von  $\delta(x, y) := 1 - \beta(x, y)$  ist auch bei der Einschließung des Realteils von  $\arccos(z)$  erforderlich. Für die Fälle  $0.75 \leq x < 1$ ,  $x = 1$  und  $x > 1$  findet man die entsprechenden Formeln in (C.147), (C.137) und (C.114). Wegen der Identität  $\arcsin(-\beta) \equiv -\arcsin(\beta)$  kann der andere Fall  $\beta \rightarrow -1$  auf den oben behandelten Fall  $\beta \rightarrow +1$  zurückgeführt werden. Die Auswertung von  $\arcsin(\beta)$  erfolgt mit der Hilfsfunktion

<sup>9</sup>Für  $\text{prec} \geq 2$  gilt:  $\text{MaxFloat}(2) \leq \text{MaxFloat}(\text{prec})$ , d.h.  $\text{MaxFloat}(\text{prec})$  ist präzisionsabhängig!

MpfiClass Asin\_beta( const MpfiClass& x, const MpfiClass& y );

Bei der intervallmäßigen Auswertung von  $\delta(x, y) \geq 0$  kann in seltenen Fällen das Infimum seiner Einschließung negativ werden, wodurch später ein NaN erzeugt wird. Diese Fehlerquelle wird in `Asin_beta(x,y)` mit der Abfrage `if (Inf(delta)<0)` beseitigt.

Wir kommen jetzt zur Berechnung einer Einschließung des Imaginärteils, wobei nach (C.77)  $\operatorname{arcosh}(\alpha)$  für das Intervallargument  $\alpha$  auszuwerten ist. Nach (C.73) sind dabei für  $x, y$  entsprechende Punktintervalle  $\mathbf{x} = [x_1, x_1]$ ,  $\mathbf{y} = [y_1, y_1]$ ,  $x_1, y_1 \in S(2, 53)$  einzusetzen:

$$\alpha := \frac{1}{2} \cdot \left\{ \sqrt{(\mathbf{x} + 1)^2 + \mathbf{y}^2} + \sqrt{(\mathbf{x} - 1)^2 + \mathbf{y}^2} \right\}$$

Bei der Auswertung von  $\operatorname{arcosh}(\alpha)$  treten jedoch für  $\operatorname{Sup}(\alpha) \rightarrow +1$  erhebliche Überschätzungen auf, da die  $\operatorname{arcosh}$ -Funktion dann in der Nähe ihrer Nullstelle 1 zu berechnen ist. Diese Überschätzungen des Imaginärteils lassen sich jedoch vermeiden, wenn man mit der folgenden Transformation  $\alpha = 1 + (\alpha - 1) = 1 + \delta$ ,  $\delta := \alpha - 1 \geq 0$  die Identität

$$(C.82) \quad \operatorname{arcosh}(\alpha) \equiv \operatorname{arcosh}(1 + \delta)$$

benutzt und dabei  $\operatorname{arcosh}(1 + \delta)$  mit Hilfe der Funktion `acoshp1(delta)` auswertet.

Bevor wir auf die Berechnung von  $\delta := \alpha - 1 \geq 0$  näher eingehen, bestimmen wir zunächst in der komplexen Ebene diejenigen Bereiche, in denen  $\operatorname{arcosh}(\alpha)$  direkt oder mit Hilfe von  $\operatorname{arcosh}(1 + \delta)$  auszuwerten ist. Dazu beweisen wir mit  $z = x + i \cdot y$ :

$$(C.83) \quad |x| \geq 2 \vee |y| \geq 2 \implies \alpha(x, y) \geq 2;$$

Zum Beweis benutzen wir:  $\alpha(x, y) \geq \alpha(x, 0) = (|x + 1| + |x - 1|)/2 =: r(x)$ . Sei zunächst  $|x| \geq 2$ , dann folgt mit einfachen Fallunterscheidungen:  $r(x) \geq 2$ . Außerdem gilt:  $\alpha(x, y) \geq \alpha(0, y) = |y|$ , und mit  $|y| \geq 2$  folgt der zweite Teil ■

Außerhalb eines Quadrats mit der Seitenlänge 4 und dem Mittelpunkt im Ursprung der komplexen Ebene benutzen wir daher wegen (C.83)  $\operatorname{arcosh}(\alpha)$  direkt und innerhalb dieses Quadrats gilt:  $\operatorname{arcosh}(\alpha) \equiv \operatorname{arcosh}(1 + \delta)$ , mit  $\delta := \alpha - 1$ .

Wir kommen jetzt zur Berechnung von  $\delta := \alpha - 1 \geq 0$ :

Man findet zunächst

$$(C.84) \quad x = \pm 1 \rightsquigarrow \delta = \left\{ \sqrt{1 + \left(\frac{y}{2}\right)^2} - 1 \right\} + \frac{|y|}{2},$$

wobei der erste Summand  $\{\dots\}$  mit Hilfe der **C-XSC** Funktion `sqrtp1m1(...)` berechnet wird. Für  $|x| \neq 1$  gilt mit  $K(x) := |x + 1| + |x - 1| - 2$ :

$$\begin{aligned} 2\delta &= |x + 1| \cdot \sqrt{1 + \left(\frac{y}{x + 1}\right)^2} - 2 + \sqrt{(x - 1)^2 + y^2} \\ &= |x + 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x + 1}\right)^2} - 1 \right\} + |x + 1| - 2 + \sqrt{(x - 1)^2 + y^2} \\ &= |x + 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x + 1}\right)^2} - 1 \right\} + |x - 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x - 1}\right)^2} - 1 \right\} \\ &\quad + K(x) \end{aligned}$$

Mit der folgenden Darstellung für  $K(x)$

$$K(x) = \begin{cases} 0 & , \text{ falls } 0 \leq |x| \leq 1 \\ 2 \cdot (|x| - 1) & , \text{ falls } |x| > 1 \end{cases} \quad \text{und mit}$$

$$V(x, y) := |x + 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x + 1}\right)^2} - 1 \right\} + |x - 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x - 1}\right)^2} - 1 \right\}$$

gilt dann insgesamt:

$$(C.85) \quad \delta = \frac{V(x, y)}{2}, \quad \text{falls } |x| < 1$$

$$(C.86) \quad \delta = \frac{V(x, y)}{2} + (|x| - 1), \quad \text{falls } |x| > 1$$

$$(C.87) \quad \delta = \left\{ \sqrt{1 + \left(\frac{y}{2}\right)^2} - 1 \right\} + \frac{|y|}{2}, \quad \text{falls } |x| = 1,$$

dabei werden die mit  $\{\dots\}$  gekennzeichneten Ausdrücke mit der **C-XSC** Funktion `sqrtp1m1()` berechnet. Bitte beachten Sie, dass innerhalb des oben beschriebenen Quadrats mit der Kantenlänge 4 die quadratischen Argumente dieser Funktion, z.B.  $(y/2)^2$  oder  $(y/(x+1))^2$ , ohne vorzeitigen Overflow berechnet werden. Beachten Sie außerdem, dass  $V(x, y)$  eine Summe aus nur positiven Summanden ist, so dass keine Auslöschungen zu befürchten sind.

Die Auswertung von  $\operatorname{arcosh}(\alpha)$  nach (C.82) mit Hilfe der Ausdrücke für  $\delta$  nach (C.85) und (C.87) können im Falle  $\delta \ll 1$  noch wesentlich vereinfacht werden. Dabei sind für  $x, y$  Punktintervalle einzusetzen, so dass in (C.85) und (C.87)  $\delta$  als Intervallargument  $\delta$  zu verstehen ist.

Wir beginnen mit dem Fall:  $|x| = 1$ ;

$$\begin{aligned} \operatorname{arcosh}(\alpha) &\equiv \operatorname{arcosh}(1 + \delta), \quad \delta = \left\{ \sqrt{1 + \left(\frac{y}{2}\right)^2} - 1 \right\} + \frac{|y|}{2}; \quad \rightsquigarrow \\ &\equiv \operatorname{arcosh}(1 + \delta), \quad \delta = \left\{ \sqrt{1 + \left(\frac{t}{2}\right)^2} - 1 \right\} + \frac{t}{2}, \quad t := |y| \geq 0; \\ &\equiv \operatorname{arcosh}(\sqrt{1 + (t/2)^2} + t/2), \quad t = |y| \geq 0; \end{aligned}$$

Die wesentliche Vereinfachung besteht nun darin, für  $t \ll 1$  die Potenzreihe des obigen Ausdrucks  $\operatorname{arcosh}(\sqrt{1 + (t/2)^2} + t/2)$  zu benutzen. Um Aussagen über diese Reihe machen zu können, geben wir zunächst die Potenzreihe für  $\operatorname{arcosh}(1 + t)$  an<sup>10</sup>:

$$(C.88) \quad \operatorname{arcosh}(1 + t) = \sqrt{2t} \cdot Q(t);$$

$$Q(t) := \sum_{k=0}^{\infty} a_k \cdot t^k = 1 - \frac{1}{12}t + \frac{3}{160}t^2 - \frac{5}{896}t^3 + \frac{35}{18432}t^4 - \frac{63}{90112}t^5 \pm \dots$$

$$a_k = \frac{1}{(2k+1) \cdot 2^k} \cdot \binom{-0.5}{k}, \quad a_0 = 1,$$

$$a_{k+1} = -a_k \cdot \frac{(2k+1)^2}{4 \cdot (2k+3)(k+1)}, \quad k = 0, 1, 2, \dots$$

Mit Hilfe des Majorantenkriteriums und der geometrischen Reihe kann leicht gezeigt werden, dass die Reihe für  $Q(t)$  sogar absolut konvergiert, wenn  $t < 1$ .

Mit `Maple7` und `Mathematica` erhält man:

$$\operatorname{arcosh}(\sqrt{1 + (t/2)^2} + t/2) = \sqrt{t} \cdot H(t), \quad H(t) = \sum_{k=0}^{\infty} b_k,$$

$$H(t) = 1 + \frac{1}{12}t - \frac{3}{160}t^2 - \frac{5}{896}t^3 + \frac{35}{18432}t^4 + \frac{63}{90112}t^5 - - + \dots,$$

<sup>10</sup>Vgl. die Ergebnisse zur Funktion  $\operatorname{arcosh}(1+x)$  im Buch zur Fehlerabschätzung; dort wird auch die folgende Ungleichung gezeigt:  $|a_{k+1}| < |a_k|$ ,  $k = 0, 1, \dots$

und es gilt vermutlich<sup>11</sup>  $|a_k| = |b_k|$ . Unter dieser Voraussetzung ist dann auch  $H(t)$  absolut konvergent, so dass wie folgt geklammert werden kann:

$$(C.89) \quad \begin{aligned} H(t) &= \left(1 + \frac{1}{12}t\right) - \left(\frac{3}{160}t^2 + \frac{5}{896}t^3\right) + \left(\frac{35}{18432}t^4 + \frac{63}{90112}t^5\right) - + \dots \\ &= h_1 - h_2 + h_3 - h_4 \pm \dots, \quad h_j \geq 0; \end{aligned}$$

Damit erhalten wir für  $H(t)$  eine *alternierende* Reihe. Wegen  $|a_{k+1}| < |a_k|$  kann unter der Voraussetzung  $t < 1$  leicht gezeigt werden:  $h_{k+1} < h_k$ , so dass die Reihe in (C.89) eine alternierende *Leibniz-Reihe* ist. Damit ergeben sich im Falle  $|x| = 1$  und  $t = |y| < 1$  mit  $\operatorname{arcosh}(\alpha) = \operatorname{arcosh}(\sqrt{1 + (t/2)^2} + t/2)$  die folgenden Abschätzungen:

$$(C.90) \quad \sqrt{t} \cdot \left[ \left(1 + \frac{t}{12}\right) - t^2 \left(\frac{3}{160} + \frac{5t}{896}\right) \right] \leq \operatorname{arcosh}(\alpha), \quad t = |y| < 1;$$

$$(C.91) \quad \operatorname{arcosh}(\alpha) = \operatorname{arcosh} \left( \sqrt{1 + \left(\frac{t}{2}\right)^2} + \frac{t}{2} \right) \leq \sqrt{t} \cdot \left(1 + \frac{t}{12}\right), \quad t = |y| < 1;$$

Für  $0 \leq t = |y| < 1$  gelten noch die folgenden Abschätzungen:

$$\begin{aligned} \sqrt{t} \left(1 - \frac{3}{80}t^2\right) &\leq \sqrt{t} \left[1 - t^2 \left(\frac{3}{160} + \frac{3}{160}\right)\right] \leq \sqrt{t} \left[1 - t^2 \left(\frac{3}{160} + \frac{5}{896} \cdot 1\right)\right] \leq \\ &\leq \sqrt{t} \left[1 - t^2 \left(\frac{3}{160} + \frac{5}{896} \cdot t\right)\right] \leq \sqrt{t} \left[\left(1 + \frac{1}{12} \cdot t\right) - t^2 \left(\frac{3}{160} + \frac{5}{896} \cdot t\right)\right] \end{aligned}$$

Mit (C.90) und (C.91) erhalten wir schließlich im Fall  $|x| = 1$ :

$$(C.92) \quad \sqrt{t} \cdot \left(1 - \frac{3}{80}t^2\right) \leq \operatorname{arcosh}(\alpha) \leq \sqrt{t} \cdot \left(1 + \frac{t}{12}\right), \quad t = |y| < 1;$$

Die Auswertung dieser Abschätzung erfolgt in der Hilfsfunktion

```
MpfiClass ACOSH_f_aux(const MpfiClass& x, const MpfiClass& y);
```

unter der Voraussetzung  $(\operatorname{expo}(\operatorname{Inf}(\operatorname{delta})) < -\operatorname{GetCurrPrecision}())$ . Bei hoher Current-Präzision kommt obige Abschätzung also nur bei entsprechend kleinen Werten von  $\delta := |y|$  zur Anwendung, um möglichst optimale Einschließungen zu erhalten.

Wir betrachten jetzt nach (C.85) den Fall:  $|x| < 1$ , mit  $\delta := V(x, y)/2$ , wobei  $V(x, y)$  definiert war durch:

$$V(x, y) := |x + 1| \cdot \left\{ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right\} + |x - 1| \left\{ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right\}.$$

Für die Einschließung des Imaginärteils  $\Im(\mathbf{w})$  von  $\mathbf{w} = \arcsin(\mathbf{z})$  ist dann nach (C.82) der Ausdruck  $\operatorname{arcosh}(1 + \delta)$  für das Intervallargument  $\delta$  auszuwerten, wobei in  $V(x, y)$  für  $x, y$  die entsprechenden Punktintervalle einzusetzen sind. Damit scheint das Problem gelöst zu sein, da man für  $y \rightarrow 0$  die Ausdrücke  $\{\dots\}$  mit Hilfe der bereits implementierten Funktion `sqrtp1m1(..)` intervallmäßig auswerten kann. Dabei übersieht man jedoch, dass wegen der Reihenentwicklung

$$(C.93) \quad \sqrt{1+t} - 1 = t \cdot \sum_{k=0}^{\infty} (-1)^{k+1} \cdot \left( \prod_{n=0}^k \frac{2n-1}{2n+2} \right) \cdot t^k$$

$$(C.94) \quad = \frac{t}{2} \cdot \left( 1 - \frac{t}{4} + \frac{t^2}{8} - \frac{5 \cdot t^3}{64} \pm \dots \right), \quad |t| < 1;$$

<sup>11</sup> $|a_k| = |b_k|$  konnte z.B. mit `Maple` bis  $k = 100$  bestätigt werden.

die obigen Ausdrücke  $\{\dots\}$  für  $y \rightarrow 0$  von der Größenordnung

$$\frac{1}{2} \cdot \left( \frac{y}{x \pm 1} \right)^2$$

sind und damit bei der späteren Intervallauswertung von  $\sqrt{V}$  zu einem vorzeitigen Unterlauf führen, womit dann starke Überschätzungen verbunden sind. Diese Überschätzungen kann man zwar vermeiden, indem man  $y$  geeignet skaliert, die entsprechenden Rechnungen sind jedoch aufwendig und lassen sich durch den folgenden sehr einfachen Algorithmus für  $|y| \rightarrow 0$  ohne Anwendung der Funktion `sqrtp1m1(...)` ganz vermeiden.

Um  $V(x, y)$  einzuschließen, berechnen wir zunächst eine **Oberschranke**: Da die Reihe in (C.94) für  $0 \leq t < 1$  eine alternierende Leibniz-Reihe ist, gilt die Abschätzung<sup>12</sup>:

$$(C.95) \quad \frac{t}{2} \cdot \left( 1 - \frac{t}{4} \right) \leq \sqrt{1+t} - 1 \leq \frac{t}{2}, \quad \text{falls } 0 \leq t < 1;$$

Mit

$$v_1 := |x+1| \cdot \left\{ \sqrt{1 + \left( \frac{y}{x+1} \right)^2} - 1 \right\}$$

folgt dann direkt

$$v_1 \leq |x+1| \cdot \frac{1}{2} \cdot \left( \frac{y}{x+1} \right)^2 = \frac{1}{2} \cdot \frac{y^2}{|x+1|}, \quad \text{falls } y^2 < (x+1)^2.$$

Mit

$$v_2 := |x-1| \cdot \left\{ \sqrt{1 + \left( \frac{y}{x-1} \right)^2} - 1 \right\}$$

folgt ganz entsprechend

$$v_2 \leq |x-1| \cdot \frac{1}{2} \cdot \left( \frac{y}{x-1} \right)^2 = \frac{1}{2} \cdot \frac{y^2}{|x-1|}, \quad \text{falls } y^2 < (x-1)^2.$$

Für  $V(x, y) = v_1 + v_2$  erhält man unter den Voraussetzungen  $|y| < |x+1|$ ,  $|y| < |x-1|$  und  $|x| < 1$  die Oberschranke:

$$(C.96) \quad V(x, y) \leq \frac{y^2}{2} \cdot \left[ \frac{1}{|x+1|} + \frac{1}{|x-1|} \right] = \frac{y^2}{1-x^2}.$$

Unter der Voraussetzung  $|x| < 1$  sind die beiden Bedingungen  $|y| < |x+1|$ ,  $|y| < |x-1|$  für jede Current-Präzision  $\text{prec} \geq 2$  erfüllt, wenn gilt:  $|y| < 1 - \text{pred}(1) = 2^{-\text{prec}}$ . Damit erhalten wir

$$(C.97) \quad V(x, y) \leq \frac{y^2}{1-x^2}, \quad \text{falls } |x| < 1 \text{ und } |y| < 2^{-\text{prec}}.$$

Um  $V(x, y)$  einzuschließen, berechnen wir jetzt eine **Unterschranke**. Mit (C.95) folgt jetzt ganz analog:

$$\frac{1}{2} \cdot \frac{y^2}{|x+1|} \cdot \left[ 1 - \frac{1}{4} \left( \frac{y}{x+1} \right)^2 \right] \leq v_1, \quad \text{falls } |y| < 2^{-\text{prec}}.$$

Wir stellen zusätzlich noch folgende Bedingung:

$$1 - \frac{1}{4} \left( \frac{y}{x+1} \right)^2 > \text{pred}(1) = 1 - 2^{-\text{prec}}, \quad \text{d.h.} \quad \left( \frac{y}{x+1} \right)^2 < 2^{-\text{prec}+2}.$$

<sup>12</sup>Diese Abschätzung lässt sich sehr einfach auch direkt beweisen.

Es gilt für  $|x| < 1$  die Abschätzung  $(x+1)^2 \geq (\text{succ}(-1)+1)^2 = 2^{-2\text{prec}}$  und damit:

$$\frac{y^2}{(x+1)^2} \leq \frac{y^2}{2^{-2\text{prec}}}.$$

Die obige Bedingung ist also erfüllt, wenn gilt

$$\frac{y^2}{2^{-2\text{prec}}} < 2^{-\text{prec}+2} \quad \text{bzw. wenn} \quad |y| < 2^{-((3\text{prec})/2)}.$$

Damit erhalten wir dann für  $\text{prec} \geq 2$  die Unterschranke:

$$(C.98) \quad \frac{1}{2} \cdot \frac{y^2}{|x+1|} \cdot \text{pred}(1) \leq v_1, \quad \text{falls } |y| < 2^{-((3\text{prec})/2)}.$$

Beim obigen Exponenten  $(3\text{prec})/2$  ist die spezielle Maschinenrundung des integer-Quotienten in Richtung Null bereits berücksichtigt!

Ganz entsprechend findet man die Abschätzung

$$(C.99) \quad \frac{1}{2} \cdot \frac{y^2}{|x-1|} \cdot \text{pred}(1) \leq v_2, \quad \text{falls } |y| < 2^{-((3\text{prec})/2)}, \quad \text{und damit:}$$

$$(C.100) \quad \frac{y^2}{2} \cdot \left[ \frac{1}{|x+1|} + \frac{1}{|x-1|} \right] \cdot \text{pred}(1) = \frac{y^2}{1-x^2} \cdot \text{pred}(1) \leq v_1 + v_2 = V(x, y);$$

Für  $V(x, y)$  gilt dann zusammen mit (C.96) für  $\text{prec} \geq 2$  die garantierte Einschließung

$$(C.101) \quad \text{pred}(1) \cdot \frac{y^2}{1-x^2} \leq V(x, y) \leq \frac{y^2}{1-x^2}, \quad \text{falls } |y| < 2^{-((3\text{prec})/2)}.$$

Damit haben wir für hinreichend kleine Werte von  $|y|$  eine schon recht einfache Einschließung für  $V(x, y)$  gefunden. Um eine Einschließung auch für  $\text{arcosh}(1+\delta)$  zu erhalten, betrachten wir die Reihenentwicklung in (C.88) auf Seite 333. Im Buch zur Fehlerabschätzung wird gezeigt, dass  $Q(t)$  für  $0 \leq t < 2$  eine alternierende Leibniz-Reihe ist, so dass sich für  $0 \leq \delta < 2$  die folgende Abschätzung ergibt:

$$(C.102) \quad \sqrt{2\delta} \cdot \left(1 - \frac{\delta}{12}\right) \leq \text{arcosh}(1+\delta) \leq \sqrt{2\delta}, \quad \delta = V(x, y)/2 < 2 \quad \rightsquigarrow \\ \sqrt{V(x, y)} \cdot \left(1 - \frac{V(x, y)}{24}\right) \leq \text{arcosh}(1+V(x, y)/2) \leq \sqrt{V(x, y)}, \quad \text{falls } V(x, y) < 4.$$

Wir stellen jetzt noch zusätzlich die Bedingung

$$1 - \frac{V}{24} > \text{pred}(1) = 1 - 2^{-\text{prec}}, \quad \text{d.h. } V(x, y) < 3 \cdot 2^{-\text{prec}+3}$$

und erhalten damit:

$$(C.103) \quad \sqrt{V} \cdot \text{pred}(1) \leq \text{arcosh}(1+V/2) \leq \sqrt{V}, \quad \text{falls } V(x, y) < 3 \cdot 2^{-\text{prec}+3}.$$

Mit (C.101) und (C.103) kann jetzt unter der Voraussetzung  $|y| < 2^{-((3\text{prec})/2)}$  für  $\text{prec} \geq 2$  eine sehr enge Einschließung für  $\text{arcosh}(1+V/2)$  angegeben werden<sup>13</sup>:

$$(C.104) \quad \text{pred}(1) \cdot \sqrt{\text{pred}(1)} \cdot \frac{|y|}{\sqrt{1-x^2}} \leq \text{arcosh}(1+V(x, y)/2) \leq \frac{|y|}{\sqrt{1-x^2}}.$$

<sup>13</sup>Man kann mit (C.101) leicht zeigen, dass unter der Voraussetzung  $|y| < 2^{-((3\text{prec})/2)}$  neben  $V(x, y) < 4$  auch die Bedingung  $V(x, y) < 3 \cdot 2^{-\text{prec}+3}$  für alle  $\text{prec} \geq 2$  erfüllt ist.



Wegen  $\sqrt{\text{pred}(1)} > \text{pred}(1)$  gilt im Fall  $|y| < 2^{-(3\text{prec})/2}$  für  $\text{prec} \geq 2$  auch die etwas gröbere Einschließung:

$$(C.105) \quad \text{pred}(1) \cdot \text{pred}(1) \cdot \frac{|y|}{\sqrt{1-x^2}} \leq \text{arcosh}(1 + V(x, y)/2) \leq \frac{|y|}{\sqrt{1-x^2}}.$$

Beachten Sie, dass jetzt für  $|x| < 1$  bei der Auswertung von  $|y|/\sqrt{1-x^2}$  ein vorzeitiger Unterlauf nicht mehr eintreten kann, wobei  $\sqrt{1-x^2}$  mit Hilfe von `sqrt_1mx2(x)` intervallmäßig auszuwerten ist. Ein Überlauf kann bei der Auswertung von  $|y|/\sqrt{1-x^2}$  für  $|x| < 1$  nur eintreten, wenn die Präzision sehr groß und damit  $\sqrt{1-x^2}$  sehr klein wird. Diese sehr großen Präzisionen haben aber wegen der damit verbundenen extremen Laufzeiten keinerlei praktische Bedeutung!

Um für  $\text{arcosh}(1 + V(x, y)/2)$  eine möglichst einfach auszuwertende Einschließung angeben zu können, beweisen wir vorher noch den folgenden Satz:

Für alle Maschinenzahlen  $x \geq 0$  des MPFR-Formats gilt:

$$(C.106) \quad \text{pred}(1) \cdot x \geq \text{pred}(x).$$

Der **Beweis** für  $x = 0$  und  $x = \text{minfloat}()$  ist trivial. Wir betrachten jetzt positive, normalisierte Zahlen  $x = m \cdot 2^{ex} > \text{minfloat}()$ , mit  $0.5 \leq m < 1$ . Nach (3.1) auf Seite 17 gilt:

$$\begin{aligned} m = 0.5 &\rightsquigarrow \text{pred}(x) = 2^{ex-1} \cdot (1 - 2^{-\text{prec}}) \\ 0.5 < m < 1 &\rightsquigarrow \text{pred}(x) = 2^{ex} \cdot (m - 2^{-\text{prec}}). \end{aligned}$$

Zum Beweis von (C.106) sind damit zwei Fälle zu unterscheiden. Für  $m = 0.5$  kann (C.106) wie folgt äquivalent umgeformt werden:

$$\begin{aligned} &\text{pred}(1) \cdot x \geq \text{pred}(x) \\ \iff &\text{pred}(1) \cdot 0.5 \cdot 2^{ex} \geq 2^{ex-1} \cdot (1 - 2^{-\text{prec}}), \\ \iff &(1 - 2^{-\text{prec}}) \cdot 2^{ex-1} \geq 2^{ex-1} \cdot (1 - 2^{-\text{prec}}). \end{aligned}$$

Für  $0.5 < m < 1$  kann (C.106) wie folgt äquivalent umgeformt werden:

$$\begin{aligned} &\text{pred}(1) \cdot x \geq \text{pred}(x) \\ \iff &\text{pred}(1) \cdot m \cdot 2^{ex} \geq 2^{ex} \cdot (m - 2^{-\text{prec}}) \\ \iff &(1 - 2^{-\text{prec}}) \cdot m \cdot 2^{ex} \geq 2^{ex} \cdot (m - 2^{-\text{prec}}) \\ \iff &(1 - 2^{-\text{prec}}) \cdot m \geq m - 2^{-\text{prec}} \\ \iff &-2^{-\text{prec}} \geq -2^{-\text{prec}} \quad \blacksquare \end{aligned}$$

Um jetzt mit (C.105) unter der Voraussetzung  $|y| < 2^{-(3\text{prec})/2}$  eine Einschließung des Ausdrucks  $\text{arcosh}(1 + \delta) = \text{arcosh}(1 + V(x, y)/2)$  zu erhalten, berechnen wir zunächst eine intervallmäßige Einschließung des Quotienten

$$|y|/\sqrt{1-x^2} \in \text{abs}(y) \diamond \text{sqrt}_1\text{mx2}(x) = [u1, u2],$$

wobei das Argument  $x = [x, x]$  in der Funktion `sqrt_1mx2(x)` als Punktintervall aufzufassen ist. Mit  $u2$  erhalten wir damit auch eine Oberschranke von  $A := \text{arcosh}(1 + V(x, y)/2) \leq u2$ .

Um für  $A$  eine Unterschranke zu berechnen, gilt mit dem Satz (C.106)

$$\begin{aligned} \text{pred}(1) \cdot \{\text{pred}(1) \cdot |y|/\sqrt{x^2-1}\} &\geq \text{pred}(1) \cdot \{\text{pred}(1) \cdot u1\} \geq \text{pred}(1) \cdot \{\text{pred}(u1)\} \\ &\geq \text{pred}\{\text{pred}(u1)\}. \end{aligned}$$

Unter der Voraussetzung  $|y| < 2^{-(3\text{prec})/2}$  erhalten wir schließlich mit (C.105) für  $\text{prec} \geq 2$  die sehr effektive Einschließung

$$(C.107) \quad \text{pred}\{\text{pred}(u1)\} \leq \text{arcosh}(1 + V(x, y)/2) \leq u2.$$

Mit  $y = m \cdot 2^{ex}$ ,  $ex = \text{expo}(y)$  und  $0.5 \leq m < 1$  ist die Voraussetzung  $|y| < 2^{-(3\text{prec})/2}$  erfüllt, wenn gilt:  $ex < -(3 \cdot \text{prec})/2$ , dabei ist die Maschinenrundung zur Null bei der integer-Division  $(3 \cdot \text{prec})/2$  schon berücksichtigt. Die Auswertung von (C.107) erfolgt in der Hilfsfunktion

```
MpfiClass ACOSH_p1(const MpfiClass& x, const MpfiClass& y);
```

die nur für Punktintervalle auszuwerten ist.

Die beiden nächsten numerischen Beispiele zeigen, wie genau die Einschließungen mit (C.107) im Vergleich zu (C.102) ausfallen. Dazu wählen wir die Current-Präzision  $\text{prec} = 70$ . Mit den beiden Maschinenzahlen  $x_1 = \text{pred}(1) = 1 - 2^{-70}$  und  $y_1 = \text{minfloat}() = 2^{-1073741824}$  ist das Eingangsintervall gegeben durch das Punktintervall  $Z = [x_1, x_1] + i \cdot [y_1, y_1]$ .

Mit (C.102) erhält man für den Imaginärteil die praktisch unbrauchbare Einschließung

$$\text{asin}(Z) = ([1.57079632675373758748, 1.57079632675373758749], \\ [0.00000000000000000000, 8.45440400895524581723e - 161614249]).$$

Dagegen erhält man mit (C.107) für den Imaginärteil die fast optimale Einschließung

$$\text{asin}(Z) = ([1.57079632675373758748, 1.57079632675373758749], \\ [5.78868064589607735286e - 323228487, 5.78868064589607735289e - 323228487]).$$



### C.2.21.1. Algorithmus

Der in [54] beschriebene Algorithmus verwendet die folgende Identität:

$$(C.108) \quad \arccos(z) \equiv \frac{\pi}{2} - \arcsin(z), \quad z = x + i \cdot y;$$

Mit  $w = \arccos(z) = \Re(w) + i \cdot \Im(w)$  wird  $\Im(w)$  daher über den inversen Sinus optimal berechnet. Im Falle  $\Re(w) \sim 0$  bzw.  $\Re(\arcsin(z)) \sim \pi/2$  wird jedoch  $\Re(w)$  nach (C.108) wegen starker Auslöschung sehr fehlerhaft berechnet, so dass  $\Re(w)$  damit nur grob eingeschlossen werden kann. Als Beispiel betrachten wir  $z = 1 + i \cdot 2^{-100}$ . Für das entsprechende Punktintervall  $\mathbf{z}$  erhalten wir mit (C.108) nur die sehr grobe Einschließung

$$\arccos(\mathbf{z}) \subset ([-8.437695E-015, 2.980234E-008], [-2.980233E-008, -0.000000E+000])$$

während z.B. *Maple7* das folgende Ergebnis liefert:

$$\arccos(\mathbf{z}) = 8.881784197001 \dots \cdot 10^{-16} - i \cdot 8.881784197001 \dots \cdot 10^{-16}$$

Um diese Überschätzung zu vermeiden, wird wie in [38] bei vorgegebenem komplexen Argumentintervall

$$Z = X + i \cdot Y = [x_1, x_2] + i \cdot [y_1, y_2], \quad x_i, y_i \text{ sind Maschinenzahlen}$$

zur Einschließung von  $\Re(\arccos(Z))$  die reelle  $\arccos$ -Funktion für den Intervall-Ausdruck  $\beta$  ausgewertet:

$$(C.109) \quad \arccos(\beta), \quad \beta := \frac{2\mathbf{x}}{\sqrt{(\mathbf{x}+1)^2 + \mathbf{y}^2} + \sqrt{(\mathbf{x}-1)^2 + \mathbf{y}^2}},$$

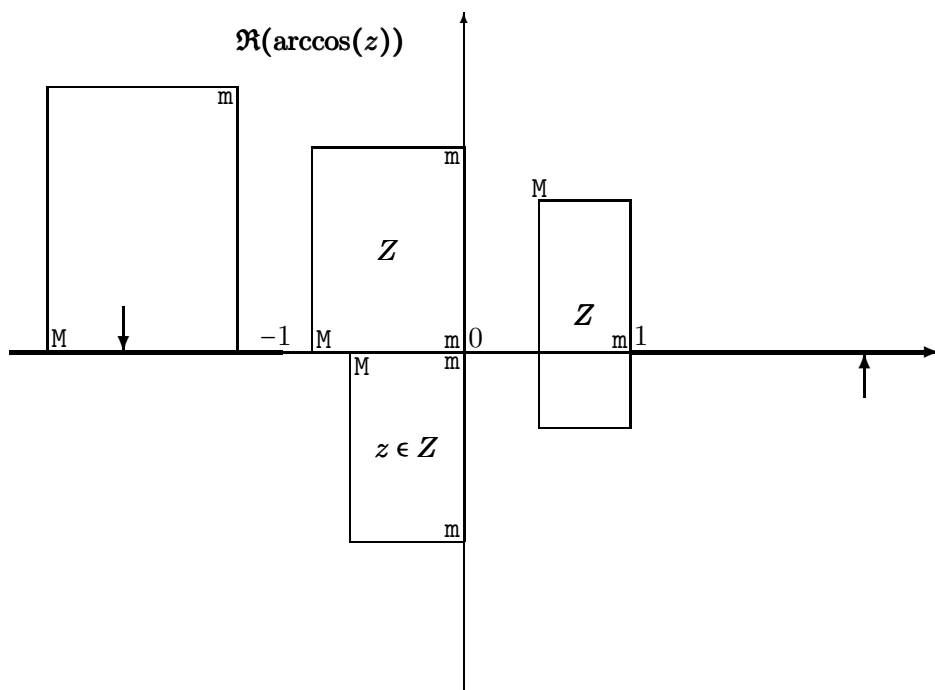
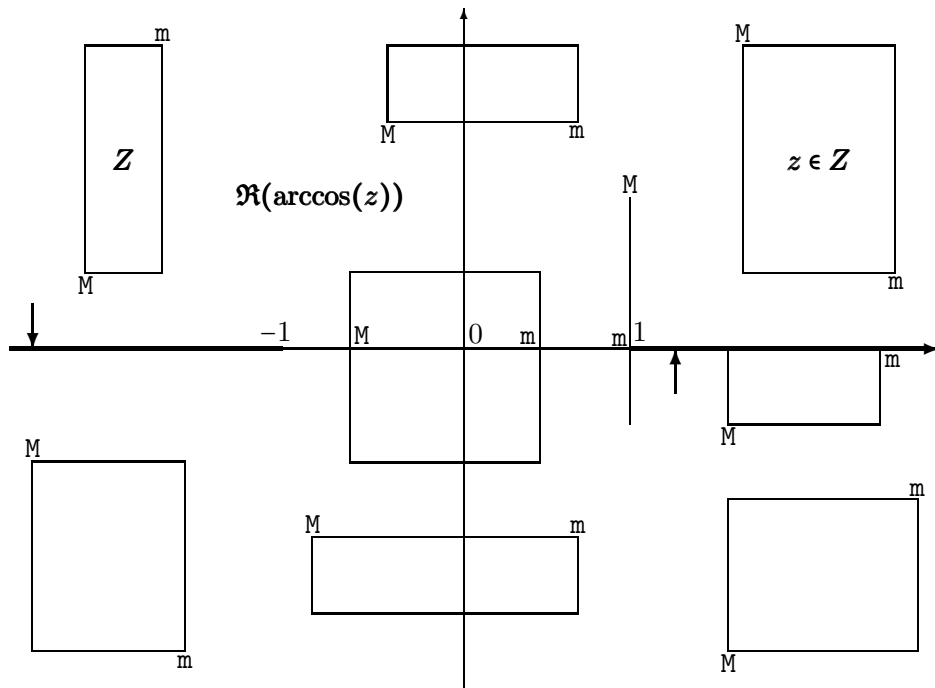
wobei die Punktintervalle  $\mathbf{x}, \mathbf{y}$  durch die Extremalpunkte auf dem Rand von  $Z$  bestimmt sind, vgl. [38, S. 135].

In den folgenden Abbildungen sind zu einigen Intervallen  $Z$ , die in der komplexen Ebene als Rechtecke dargestellt werden, die Extremalpunkte angegeben. Dabei sind die Koordinaten von  $\mathbf{m}$  die Koordinaten des Extremalpunktes, an dem das Minimum von  $\Re(w)$  angenommen wird. Entsprechend sind die Koordinaten von  $\mathbf{M}$  die Koordinaten des Extremalpunktes, an dem das Maximum von  $\Re(w)$  angenommen wird. Man erkennt, dass  $\mathbf{m}$  und  $\mathbf{M}$  fast immer die Eckpunkte von  $Z$  sind. Nur wenn der Koordinatenursprung im Innern von  $Z$  enthalten ist, liegen die Extremalpunkte  $\mathbf{m}, \mathbf{M}$  auf der reellen Achse, d.h. also nicht auf den Eckpunkten von  $Z$ .

Einige Besonderheiten liegen vor, wenn  $Z$  z.B. in der rechten Halbebene liegt und die reelle Achse durch das Innere von  $Z$  läuft. Dann liegt  $\mathbf{M}$  auf demjenigen linken Eckpunkt von  $Z$ , der die betragsmäßig maximale  $y$ -Koordinate besitzt. Liegt jedoch  $Z$  in der linken Halbebene und läuft die reelle Achse wieder durch das Innere von  $Z$ , dann liegt  $\mathbf{m}$  auf demjenigen rechten Eckpunkt von  $Z$ , der die betragsmäßig maximale  $y$ -Koordinate besitzt.

Zum Verständnis der Abbildung ist weiter zu beachten, dass im Fall zweier Punkte  $\mathbf{m}$  auf dem Rand von  $Z$  auch alle Zwischenpunkte Extremalpunkte sind, auf denen  $\Re(w)$  sein Minimum über  $Z$  annimmt. Entsprechendes gilt im Falle zweier verschiedener Punkte  $\mathbf{M}$  auf dem Rand von  $Z$ .

In den folgenden Abbildungen sind einige mögliche Lagen der Intervalle  $Z$  angegeben. Die Pfeile auf die Verzweigungsschnitte beschreiben, aus welcher Richtung die Funktionswerte des Hauptzweiges auf die Schnitte analytisch fortzusetzen sind.



Nach Seite 330 gilt  $\text{Sup}(\beta) \leq 1$ . Für  $\text{Inf}(\beta) \rightarrow 1$  wird daher die reelle arccos-Funktion in der Nähe ihrer Nullstelle  $x_0 = +1$  ausgewertet, so dass bei der Einschließung von  $\arccos(\beta)$  große Überschätzungen zu erwarten sind. Mit dem komplexen Argument  $z = 1 + i \cdot 2^{-100}$  erhält man in der Tat bei Anwendung von (C.109) für den Realteil die sehr grobe Einschließung:

$$\arccos(z) \in ([0.000000E+000, 2.980233E-008], [-8.881785E-016, -8.881784E-016])$$

Diese Überschätzungen lassen sich vermeiden, wenn man mit Hilfe der Transformation

$$(C.110) \quad \beta := \frac{2x}{\sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2}} = 1 - \delta, \quad \delta \geq 0$$

im Falle  $\beta \geq 0.75$  für  $\arccos(\beta)$  die folgende Darstellung benutzt:

$$(C.111) \quad \arccos(\beta) = \arccos(1 - \delta) = \arcsin(\sqrt{2\delta - \delta^2})$$

Bevor wir für die Fälle:  $x < 1$ ,  $x = 1$ ,  $x > 1$  zur Auswertung von  $\delta(x, y)$  die jeweils geeigneten Terme angeben, soll zunächst in der komplexen Ebene der Bereich bestimmt werden, in dem  $\beta \geq 0.75$  gilt. Aus (C.110) folgt unmittelbar, dass  $\beta(x, y)$  für festes  $x > 0$  sein relatives Maximum bez.  $y$  auf der reellen Achse, d.h. für  $y = 0$  annimmt. Aus (C.110) erhält man für diese Maximumwerte direkt:

$$(C.112) \quad \beta(x, 0) = \begin{cases} x & \text{falls } 0 \leq x \leq 1 \\ 1 & \text{falls } x > 1 \quad \rightsquigarrow \quad \delta = 0; \end{cases}$$

Dass für  $x > 0$  die obigen Werte tatsächlich relative Extrema bez.  $y$  sind, ergibt sich auch direkt aus der nachfolgenden partiellen Ableitung  $\partial\beta(x, y)/\partial y$ , die für  $y = 0$  verschwindet.

$$\frac{\partial\beta(x, y)}{\partial y} = -2xy \cdot \frac{\frac{1}{\sqrt{(x+1)^2 + y^2}} + \frac{1}{\sqrt{(x-1)^2 + y^2}}}{\left(\sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2}\right)^2}$$

Da die partielle Ableitung in der rechten Halbebene für  $y > 0$  negativ und für  $y < 0$  positiv ist, liegt in der komplexen Ebene der gesuchte Bereich, in dem  $\beta(x, y) \geq 0.75$  gilt, innerhalb der Höhenlinie

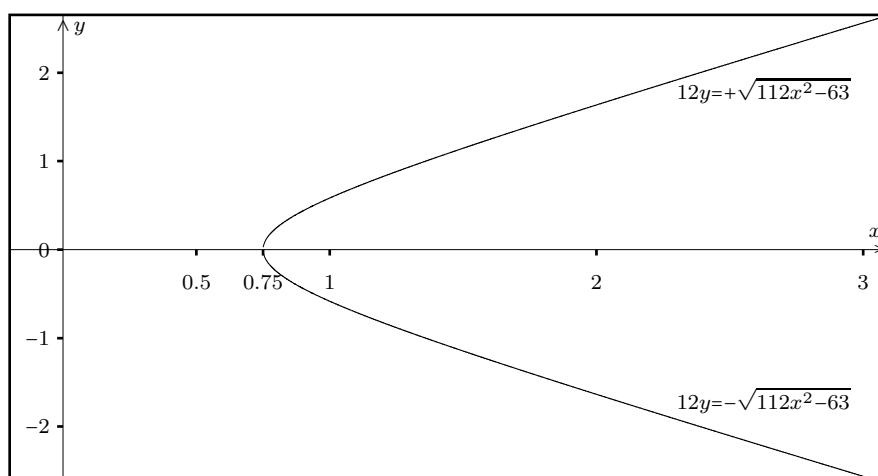
$$\beta(x, y) = \frac{3}{4} \quad \Longleftrightarrow \quad |y| = \frac{\sqrt{112x^2 - 63}}{12};$$

wobei die positive reelle Achse als Symmetrieachse durch diesen Bereich läuft. Beachten Sie bitte, dass durch die Funktionsgleichungen

$$y = \pm \frac{\sqrt{112x^2 - 63}}{12}$$

eine Hyperbel mit ihrem rechten Hyperbelast beschrieben wird, dessen Scheitelpunkt durch  $x = 3/4$  und  $y = 0$  gegeben ist. Der gesuchte Bereich  $\beta(x, y) \geq 0.75$  ist in der folgenden Abbildung innerhalb dieses Hyperbelastes dargestellt:

Bereich in der komplexen Ebene bez.  $\beta \geq 0.75$



Um (C.111) im obigen Bereich  $\beta(x, y) \geq 0.75$  anwenden zu können, benötigt man zunächst eine Darstellung für  $\delta = \delta(x, y)$ . Man findet direkt

$$(C.113) \quad 2\delta = 2 - \sqrt{(x+1)^2 + y^2} + \sqrt{(x-1)^2 + y^2} \geq 0.$$

Aber (C.113) ist offensichtlich für  $\delta \rightarrow 0$  wegen starker Auslöschung zur Auswertung auf der Maschine nicht geeignet. Wir betrachten daher unter der Voraussetzung  $\beta \geq 0.75$  zunächst den Fall  $x > 1$ :

Wegen  $y = 0 \rightsquigarrow \delta = 0$ , vgl. (C.112), kann man sich jetzt auf  $y \neq 0$  beschränken. Aus (C.113) ergibt sich zunächst:

$$2\delta = -(x+1) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right\} + (x-1) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right\}$$

Die beiden obigen Summanden liefern jedoch für  $x \rightarrow +\infty$  ebenfalls starke Auslöschung, so dass weitere Umformungen nötig sind:

$$\begin{aligned} & -x \cdot \left\{ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right\} + x \cdot \left\{ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right\} \\ &= -x \cdot \sqrt{1 + \left(\frac{y}{x+1}\right)^2} + x \cdot \sqrt{1 + \left(\frac{y}{x-1}\right)^2} \\ &= x \cdot \frac{1 + \left(\frac{y}{x-1}\right)^2 - 1 - \left(\frac{y}{x+1}\right)^2}{\sqrt{1 + \left(\frac{y}{x+1}\right)^2} + \sqrt{1 + \left(\frac{y}{x-1}\right)^2}} \\ &= \frac{4y^2 \cdot x^2}{(x^2 - 1)^2 \cdot N}, \quad N := \sqrt{1 + \left(\frac{y}{x+1}\right)^2} + \sqrt{1 + \left(\frac{y}{x-1}\right)^2}; \end{aligned}$$

Für  $2\delta$  erhält man damit:

$$\begin{aligned} 2\delta &= \frac{4y^2 \cdot x^2}{(x^2 - 1)^2 \cdot N} + 2 - N = \frac{4y^2 \cdot x^2 + (2 - N)(x^2 - 1)^2 \cdot N}{(x^2 - 1)^2 \cdot N} \\ &= \frac{Z}{(x^2 - 1)^2 \cdot N}, \quad Z = 4y^2 \cdot x^2 + (2 - N)(x^2 - 1)^2 \cdot N; \\ & \quad Z = -2 \cdot (x^2 - 1) \cdot F; \end{aligned}$$

$$\begin{aligned} F &= x^2 - y^2 - 1 + \sqrt{1 + \left(\frac{y}{x+1}\right)^2} \cdot \sqrt{1 + \left(\frac{y}{x-1}\right)^2} \cdot (x^2 - 1) - \\ & \quad - \left( \sqrt{1 + \left(\frac{y}{x+1}\right)^2} + \sqrt{1 + \left(\frac{y}{x-1}\right)^2} \right) \cdot (x^2 - 1) \\ &= (x^2 - 1) \cdot \left( \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right) \cdot \left( \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right) - y^2 \end{aligned}$$

Für  $\delta = \delta(x, y)$  ergibt sich damit schließlich im Fall  $x > 1$ :

$$(C.114) \quad 2\delta = \frac{\frac{y^2}{x^2 - 1} - \left[ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right] \cdot \left[ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right]}{1 + \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right) + \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right)}$$

Der Nenner in (C.114) lässt sich zwar noch weiter vereinfachen zu

$$\frac{1}{2} \cdot \sqrt{1 + \left(\frac{y}{x-1}\right)^2} + \frac{1}{2} \cdot \sqrt{1 + \left(\frac{y}{x+1}\right)^2},$$

die Darstellung in (C.114) hat jedoch den Vorteil, dass zur Berechnung von  $\delta(x, y)$  im Wesentlichen nur zwei Terme mit Hilfe der Funktion `sqrtp1m1()` auszuwerten sind. Außerdem kann  $2\delta$  für  $y \rightarrow 0$  mit Hilfe von (C.114) sehr einfach eingeschlossen werden. Für eine Oberschranke erhält man aus (C.114) direkt:

$$(C.115) \quad 2 \cdot \delta(x, y) \leq \frac{y^2}{x^2 - 1}.$$

Die Berechnung einer Unterschranke für  $2 \cdot \delta(x, y)$  ist etwas komplizierter. Nach (C.114) schreiben wir  $2\delta = Z/N$  und bestimmen zunächst eine Oberschranke des Nenners  $N$ . Mit Hilfe der Reihenentwicklung für  $\sqrt{1+t} - 1$  von Seite 334 gelten die Abschätzungen:

$$\begin{aligned} \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right) &\leq \frac{1}{4} \cdot \left(\frac{y}{x+1}\right)^2, \quad \text{falls } \frac{|y|}{x+1} < 1 \\ \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right) &\leq \frac{1}{4} \cdot \left(\frac{y}{x-1}\right)^2, \quad \text{falls } \frac{|y|}{x-1} < 1. \end{aligned}$$

Für den Nenner  $N$  folgt damit:

$$\begin{aligned} N &\leq 1 + \frac{y^2}{4} \cdot \left[ \frac{1}{(x+1)^2} + \frac{1}{(x-1)^2} \right] \\ &= 1 + \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2}, \quad \text{falls } \frac{|y|}{x-1} < 1. \end{aligned}$$

Beachten Sie bitte, dass wegen  $x > 1$  die Bedingung  $|y| < x + 1$  automatisch erfüllt ist, falls  $|y| < x - 1$  gilt. Für den Kehrwert von  $N$  erhalten wir:

$$\frac{1}{N} \geq \frac{1}{1 + \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2}}, \quad \text{falls } \frac{|y|}{x-1} < 1.$$

Da die Reihe

$$\frac{1}{1 + \varepsilon} = 1 - \varepsilon + \varepsilon^2 - \varepsilon^3 \pm \dots$$

für  $0 \leq \varepsilon < 1$  eine alternierende Leibniz-Reihe ist, folgt direkt:

$$(C.116) \quad \frac{1}{N} \geq 1 - \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2}, \quad \text{falls } \frac{|y|}{x-1} < 1 \wedge \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2} < 1.$$

Wir zeigen jetzt, dass wegen unserer Voraussetzung  $x > 1$  mit  $|y| < x - 1$  die obige zweite Bedingung automatisch erfüllt ist. Es gilt

$$\begin{aligned} \frac{x^2 + 1}{(x+1)^2} &< \frac{x^2 + 1}{x^2} = 1 + \frac{1}{x^2} < 2 \quad \text{und damit:} \\ \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2} &= \frac{y^2}{2(x-1)^2} \cdot \frac{x^2 + 1}{(x+1)^2} < \frac{y^2}{(x-1)^2}, \end{aligned}$$

d.h. die obige zweite Bedingung ist erfüllt, wenn

$$\frac{y^2}{(x-1)^2} < 1 \quad \iff \quad \frac{|y|}{x-1} < 1,$$



und die letzte Ungleichung ist gerade die erste Bedingung in (C.116). Es gilt damit:

$$(C.117) \quad \frac{1}{N} \geq 1 - \frac{y^2}{2} \cdot \frac{x^2 + 1}{(x^2 - 1)^2}, \quad \text{falls } \frac{|y|}{x-1} < 1.$$

Wir müssen jetzt noch eine Unterschranke des Zählers  $Z$  in (C.114) bestimmen. Zunächst gilt wieder die Abschätzung:

$$\begin{aligned} [\dots] \cdot [\dots] &\leq \frac{1}{4} \cdot \left(\frac{y}{x-1}\right)^2 \cdot \left(\frac{y}{x+1}\right)^2 \\ &= \frac{1}{4} \cdot \frac{y^4}{(x^2-1)^2}, \quad \text{falls } \frac{|y|}{x-1} < 1, \end{aligned}$$

und damit erhalten wir:

$$\begin{aligned} Z &\geq \frac{y^2}{x^2-1} - \frac{1}{4} \cdot \frac{y^4}{(x^2-1)^2} \\ &= \frac{y^2}{x^2-1} \cdot \left[1 - \frac{y^2}{4 \cdot (x^2-1)}\right], \quad \text{falls } \frac{|y|}{x-1} < 1. \end{aligned}$$

Zusammen mit (C.117) ergibt sich für  $2 \cdot \delta(x, y)$  die Abschätzung:

$$(C.118) \quad 2\delta \geq \frac{y^2}{x^2-1} \cdot \left\{1 - \frac{y^2}{4(x^2-1)}\right\} \cdot \left\{1 - \frac{y^2}{2} \cdot \frac{x^2+1}{(x^2-1)^2}\right\}, \quad \text{falls } \frac{|y|}{x-1} < 1.$$

$$(C.119) \quad \begin{aligned} \{\dots\} \cdot \{\dots\} &\geq 1 - \frac{y^2}{4(x^2-1)} - \frac{y^2}{2} \cdot \frac{x^2+1}{(x^2-1)^2} \\ &= 1 - \frac{y^2}{2(x^2-1)} \left(\frac{1}{2} + \frac{x^2+1}{x^2-1}\right). \end{aligned}$$

$$\begin{aligned} \left(\frac{1}{2} + \frac{x^2+1}{x^2-1}\right) &= \frac{3x^2+1}{2(x+1)} \cdot \frac{1}{x-1} \leq \frac{1}{x-1} \cdot \frac{3x^2+1}{2x} \\ &= \frac{1}{x-1} \cdot \left[\frac{3}{2} \cdot x + \frac{1}{2x}\right] \leq \frac{1}{x-1} \cdot \left[\frac{3}{2} \cdot x + \frac{1}{2}\right] \\ &= \frac{3x+1}{2(x-1)} \quad \rightsquigarrow \end{aligned}$$

$$\{\dots\} \cdot \{\dots\} \geq 1 - \frac{y^2}{4(x-1)^2} \cdot \frac{3x+1}{x+1}.$$

Wegen  $x > 1$  gilt noch

$$\begin{aligned} \frac{3x+1}{x+1} &\leq \frac{3x+1}{x} = 3 + \frac{1}{x} < 4 \quad \text{und damit} \\ \{\dots\} \cdot \{\dots\} &\geq 1 - \frac{y^2}{(x-1)^2}. \end{aligned}$$

Zusammen mit (C.115) und (C.118) erhalten wir unter der Voraussetzung  $x > 1$  die Einschließung:

$$(C.120) \quad \left[1 - \frac{y^2}{(x-1)^2}\right] \cdot \frac{y^2}{x^2-1} \leq 2 \cdot \delta(x, y) \leq \frac{y^2}{x^2-1}, \quad \text{falls } \frac{|y|}{x-1} < 1.$$

Wir verlangen zusätzlich noch

$$1 - \frac{y^2}{(x-1)^2} > \text{pred}(1) = 1 - 2^{-\text{prec}} \iff \frac{|y|}{x-1} < \sqrt{2^{-\text{prec}}}$$

und erhalten damit im Falle  $x > 1$  für  $2 \cdot \delta(x, y)$  die gesuchte Einschließung:

$$(C.121) \quad \text{pred}(1) \cdot \frac{y^2}{x^2-1} \leq 2 \cdot \delta(x, y) \leq \frac{y^2}{x^2-1}, \quad \text{falls } \frac{|y|}{x-1} < \sqrt{2^{-\text{prec}}};$$

Mit Hilfe der gefundenen Doppelungleichung (C.121) kann man unter den Voraussetzungen  $x > 1$  und  $|y|/(x-1) \ll 1$  eine sehr effektive Einschließung für  $2 \cdot \delta$  und damit auch für  $\arccos(1-\delta) \approx \sqrt{2\delta}$  berechnen. Im nächsten Schritt wird eine Einschließung für  $\arccos(1-\delta)$  angegeben. Nach Abramowitz gilt, [1, S. 81]:

$$\begin{aligned} \arccos(1-\delta) &= \sqrt{2\delta} \cdot \left[ 1 + \sum_{k=1}^{\infty} \frac{1 \cdot 3 \cdot 5 \cdots (2k-1)}{2^{2k} (2k+1) k!} \cdot \delta^k \right] \\ &= \sqrt{2\delta} \cdot \left[ 1 + \frac{1}{12} \delta + \frac{3}{160} \delta^2 + \frac{5}{896} \delta^3 + \dots \right], \quad 0 \leq \delta < 2; \end{aligned}$$

dabei sind die Taylorkoeffizienten bis aufs Vorzeichen identisch mit den Taylorkoeffizienten der Reihe für  $\text{arcosh}(1+\delta)$ , so dass die vorliegende Reihe mit Hilfe der geometrischen Reihe abgeschätzt werden kann. Man findet unmittelbar:

$$(C.122) \quad \sqrt{2\delta} \leq \arccos(1-\delta) \leq \sqrt{2\delta} \cdot \frac{1}{1-\delta}, \quad 0 \leq \delta < 1.$$

Wir verlangen noch zusätzlich

$$\frac{1}{1-\delta} < \text{succ}(1) = 1 + 2^{-\text{prec}+1} \iff \delta < \frac{2^{-\text{prec}+1}}{1 + 2^{-\text{prec}+1}}$$

und die letzte Ungleichung ist erfüllt, wenn  $\delta < 2^{-\text{prec}}$ , d.h. es gilt

$$(C.123) \quad \sqrt{2\delta} \leq \arccos(1-\delta) \leq \sqrt{2\delta} \cdot \text{succ}(1), \quad \text{falls } \delta < 2^{-\text{prec}}.$$

Zusammen mit (C.121) ergibt sich daraus:

$$(C.124) \quad \sqrt{\text{pred}(1)} \cdot \frac{|y|}{\sqrt{x^2-1}} \leq \arccos(1-\delta) \leq \frac{|y|}{\sqrt{x^2-1}} \cdot \text{succ}(1).$$

Die bisherigen Bedingungen für (C.124) lauten:

$$(C.125) \quad x > 1$$

$$(C.126) \quad \frac{|y|}{x-1} < \sqrt{2^{-\text{prec}}}$$

$$(C.127) \quad \delta < 2^{-\text{prec}}.$$

Wir zeigen jetzt, dass mit (C.125) und (C.126) die Bedingung (C.127) automatisch erfüllt ist. Nach (C.121) gilt mit (C.126) die Abschätzung:

$$\begin{aligned} \delta &\leq \frac{y^2}{2(x^2-1)}, \quad \text{und damit} \\ \delta &\leq \frac{y^2}{(x-1)^2} \cdot \frac{x-1}{2(x+1)}, \quad \text{und wegen } \frac{x-1}{2(x+1)} < \frac{1}{2} \quad \text{folgt} \\ \delta &< 2^{-\text{prec}} \cdot \frac{1}{2} = 2^{-\text{prec}-1} < 2^{-\text{prec}} \quad \blacksquare \end{aligned}$$

Für  $\arccos(1 - \delta)$  gilt damit im Falle  $x > 1$  und  $|y|/(x - 1) < \sqrt{2^{-\text{prec}}}$  die Einschließung

$$(C.128) \quad \sqrt{\text{pred}(1)} \cdot \frac{|y|}{\sqrt{x^2 - 1}} \leq \arccos(1 - \delta) \leq \frac{|y|}{\sqrt{x^2 - 1}} \cdot \text{succ}(1).$$

Um für  $\arccos(1 - \delta)$  eine möglichst einfach auszuwertende Einschließung angeben zu können, beweisen wir jetzt noch den folgenden **Satz**:

Für alle Maschinenzahlen  $x \geq 0$  gilt:

$$(C.129) \quad x \cdot \text{succ}(1) \leq \text{succ}(\text{succ}((x))).$$

Der **Beweis** für  $x = 0$  ist trivial. Wir betrachten daher jetzt positive Zahlen  $x = m \cdot 2^{ex}$ , mit  $0.5 \leq m < 1$ . Nach (3.2) von Seite 17 gilt  $\text{succ}(x) = x + 2^{ex - \text{prec}}$  und daraus folgt direkt:

$$\begin{aligned} x \cdot \text{succ}(1) \leq \text{succ}(\text{succ}((x))) &\iff m \cdot 2^{ex} \cdot (1 + 2^{1 - \text{prec}}) \leq \text{succ}[x + 2^{-\text{prec} + ex}] = \\ &= x + 2^{-\text{prec} + ex} + 2^{-\text{prec} + \text{expo}[\dots]}. \text{ Zu zeigen ist also} \\ x \cdot (1 + 2^{1 - \text{prec}}) &\leq x + 2^{-\text{prec} + ex} + 2^{-\text{prec} + \text{expo}[\dots]} \\ \iff x \cdot 2^{1 - \text{prec}} &\leq 2^{-\text{prec} + ex} + 2^{-\text{prec} + \text{expo}[\dots]} \\ \iff m \cdot 2^{ex - \text{prec} + 1} &\leq 2^{-\text{prec} + ex} + 2^{-\text{prec} + \text{expo}[\dots]} \\ \iff m \cdot 2^{ex + 1} &\leq 2^{ex} + 2^{\text{expo}[\dots]}. \end{aligned}$$

Mit den zusätzlichen Abschätzungen:  $m \cdot 2^{ex + 1} < 2^{ex + 1}$  und  $2^{ex} + 2^{\text{expo}[\dots]} \geq 2^{ex} + 2^{ex} = 2^{ex} + 1$  bleibt also zu zeigen  $2^{ex + 1} \leq 2^{ex + 1}$ . ■

Man beweist sehr einfach  $\sqrt{\text{pred}(1)} > \text{pred}(1)$ , und mit dem **Satz**

$$\text{pred}(1) \cdot x \geq \text{pred}(x), \quad \text{falls } x \geq 0$$

von Seite 337 folgt zusammen mit (C.128) zur Einschließung von  $\arccos(1 - \delta)$  unter der Bedingung  $|y|/(x - 1) < \sqrt{2^{-\text{prec}}}$  der folgende sehr einfache **Algorithmus**:

- Berechne im Falle  $|y|/(x - 1) < \sqrt{2^{-\text{prec}}}$  und  $x > 1$  mit  $u = \text{abs}(y)/\text{sqr}tx2m1(x)$  eine garantierte Einschließung für:  $|y|/\sqrt{x^2 - 1} \subseteq u := [u1, u2]$ ;
- Es gilt dann:  $\text{pred}(u1) \leq \arccos(1 - \delta) \leq \text{succ}(\text{succ}(u2))$ .

Beachten Sie, dass bei Anwendung der Funktion  $\text{sqr}tx2m1(\dots)$  ein vorzeitiger Overflow bei der Berechnung von  $\sqrt{x^2 - 1}$  nicht eintreten kann.

Die Bedingung  $|y|/(x - 1) < \sqrt{2^{-\text{prec}}}$  wird im Programm wie folgt realisiert: Zunächst wird mit  $\text{tm} = y/(x - 1)$  der Quotient  $y/(x - 1)$  intervallmäßig ausgewertet<sup>15</sup>, d.h. es gilt:  $y/(x - 1) \in \text{tm}$ . Verlangt man danach  $\text{expo}(\text{Sup}(\text{tm})) \leq -\text{prec} \oslash 2 - 1$ , so ist die Bedingung  $|y|/(x - 1) < \sqrt{2^{-\text{prec}}}$  erfüllt.  $\oslash$  bedeutet dabei die zur Null gerundete integer-Division durch 2, wenn die Current-Präzision  $\text{prec} \geq 2$  ungerade ist.

Für  $\text{prec} \geq 2$  und  $x > 1$  müssen wir also noch beweisen:

$$(C.130) \quad \text{expo}(\text{Sup}(\text{tm})) \leq -\text{prec} \oslash 2 - 1 \implies \frac{|y|}{x - 1} < \sqrt{2^{-\text{prec}}}.$$

Zum **Beweis** sei zunächst  $y \geq 0$ . Dann gilt  $y/(x - 1) \leq \text{Sup}(\text{tm})$  und die Behauptung rechts in (C.130) ist erfüllt, falls

$$\begin{aligned} \text{Sup}(\text{tm}) &< \sqrt{2^{-\text{prec}}} \\ \iff m \cdot 2^{ex} &< \sqrt{2^{-\text{prec}}}, \quad 0.5 \leq m < 1, \quad ex := \text{expo}(\text{Sup}(\text{tm})) \\ \iff m^2 \cdot 2^{2ex} &< 2^{-\text{prec}}, \end{aligned}$$

<sup>15</sup> $x, y$  sind Punktintervalle, welche die Maschinenzahlen  $x, y$  einschließen.

und wegen  $m^2 < 1$  ist die letzte Ungleichung erfüllt, falls  $2^{2ex} < 2^{-\text{prec}}$ , bzw. falls  $ex < -\text{prec}/2$ .  
Damit ist bewiesen:

$$(C.131) \quad \text{expo}(\text{Sup}(\text{tm})) \leq -\text{prec}/2 \implies \frac{|y|}{x-1} < \sqrt{2^{-\text{prec}}}.$$

Da auf der Maschine integer-Divisionen zur Null gerundet werden, gilt  $-\text{prec}/2 \geq -\text{prec} \oslash 2 - 1$ , womit dann (C.130) für  $y \geq 0$  bewiesen ist.

Für den Rest des Beweises sei jetzt  $y < 0$ . Dann gilt  $0 \geq \text{Sup}(\text{tm}) \geq y/(x-1)$  bzw.

$$(C.132) \quad -\text{Sup}(\text{tm}) \leq \frac{-y}{x-1} = \frac{|y|}{x-1},$$

und die Behauptung rechts in (C.130) ist erfüllt, falls

$$\begin{aligned} & -\text{Sup}(\text{tm}) < \sqrt{2^{-\text{prec}}} \\ \iff & -(m \cdot 2^{ex}) < \sqrt{2^{-\text{prec}}}, \quad -1 < m \leq -0.5, \quad ex := \text{expo}(\text{Sup}(\text{tm})) \\ \iff & m^2 \cdot 2^{2ex} < 2^{-\text{prec}}, \end{aligned}$$

und wegen  $m^2 < 1$  ist die letzte Ungleichung erfüllt, falls  $2^{2ex} < 2^{-\text{prec}}$ , bzw. falls  $ex < -\text{prec}/2$ .  
Der Rest des Beweises verläuft dann wie auf Seite 348 ■

Im Fall  $x > 1$  und  $|y|/(x-1) \geq \sqrt{2^{-\text{prec}}}$  wird  $\delta$  nach (C.114) ausgewertet:

$$(C.133) \quad 2\delta = \frac{\frac{y^2}{x^2-1} - \left[ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right] \cdot \left[ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right]}{1 + \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right) + \frac{1}{2} \cdot \left( \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right)}$$

Dabei entsteht sofort wieder die Frage, ob der Zähler in (C.133) als Differenz zweier positiver Größen ohne merkliche Auslöschung berechnet werden kann. Zur Beantwortung dieser Frage betrachten wir die für  $y \neq 0$  äquivalenten Gleichungen

$$\begin{aligned} & \frac{y^2}{x^2-1} - \left[ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right] \cdot \left[ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right] = r \cdot \frac{y^2}{x^2-1} \\ & 1 - r = \frac{1}{y^2} \cdot \left[ \sqrt{(x-1)^2 + y^2} - (x-1) \right] \cdot \left[ \sqrt{(x+1)^2 + y^2} - (x+1) \right] \end{aligned}$$

und suchen im ganzen Bereich  $\beta \geq 0.75$  im Falle  $x > 1$  für  $r = r(x, y)$  eine möglichst große Unterschranke bzw. für  $1 - r$  eine möglichst kleine Oberschranke. Lässt sich dies realisieren, so kann  $2\delta$  nach (C.133) ohne Auslöschung ausgewertet werden. Zur Berechnung einer möglichst kleinen Oberschranke von  $1 - r$  kann man sich wegen der Symmetrie zur reellen Achse, d.h. wegen  $r(x, y) \equiv r(x, -y)$  auf  $y > 0$  beschränken. Es gilt dann:

$$(C.134) \quad 1 - r = \left[ \sqrt{\left(\frac{x-1}{y}\right)^2 + 1} - \frac{x-1}{y} \right] \cdot \left[ \sqrt{\left(\frac{x+1}{y}\right)^2 + 1} - \frac{x+1}{y} \right]$$

Im nächsten Schritt zeigen wir, dass  $\partial(1-r)/\partial y$  für  $y > 0$  positiv ist, so dass das Maximum von  $1 - r$  für  $x \geq 1$  daher auf dem folgenden Hyperbelast liegt:

$$(C.135) \quad y = \frac{1}{12} \cdot \sqrt{112x^2 - 63}$$

Für  $y > 0$  gilt mit

$$A := \sqrt{(x-1)^2 + y^2} > x-1; \quad B := \sqrt{(x+1)^2 + y^2} > x+1;$$

$$\frac{\partial(1-r)}{\partial y} = \frac{[A - (x-1)] \cdot [B - (x+1)] \cdot [B(x-1) + A(x+1)]}{y^3 \cdot A \cdot B}$$

und daraus ergibt sich für  $x \geq 1$  die gewünschte Beziehung  $\partial(1-r)/\partial y > 0$ . Das gesuchte Maximum von  $1-r(x,y)$  liegt daher auf dem Hyperbelast. Um dieses Maximum zu berechnen, setzen wir daher  $y$  aus (C.135) in (C.134) ein. Nach einigen Umrechnungen erhält man dann das erstaunlicher Ergebnis:

$$(C.136) \quad 1-r \equiv \frac{1}{7} = 0.142857\dots, \quad \text{falls } x \geq 1,$$

d.h. auf dem Hyperbelast ist  $1-r(x,y)$  eine konstante Funktion. Man zeigt nun leicht, dass das Maximum von  $1-r(x,y)$  für  $x=1$  und  $0 < y \leq 7/12$ , d.h. also auf der Parallelen zur imaginären Achse durch  $x=1$  ebenfalls durch den Wert  $1/7$  gegeben ist. Der Nachweis bleibt dem Leser überlassen. Für  $1-r(x,y) \leq 1/7$  haben wir damit im Bereich  $\beta \geq 0.75$  eine hinreichend kleine Oberschranke berechnet, so dass  $2\delta$  nach (C.133) trotz der Differenz im Zähler auf der Maschine stabil, d.h. ohne Auslöschung ausgewertet werden kann.

Wir betrachten jetzt den Fall  $x = +1$ . Mit (C.113) von Seite 343 erhält man:

$$(C.137) \quad \delta = \left\{ 1 - \sqrt{1 + \left(\frac{y}{2}\right)^2} \right\} + \frac{|y|}{2}, \quad \text{falls } x = 1;$$

Wir zeigen zunächst, dass (C.137) im Bereich  $\beta \geq 0.75$  nur für  $|y| \leq 7/12$  auszuwerten ist. Im Falle  $x=1$  gilt nach (C.113)

$$\begin{aligned} \beta &= \frac{2}{\sqrt{4+y^2} + |y|} \geq \frac{3}{4} \\ \iff \sqrt{4+y^2} + |y| &\leq \frac{8}{3} \iff \sqrt{4+y^2} = \frac{8}{3} - |y| \\ \iff |y| &\leq \frac{7}{12} = 0.58333\dots \blacksquare \end{aligned}$$

In (C.137) wird der erste Summand  $\{\dots\} \leq 0$  mit Hilfe der Funktion `sqrtp1m1()` ausgewertet, die Auslöschung innerhalb von  $\{\dots\}$  vermeidet. Man könnte jedoch einwenden, dass in (C.137) die Addition eines negativen und positiven Summanden zur Auslöschung führen kann. Um dies auszuschließen, zeigen wir mit  $t := |y|/2$ , dass in

$$(C.138) \quad t - \left\{ \sqrt{1+t^2} - 1 \right\} = r \cdot t$$

der Faktor  $r < 1$  für  $|y| \leq 7/12$  bzw. für  $t \leq 7/24$  hinreichend groß ist. Mit

$$\begin{aligned} r(t) &:= 1 - \frac{1}{t} \cdot \left\{ \sqrt{1+t^2} - 1 \right\} \quad \text{folgt} \\ r'(t) &= -\frac{\sqrt{1+t^2} - 1}{t^2 \cdot \sqrt{1+t^2}} < 0 \end{aligned}$$

so dass  $r(t)$  streng monoton fallend ist. Im Bereich  $0 \leq t \leq 7/24$  gilt daher

$$(C.139) \quad r(t) \geq r(7/24) = \frac{6}{7} = 0.85714\dots$$

und mit dieser Unterschranke für  $r(t)$  folgt aus (C.138), dass für  $|y| \leq 7/12$  bei der Auswertung von (C.137) keine wesentliche Auslöschung auftreten kann. Damit könnte man den Fall  $x=1$  als erledigt ansehen. Für  $|y| \ll 1$  lässt sich die Auswertung nach (C.137) jedoch noch wesentlich

vereinfachen. Mit  $t := |y|/2$  folgt nach (C.137) zunächst  $\delta = 1 + t - \sqrt{1 + t^2}$ , und damit ergibt sich ganz elementar die folgende Einschließung für  $\delta$ :

$$(C.140) \quad t \cdot \left(1 - \frac{t}{2}\right) \leq \delta \leq t, \quad t := \frac{|y|}{2}.$$

Mit  $0 \leq \delta < 1$  gilt nach (C.122) für  $\arccos(1 - \delta)$  die Einschließung

$$(C.141) \quad \sqrt{2\delta} \leq \arccos(1 - \delta) \leq \sqrt{2\delta} \cdot \frac{1}{1 - \delta}, \quad 0 \leq \delta < 1,$$

und zusammen mit (C.140) folgt direkt:

$$(C.142) \quad \sqrt{2t} \cdot \sqrt{1 - \frac{t}{2}} \leq \arccos(1 - \delta) \leq \frac{\sqrt{2t}}{1 - t}, \quad \text{falls } 0 \leq t = \frac{|y|}{2} < 1.$$

Wir verlangen zusätzlich:

$$\begin{aligned} & \sqrt{1 - \frac{t}{2}} \geq \text{pred}(1) = 1 - 2^{-\text{prec}} \\ \Leftrightarrow & 1 - \frac{t}{2} \geq 1 - 2 \cdot 2^{-\text{prec}} + 2^{-2\text{prec}} \\ \Leftrightarrow & |y| \leq 2^{-\text{prec}+3} - 2^{-2\text{prec}+2} = 2^{-\text{prec}}(8 - 2^{-\text{prec}+2}), \end{aligned}$$

Wegen  $|y| = m \cdot 2^{ex} < 2^{ex}$  und wegen  $8 - 2^{-\text{prec}+2} > 2^2$  ist die letzte Ungleichung erfüllt, wenn gilt:  $\text{expo}(y) = ex \leq -\text{prec} + 2$ . Zusammen mit (C.142) erhält man:

$$(C.143) \quad \text{pred}(1) \cdot \sqrt{|y|} \leq \arccos(1 - \delta) \leq \frac{\sqrt{|y|}}{1 - t}, \quad \text{falls } \text{expo}(y) \leq -\text{prec} + 2.$$

Wie auf Seite 346 verlangen wir außerdem

$$\frac{1}{1 - t} < \text{succ}(1) = 1 + 2^{1-\text{prec}} \quad \Leftrightarrow \quad t < \frac{2^{1-\text{prec}}}{1 + 2^{1-\text{prec}}},$$

und die letzte Ungleichung ist erfüllt, wenn  $t < 2^{-\text{prec}}$ , bzw. wenn  $|y| < 2^{-\text{prec}+1}$  oder wenn gilt  $\text{expo}(y) \leq -\text{prec} + 1$ . Mit (C.143) ergibt sich daraus die Einschließung:

$$(C.144) \quad \text{pred}(1) \cdot \sqrt{|y|} \leq \arccos(1 - \delta) \leq \sqrt{|y|} \cdot \text{succ}(1), \quad \text{falls } \text{expo}(y) \leq -\text{prec} + 1.$$

Um die Einschließung in (C.144) auf der Maschine effizient realisieren zu können, berechnet man zu gegebener Maschinenzahl  $y$  mit  $\mathbf{u} = [\mathbf{u1}, \mathbf{u2}]$  zunächst eine Einschließung für  $\sqrt{|y|}$ :

$$\mathbf{u} = \text{sqrt}(\text{MpfiClass}(\text{abs}(y))) \quad \rightsquigarrow \quad \sqrt{|y|} \in [\mathbf{u1}, \mathbf{u2}].$$

Mit (C.144) erhält man daher

$$(C.145) \quad \text{pred}(1) \cdot \mathbf{u1} \leq \arccos(1 - \delta) \leq \mathbf{u2} \cdot \text{succ}(1), \quad \text{falls } \text{expo}(y) \leq -\text{prec} + 1;$$

Bei Anwendung von (C.106) und (C.129) kann man außerdem noch die beiden Multiplikationen in (C.145) vermeiden:

$$(C.146) \quad \text{pred}(\mathbf{u1}) \leq \arccos(1 - \delta) \leq \text{succ}(\text{succ}(\mathbf{u2})), \quad \text{falls } \text{expo}(y) \leq -\text{prec} + 1.$$

(C.146) liefert damit unter den Voraussetzungen  $\beta \geq 0.75$ ,  $x = 1$  und  $|y| < 2^{-\text{prec}+1}$ , bzw.  $\text{expo}(y) \leq -\text{prec} + 1$  einen sehr effektiven Algorithmus zur fast optimalen Einschließung von

$\arccos(1 - \delta)$ . Für  $\text{expo}(y) > -\text{prec} + 1$  benutzt man (C.137) zusammen mit der Beziehung  $\arccos(1 - \delta) \equiv \arcsin(\sqrt{\delta \cdot (2 - \delta)})$ .

Wir müssen jetzt noch im Bereich  $\beta \geq 0.75$  den letzten Fall  $0.75 \leq x < +1$  betrachten: Ausgehend von (C.113) auf Seite 343 findet man nach einigen Umrechnungen:

$$\begin{aligned}
 2\delta &= -(x+1) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right\} + \\
 &\quad + (1-x) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} - 1 \right\} + 2 \cdot (1-x) \\
 \text{(C.147)} \quad 2 \cdot \delta(x, y) &= -(x+1) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x+1}\right)^2} - 1 \right\} + \\
 &\quad + (1-x) \cdot \left\{ \sqrt{1 + \left(\frac{y}{x-1}\right)^2} + 1 \right\}, \quad \text{falls } x < 1;
 \end{aligned}$$

Da in (C.147) zwei Summanden mit unterschiedlichen Vorzeichen zu addieren sind, stellt sich auch jetzt wieder die Frage, ob diese Addition in  $0.75 \leq x < +1$  ohne wesentliche Auslöschung durchführbar ist. Zur Beantwortung dieser Frage betrachten wir die beiden äquivalenten Gleichungen

$$\begin{aligned}
 (1-x) \cdot \left[ \sqrt{1 + \left(\frac{y}{1-x}\right)^2} + 1 \right] - (1+x) \cdot \left[ \sqrt{1 + \left(\frac{y}{1+x}\right)^2} - 1 \right] \\
 = r \cdot (1-x) \cdot \left[ \sqrt{1 + \left(\frac{y}{1-x}\right)^2} + 1 \right] \\
 \text{(C.148)} \quad 1-r = \frac{(1+x) \cdot \left[ \sqrt{1 + \left(\frac{y}{1+x}\right)^2} - 1 \right]}{(1-x) \cdot \left[ \sqrt{1 + \left(\frac{y}{1-x}\right)^2} + 1 \right]}
 \end{aligned}$$

und suchen im Bereich  $\beta \geq 0.75$  im Falle  $0.75 \leq x < 1$  für  $r = r(x, y)$  eine möglichst große Unterschranke bzw. für  $1 - r$  eine möglichst kleine Oberschranke. Lässt sich dies realisieren, so kann  $2\delta$  nach (C.147) ohne Auslöschung ausgewertet werden. Zur Berechnung einer möglichst kleinen Oberschranke von  $1 - r$  kann man sich wegen der Symmetrie zur reellen Achse, d.h. wegen  $r(x, y) \equiv r(x, -y)$  auf  $y > 0$  beschränken. Wir zeigen jetzt wieder, dass  $\partial(1 - r)/\partial y$  für  $y > 0$  positiv ist, so dass das Maximum von  $1 - r$  für  $x < 1$  auf dem folgenden Hyperbelast liegt:

$$\text{(C.149)} \quad y = \frac{1}{12} \cdot \sqrt{112x^2 - 63}$$

Für  $0.75 \leq x < 1$  ist die partielle Ableitung  $\partial(1 - r)/\partial y$  mit

$$A := \sqrt{(1-x)^2 + y^2} > 1-x, \quad B := \sqrt{(1+x)^2 + y^2} > 1+x \quad \text{gegeben durch:}$$

$$\text{(C.150)} \quad \frac{\partial(1-r)}{\partial y} = \frac{y \cdot \{-4x + A \cdot (1-x) + B \cdot (1+x)\}}{A \cdot B \cdot (A+1-x)^2}$$

und für  $y > 0$  ist  $\partial(1 - r)/\partial y > 0$ , wenn in (C.150) gilt:  $\{\dots\} > 0$ , d.h. wenn

$$\text{(C.151)} \quad C := (1-x) \cdot \sqrt{(1-x)^2 + y^2} + (1+x) \cdot \sqrt{(1+x)^2 + y^2} > 4x.$$

Wegen  $C \geq (1-x)^2 + (1+x)^2 = 2 + 2x^2$  ist (C.151) erfüllt, wenn gilt

$$2 + 2x^2 > 4x \iff (1-x)^2 > 0.$$

Das Maximum von  $1 - r(x, y)$  liegt damit für  $0.75 \leq x < 1$  auf dem Hyperbelast

$$(C.152) \quad y = \frac{1}{12} \cdot \sqrt{112x^2 - 63}.$$

Um das Maximum zu berechnen, setzen wir  $y$  aus (C.152) ein in (C.148) und erhalten

$$\begin{aligned} 1 - r(x, y(x)) &= \frac{\sqrt{\frac{16}{9}x^2 + 2x + \frac{9}{16}} - x - 1}{\sqrt{\frac{16}{9}x^2 - 2x + \frac{9}{16}} - x + 1} \\ &= \frac{\sqrt{\frac{1}{144}(16x+9)^2} - (x+1)}{\sqrt{\frac{1}{144}(16x-9)^2} - (x-1)} = \frac{4x-3}{4x+3}. \end{aligned}$$

Da der letzte Term in  $x$  monoton wächst, gilt für  $0.75 \leq x < 1$

$$1 - r(x, y(x)) \leq \frac{4 \cdot 1 - 3}{4 \cdot 1 + 3} = \frac{1}{7}$$

Im Fall  $0.75 \leq x < 1$  haben wir damit im Bereich  $\beta(x, y) \geq 0.75$  für  $1 - r(x, y)$  mit  $1/7$  eine hinreichend kleine Oberschranke gefunden, so dass  $2\delta$  nach (C.147) ohne Auslöschung ausgewertet werden kann. Bitte beachten Sie, dass wir in den verschiedenen Fällen  $x < 1$ ,  $x = 1$ ,  $x > 1$  für den Ausdruck  $1 - r$  jeweils die gleiche, hinreichend kleine Oberschranke  $1/7$  gefunden haben, obwohl in (C.147), (C.137), (C.114) verschieden Terme zur Berechnung von  $\delta$  definiert wurden. Vermutlich ist die Gleichheit der gefundenen Oberschranken ein Hinweis dafür, dass die gewählten Terme zur Auswertung von  $\delta(x, y)$  für die numerische Stabilität schon optimal gewählt wurden.

Zur Einschließung von  $\Re(\arccos(\mathbf{Z}))$  ist nach Seite 340 die reelle arccos-Funktion für den Intervall-Ausdruck  $\beta$  auszuwerten:

$$(C.153) \quad \arccos(\beta), \quad \beta := \frac{2\mathbf{x}}{\sqrt{(\mathbf{x}+1)^2 + \mathbf{y}^2} + \sqrt{(\mathbf{x}-1)^2 + \mathbf{y}^2}},$$

wobei die Punktintervalle  $\mathbf{x}, \mathbf{y}$  durch die Extrempunkte auf dem Rand von  $\mathbf{Z}$  bestimmt sind. Dabei haben wir zunächst den Fall  $\text{Inf}(\beta) \rightarrow +1$  betrachtet, um Überschätzungen in der Nähe der Nullstelle der reellen arccos-Funktion zu vermeiden. Überschätzungen treten jedoch auch im Fall  $\text{Sup}(\beta) \rightarrow -1$  auf, da die reelle arccos-Funktion am linken Definitionsrand nahezu senkrechte Tangenten besitzt, so dass kleine Überschätzungen bei der Berechnung des Arguments  $\beta$  deutliche Überschätzungen bei den Funktionswerten verursachen. Mit dem Punktargument  $z = -1 + i \cdot 2^{-200}$  und  $\text{prec} = 53$  erhält man z.B. für  $\Re(w) := \Re(\arccos(z))$  die folgende recht grobe Einschließung:

$$\Re(w) \in [3.1415926\mathbf{237874627}, 3.1415926\mathbf{535898020}]$$

mit nur 8 korrekten Dezimalziffern. Diese zu groben Einschließungen lassen sich durch Anwendung folgender, für  $0 \leq \delta \leq 2$  geltenden Identität vermeiden:

$$(C.154) \quad \arccos(-1 + \delta) \equiv \pi - \arccos(1 - \delta) \equiv \pi - \arcsin(\sqrt{\delta \cdot (2 - \delta)}).$$

Bei vorgegebenem  $\beta = -1 + \delta \leq -0.75$  gilt dann  $\delta = 1 - (-\beta)$ , mit  $(-\beta) \geq 0$ , so dass  $\delta \geq 0$  mit Hilfe der Gleichungen (C.114), (C.137) und (C.147) intervallmäßig berechnet werden kann. Die intervallmäßige Auswertung der rechten Seite von (C.154) liefert dann eine fast optimale



Einschließung des Funktionswertes  $\arccos(\beta) = \arccos(-1 + \delta)$ . Der beschriebene Algorithmus liefert dann mit dem Punktargument  $z = -1 + i \cdot 2^{-200}$  für  $\Re(w) := \Re(\arccos(z))$  und `prec = 53` die folgende fast optimal Einschließung:

$$\Re(w) \in [3.141592653589792, 3.141592653589794].$$

Damit wird der Realteil  $\Re(\arccos(z)) = \arccos(\beta)$  im ganzen Bereich  $|\beta| \leq +1$  in ausreichender Genauigkeit eingeschlossen.

### C.2.22. $\arctan(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - \{i \cdot (-\infty, -1) \cup i \cdot (+1, +\infty)\}$  liefert die Funktion

```
MpfciClass atan(const MpfciClass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\text{atan}(Z)$  für die komplexen Funktionswerte  $\arctan(z)$ , mit  $z \in Z$ .

$$\{\arctan(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \text{atan}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der imaginären Achse von  $-i\infty$  bis  $-i$  bzw. von  $+i$  bis  $+i\infty$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben.

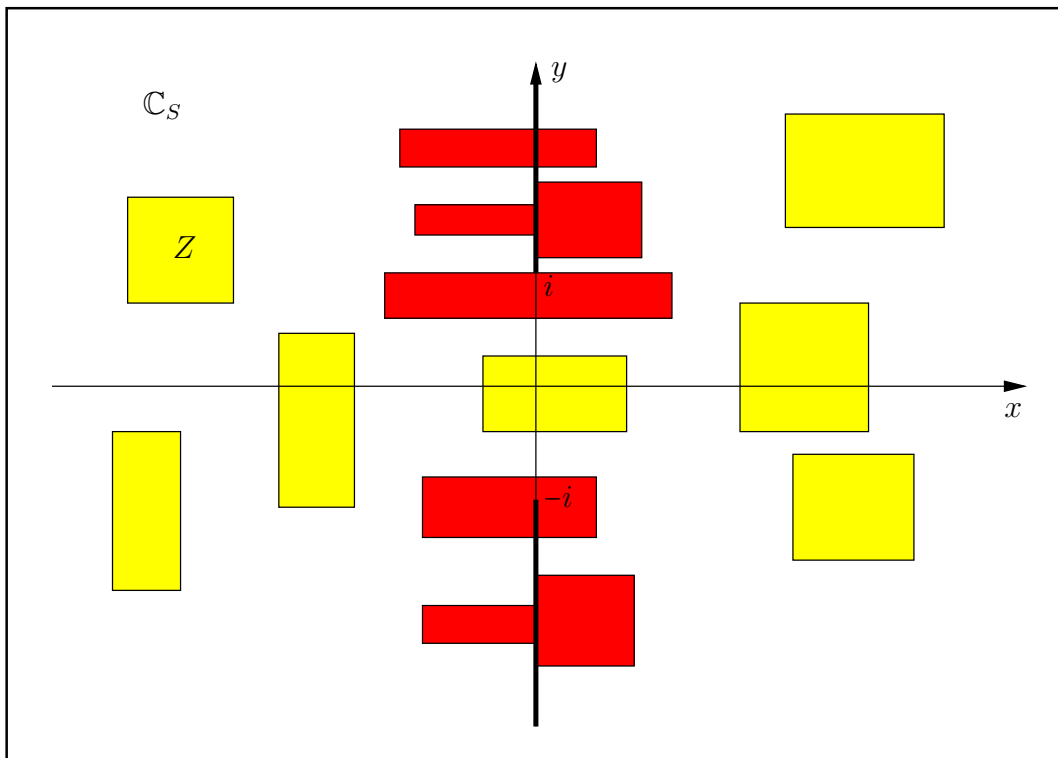


Abbildung C.41.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\text{atan}(Z)$ .

Nach Abbildung C.41 darf also kein Punkt der beiden Verzweigungsschnitte ein Element eines erlaubten, achsenparallelen Intervalls  $Z$  sein.

Mit  $Z = [2^{-1073741824}, 2^{-1073741824}] + i \cdot [1, 1]$  erhält man mit der Current-Präzision  $\text{prec} = 70$  die Einschließung<sup>16</sup>

$$\begin{aligned} \text{atan}(Z) = & ([7.85398163397448309613e - 1, 7.85398163397448309619e - 1], \\ & [3.72130559324020099216e8, 3.72130559324020099218e8]). \end{aligned}$$

Mit  $Z = [0.5, 1.5] + i \cdot [-1, 1]$  erhält man mit der Current-Präzision  $\text{prec} = 80$  die Einschließung

$$\begin{aligned} \text{atan}(Z) = & ([4.63647609000806116214256e - 1, 1.10714871779409050301707], \\ & [-7.08303336014054020062385e - 1, 7.08303336014054020062385e - 1]). \end{aligned}$$

<sup>16</sup>Die Maschinenzahl  $2^{-1073741824} = \text{minfloat}()$  ist präzisionsunabhängig die kleinste positive Maschinenzahl.

### C.2.22.1. Algorithmus

In [54] wird für den Imaginärteil der arctan-Funktion u.a. folgender Term ausgewertet:

$$T(x) := \ln \left[ 1 + \frac{4 \cdot \sqrt{1+x^2}}{x^2 + (1 - \sqrt{1+x^2})^2} \right], \quad x = [x] : \text{Punktintervall}$$

dabei wurde vermutlich übersehen, dass sich  $T(x)$  noch wesentlich vereinfachen lässt, wodurch sowohl die Güte der Einschließung als auch die Laufzeit deutlich verbessert werden können. Nach entsprechender Vereinfachung sind dann für den Imaginärteil folgende Terme auszuwerten:

$$(C.155) \quad T(x) := \ln \left[ 1 + \frac{2}{\sqrt{1+x^2} - 1} \right], \quad x = [x] : \text{Punktintervall}$$

$$(C.156) \quad Q_{1,2}(x, y) := \ln \left[ 1 \pm \frac{4y}{x^2 + (1 \mp y)^2} \right], \quad \text{nur } y = [y] \text{ ist ein Punktintervall.}$$

In (C.156) bedeuten die Indices 1 bzw. 2 das jeweilige obere bzw. untere Vorzeichen.

### Auswertung von $T(x)$

Der Term  $T(x)$  ist nur für Punktintervalle  $x = [x]$  auszuwerten. Zur Vermeidung von Overflow müssen die beiden Fälle  $x \rightarrow 0$ , d.h. Nenner  $\rightarrow 0$  und  $x \rightarrow \infty$ , d.h.  $x^2 \rightarrow \infty$  gesondert betrachtet werden.

Für  $x \rightarrow 0$  benutzen wir die folgende Darstellung:

$$(C.157) \quad T(x) = \ln \left[ 2 + x^2 + 2 \cdot \sqrt{1+x^2} \right] - 2 \cdot \ln(x)$$

Der obige Term  $T(x)$  kann jetzt intervallmäßig für  $x \rightarrow 0$  problemlos ausgewertet werden. Die Laufzeit kann jedoch durch die folgende Vereinfachung noch mehr als halbiert werden:

Zunächst gilt:

$$\begin{aligned} \alpha &:= 2 + x^2 + 2 \cdot \sqrt{1+x^2} \\ &= 4 + \left[ 2x^2 + \sum_{k=2}^{\infty} (-1)^{k+1} \frac{1 \cdot 3 \cdot 5 \cdot \dots \cdot (2k-3)}{4 \cdot 6 \cdot 8 \cdot \dots \cdot 2k} \cdot (x^2)^k \right], \end{aligned}$$

und weil die Reihe in [...] für  $x^2 < 1$  eine alternierende Leibniz-Reihe ist, erhält man für  $\alpha$  die Einschließung:

$$4 < \alpha < 4 + 2x^2, \quad x^2 < 1;$$

Wegen der Monotonie der ln-Funktion folgt daraus direkt die Einschließung:

$$\ln(4) < \ln(\alpha) < \ln(4 + 2x^2) = \ln(4) + \ln\left(1 + \frac{1}{2} \cdot x^2\right) < \ln(4) + \frac{1}{2} \cdot x^2, \quad \text{d.h.}$$

$$(C.158) \quad \ln(4) < \ln(\alpha) < \ln(4) + \frac{1}{2} \cdot x^2, \quad x^2 < 1.$$

Es gilt also für das folgende Nicht-Maschinenintervall  $\text{Nm}(x)$

$$(C.159) \quad \ln(\alpha) \in \text{Nm}(x) := \left[ \ln(4), \ln(4) + \frac{1}{2} \cdot x^2 \right], \quad x^2 < 1.$$

Das Ziel ist nun,  $x \in ]0, 1[$  in (C.159) so klein zu wählen, dass für jede Current-Präzision  $\text{prec}$ , mit  $2 \leq \text{prec} \leq \text{prec}_0$ , ein Maschinenintervall  $U$  so angegeben werden kann, dass gilt:

$$(C.160) \quad \text{Nm}(x) \subseteq U \quad \text{und} \quad \sup(U) = \text{succ}(\text{Inf}(U)).$$

Es ist klar, dass die zweite Bedingung in (C.160) bei festem  $\text{prec}$  nicht erfüllt sein kann, wenn  $x = 2^k \in ]0, 1[$  zu groß gewählt wird. In einer Schleife, d.h. mit der Funktion

```

int Schleife1(const int prec)
{
    int k(-prec/2 -9);
    MpfClass::SetCurrPrecision(prec);
    MpfClass U(0);
    MpfClass ln4(MpfClass::Ln2(2*prec));
    MpfClass x, u;
    times2pown(ln4,1); // ln(4) mit Praezision 2*prec;
    MPFR::MpfrClass::MpfrClass LrI, LrS;

    do
    {
        k++;
        MpfClass::SetCurrPrecision(2*prec);
        x = MpfClass( comp(MPFR::MpfrClass::MpfrClass(0.5),k) );
        LrI = Inf(ln4);
        LrS = Sup(ln4 + sqr(x)/2);
        u = MpfClass(LrI, LrS);
        MpfClass::SetCurrPrecision(prec);
        U = u;
        U.RoundPrecision(prec);
    } while( Sup(U) == succ(Inf(U)) );

    return k-1;
}

```

berechnet man daher, beginnend mit einem hinreichend kleinen  $k < 0$ , in stark vergrößerter Präzision<sup>17</sup> eine Einschließung von  $\ln(4) + x^2/2 \in u$  und rundet dieses  $u$  in ein Maschinenintervall  $U \supseteq u$ . In der Schleife wird dann das maximale  $k$  berechnet, für das die obige Bedingung  $\text{sup}(U) = \text{succ}(\text{Inf}(U))$  erfüllt ist, und dieses  $k$  wird von `Schleife1(prec)` zurückgegeben. Für das in `Schleife1(prec)` in der Current-Präzision `prec` berechnete  $U$  und für das zurückgegeben  $k$  gilt damit:

$$0 \leq x \leq 0.5 \cdot 2^k \implies \ln(4) + \frac{x^2}{2} \in U, \ln(4) \in U \text{ und } \text{sup}(U) = \text{succ}(\text{Inf}(U)).$$

Berechnet man daher in der Current-Präzision `prec` mit  $\hat{U} := 2 \diamond \text{MpfClass::Ln2}(\text{prec})$  eine nicht notwendig optimale Maschineneinschließung von  $\ln(4)$ , so gilt  $U \subseteq \hat{U}$ , wobei  $U = \hat{U}$  nur dann gilt, wenn  $\hat{U}$  eine optimale Einschließung von  $\ln(4)$  ist. Bedeutet dann in der Current-Präzision `prec`  $[x]$  das Punktintervall, das  $x$  einschließt, so gilt

$$0 \leq x \leq 0.5 \cdot 2^k \implies T(x) \in -2 \diamond \ln([x]) \diamond U \subseteq -2 \diamond \ln([x]) \diamond \hat{U},$$

wobei man jetzt das obige Intervall  $-2 \diamond \ln([x]) \diamond \hat{U}$  mit minimalem Durchmesser und optimaler Laufzeit berechnen kann.

Da der maximale  $k$ -Wert von der gewählten Current-Präzision `prec`  $\geq 2$  unmittelbar abhängt, wäre es äußerst sinnvoll, wenn man eine von `prec` abhängige Unterschranke  $k_0(\text{prec}) < k$  so angeben könnte, dass für einen möglichst großen Bereich  $2 \leq \text{prec} \leq \text{prec}_0$  folgendes gilt:

$$0 \leq x \leq 0.5 \cdot 2^{k_0} \implies T(x) \in -2 \diamond \ln([x]) \diamond \hat{U}, \quad 2 \leq \text{prec} \leq \text{prec}_0.$$

Mit einer zweiten Schleife (Programmteil `Schleife2`) findet man dazu<sup>18</sup>

$$(C.161) \quad 0 \leq x \leq 0.5 \cdot 2^{-(\text{prec} \diamond 2 + 9)}, \quad 2 \leq \text{prec} \leq \text{prec}_0 = 677370 \implies T(x) \in -2 \diamond \ln([x]) \diamond \hat{U}.$$

<sup>17</sup>Es zeigt sich, dass die berechneten Einschließungen in doppelter Präzision eng genug sind.

<sup>18</sup> $\diamond$  bedeutet die integer-Division auf der Maschine, wobei zur Null hin gerundet wird, wenn der exakte Quotient keine integer-Zahl ist.

```

// Schleife2
int k, prec = 1;
do
{
    prec++;
    cout << "prec == " << prec << endl;
    k = Schleife1(prec);
} while( k >= -(prec/2+9) );

```

Damit können wir in nahezu allen praktischen Fällen mit (C.161) eine optimale Einschließung berechnen. Abschließend geben wir für  $x = m \cdot 2^{ex}$  noch an, wie die erste Bedingung zu realisieren ist. Wegen  $m \cdot 2^{ex} < 2^{ex}$  verlangen wir  $2^{ex} \leq 0.5 \cdot 2^{-(\text{prec} \oslash 2 + 9)} \iff ex \leq -\text{prec} \oslash 2 - 10$ , wobei  $ex$  durch  $\text{expo}(x)$  definiert ist.

Im sehr seltenen Fall  $\text{prec} > \text{prec}_0$  wird die Einschließung von  $T$  mit (C.157) jedoch etwas aufwendiger realisiert. Der Präzision von  $\text{prec}_0 = 677370$  Bits entsprechen ca. 203908 Dezimalstellen, womit alle praktischen Fälle abgedeckt sein sollten.

Wir kommen jetzt zur Auswertung von  $T(x)$  für  $x \rightarrow +\infty$ . Bei der normalen Berechnung von  $T(x)$  nach (C.155) wird der Nenner  $\sqrt{1+x^2} - 1$  mit Hilfe der Funktion `sqrt1pm1(sqr([x]))` ausgewertet. Um dabei einen Überlauf zu vermeiden, verlangen wir

$$x^2 = m^2 \cdot 2^{2ex} < 2^{2ex} < 2^{+1073741820} \iff \text{expo}(x) = ex < 536870910,$$

so dass im Fall  $\text{expo}(x) \geq 536870910$  ein Überlauf eintreten kann. Um in diesem Fall einen solchen Überlauf zu vermeiden, schreiben wir den Bruch in (C.155) wie folgt um:

$$\beta := \frac{2}{\sqrt{1+x^2} - 1} = \frac{\frac{2}{x}}{\sqrt{1 + \frac{1}{x} \cdot \frac{1}{x} - \frac{1}{x}}}$$

Auch jetzt könnte man den Term rechts ohne Overflow für  $x \rightarrow +\infty$  auf dem Rechner auswerten. Die Laufzeit kann jedoch durch die folgenden Überlegungen wesentlich reduziert werden. Mit  $r := 1/x$  erhält man zunächst:

$$\begin{aligned} \beta(r) &:= 2r \cdot \frac{1}{\sqrt{1+r^2} - r} = 2r \cdot \frac{\sqrt{1+r^2} + r}{(\sqrt{1+r^2} - r) \cdot (\sqrt{1+r^2} + r)} \\ &= 2r \cdot (\sqrt{1+r^2} + r) \\ &= 2r + 2r^2 + \left( r^3 - \frac{r^5}{4} + \frac{r^7}{8} - \frac{5}{64}r^9 + \dots \right) \end{aligned}$$

und da der letzte Klammerausdruck ( ) wegen  $\sqrt{1+r^2}$  selbst wieder eine alternierende Leibniz-Reihe ist, ergibt sich für  $\beta(r)$  die Einschließung:

$$(C.162) \quad 2r + 2r^2 < \beta(r) < 2r + 2r^2 + r^3, \quad 0 < r < 1.$$

Da die Taylorreihe von

$$\ln(1+x) = x - \frac{1}{2}x^2 + \frac{1}{3}x^3 - + \dots, \quad -1 < x \leq 1$$

für  $0 < x \leq 1$  eine alternierende Leibniz-Reihe ist, folgt die Ungleichung

$$(C.163) \quad x - \frac{1}{2}x^2 < \ln(1+x) < x, \quad 0 < x \leq 1;$$

Mit den rechten Ungleichungen aus (C.162) und (C.163) folgt nun direkt:

$$\ln(1 + \beta(r)) < 2r + 2r^2 + r^3, \quad \text{falls } 2r + 2r^2 + r^3 < 1,$$

wobei die Bedingung  $2r + 2r^2 + r^3 < 1$  wegen  $x \rightarrow \infty$ , d.h.  $r \rightarrow 0$  sicher erfüllt sein wird. Wir benötigen jetzt noch eine Unterschranke von  $\ln(1 + \beta(r))$ . Dazu gilt wieder nach (C.162) und (C.163)

$$\ln(1 + \beta(r)) > \beta(r) - \frac{1}{2}\beta(r)^2 > 2r + 2r^2 - \frac{1}{2}\beta(r)^2$$

Aus (C.162) folgt für  $0 < r < 1$

$$-\frac{1}{2}\beta(r)^2 > -2r^2 - 4r^3 - 4r^4 - 2r^5 - \frac{1}{2}r^6 > -2r^2 - 4r^3 - 4r^4 - 2r^4 - 1r^4$$

und das ergibt dann die gesuchte Einschließung:

$$\begin{aligned} 2r - 4r^3 - 7r^4 < \ln(1 + \beta(r)) < 2r + 2r^2 + r^3, \quad 2r + 2r^2 + r^3 < 1; \\ 2r \cdot (1 - 2r^2 - \frac{7}{2}r^3) < \ln(1 + \beta(r)) < 2r \cdot (1 + r + \frac{1}{2}r^2), \end{aligned}$$

die noch etwas vereinfacht werden kann:

$$(C.164) \quad 2r \cdot (1 - 2r^2 - 4r^3) < \ln(1 + \beta(r)) < 2r \cdot (1 + r + \frac{1}{2}r^2), \quad 2r + 2r^2 + r^3 < 1.$$

Für hinreichend kleine  $r$  kann diese Einschließung noch wesentlich vereinfacht werden. Dazu verlangen wir für die Unterschranke

$$(C.165) \quad \begin{aligned} 1 - 2r^2 - 4r^3 > \text{pred}(1) = 1 - 2^{-\text{prec}} \\ \iff 2r^2 \cdot (1 + 2r) < 2^{-\text{prec}}. \end{aligned}$$

Mit  $r = 1/x = (1/m) \cdot 2^{-ex}$ ,  $ex \gg 1$ , folgt weiter

$$1 + 2r = 1 + \frac{1}{m} \cdot 2^{-ex+1} \leq 1 + 2^{-ex+2}, \quad \text{d.h. (C.165) ist erfüllt, wenn gilt: } 2r^2 \cdot 2 < 2^{-\text{prec}}.$$

Es gilt zusätzlich

$$2r^2 \cdot 2 = \frac{1}{m^2} \cdot 2^{-2ex+2} \leq 4 \cdot 2^{-2ex+2} = 2^{-2ex+4},$$

d.h. (C.165) ist erfüllt, wenn gilt

$$2^{-2ex+4} < 2^{-\text{prec}} \iff ex > \text{prec}/2 + 2.$$

Symbolisiert  $\oslash$  den Operator der integer-Division, der zur Null rundet, wenn der Quotient keine integer-Zahl ist, so gilt  $\text{prec}/2 < \text{prec} \oslash 2 + 1$ , d.h. (C.165) ist erfüllt, wenn gilt

$$(C.166) \quad ex > \text{prec} \oslash 2 + 3 \implies 1 - 2r^2 - 4r^3 > \text{pred}(1).$$

Zur Verbesserung der Oberschranke in (C.164) verlangen wir jetzt

$$(C.167) \quad \begin{aligned} 1 + r + \frac{r^2}{2} < \text{succ}(1) = 1 + 2^{-\text{prec}+1} \\ \iff r + \frac{r^2}{2} < 2^{-\text{prec}+1}. \end{aligned}$$

Es gilt zusätzlich noch die folgende Abschätzung

$$r + \frac{r^2}{2} = \frac{1}{m} \cdot 2^{-ex} + \frac{1}{m^2} \cdot 2^{-2ex-1} \leq 2 \cdot 2^{-ex} + 4 \cdot 2^{-2ex-1} = 2^{-ex+1}(1 + 2^{-ex}) < 2^{-ex+2},$$

d.h. (C.167) ist erfüllt, wenn gilt:  $2^{-ex+2} < 2^{-\text{prec}+1} \iff ex > \text{prec} + 1$ . Damit erhalten wir

$$(C.168) \quad ex > \text{prec} + 1 \implies 1 + r + \frac{r^2}{2} < \text{succ}(1).$$

Zusammen mit (C.164) folgt jetzt unter den Bedingungen  $ex > \text{prec} + 1$  und  $ex > \text{prec} \oslash 2 + 3$

$$2r \cdot \text{pred}(1) < \ln(1 + \beta(r)) < 2r \cdot \text{succ}(1), \quad 2r + 2r^2 + r^3 < 1.$$

Man kann noch einfach nachweisen, dass die beiden obigen Bedingungen erfüllt sind, wenn gilt  $ex > \text{prec} + 2$ . Zusammen mit (C.106) und (C.129) erhalten wir mit  $x = m \cdot 2^{ex}$ ,  $r = 1/x$ ,  $t = [2r]$ :

$$(C.169) \quad ex > \text{prec} + 2 \implies \text{pred}(\text{Inf}(t)) < \ln(1 + \beta(r)) < \text{succ}(\text{succ}(\text{Sup}(t))).$$

Im Fall  $ex \leq \text{prec} + 2$ , der wegen  $ex \geq 536870910$  in der Praxis kaum auftreten wird, benutzen wir:

$$(C.170) \quad 2r \cdot [1 - 2r^2(1 + 2r)] < \ln(1 + \beta(r)) < 2r \cdot [1 + r(1 + r)].$$

Die Auswertung von  $T(x)$  erfolgt mithilfe der Funktion

```
MpfiClass Aux_1_atan(const MPFR::MpfrClass::MpfrClass& x);
```

### Auswertung von $Q_{1,2}(x, y)$

Wir beschränken uns auf den Index 1, und nach (C.156) ist auszuwerten

$$(C.171) \quad Q_1(x, y) := \ln \left[ 1 + \frac{4y}{x^2 + (1-y)^2} \right],$$

dabei ist  $Q_1(x, y)$  einzuschließen, wobei  $y = [y]$  als Punktintervall und  $x = [x]$  als echtes Intervall aufzufassen ist. Bei der Auswertung des Bruches in (C.171) ist zu vermeiden, dass Zähler oder Nenner oder der Bruch selbst einen Überlauf liefern.

Wir betrachten zunächst den Fall:  **$y = 1$**

Wegen

$$\ln \left( 1 + \frac{4}{x^2} \right) = \begin{cases} \ln \left( 1 + \frac{2}{x} \cdot \frac{2}{x} \right), & x > 1 \\ \ln(4 + x^2) - 2 \cdot \ln(x), & 0 < x \leq 1 \end{cases}$$

sind drei Unterfälle zu behandeln:

1.  $Q_1([x]) \subseteq \ln \left( 1 + \frac{2}{[x]} \cdot \frac{2}{[x]} \right)$ , wenn  $\text{Inf}([x]) \in [1, \text{MaxFloat}())$
2.  $Q_1([x]) \subseteq \ln(4 + [x]^2) - 2 \cdot \ln([x])$ , wenn  $\text{Sup}([x]) < 1$ ;
3. Es bleibt der Fall, dass 1 Innen- oder rechter Randpunkt des echten Intervalls  $[x] = [x_1, x_2]$  ist. Mit der Auswertung der Funktionen

$$\begin{aligned} H_1([x_2]) &:= \ln \left( 1 + \frac{2}{[x_2]} \cdot \frac{2}{[x_2]} \right) \\ H_2([x_1]) &:= \ln(4 + [x_1]^2) - 2 \cdot \ln([x_1]) \end{aligned}$$

erhält man die Einschließung:

$$Q_1([x]) \subseteq [\text{Inf}(H_1([x_2])), \text{Sup}(H_2([x_1]))]$$

Im Falle  $y = 1$  kann man so für beliebige Intervalle  $[x] = [x_1, x_2]$ , mit  $x_1 > 0$  eine garantierte Einschließung für  $Q_1([x])$  ohne Auftreten eines Überlaufs berechnen.

Es bleibt der Fall:  $y \neq 1$ .

Zunächst gilt folgendes: Für  $x \rightarrow 0$  und  $y \rightarrow 1$ , mit  $y \neq 1$  bleibt auch bei großen Präzisionen von  $y$

$$(C.172) \quad b := \frac{y}{x^2 + (1 - y)^2}$$

beschränkt, d.h.  $b < \text{MaxFloat}()$ , so dass kein Überlauf entstehen kann. Zur Vermeidung eines vorzeitigen Überlaufs im Nenner muss daher noch der Fall:  $x \rightarrow +\infty$  oder  $y \rightarrow +\infty$  betrachtet werden. Falls der Nenner  $Ne := x^2 + (1 - y)^2$  zum Überlauf führt, muss mit einer geeigneten Zweierpotenz  $2^s$ ,  $s < 0$  so multipliziert werden, dass gerade kein Überlauf mehr auftritt:

$$2^{2s} \cdot Ne = (2^s \cdot x)^2 + (s^s - 2^s \cdot y)^2 =: Nes < \text{MaxFloat}()$$

Beachten Sie bitte, dass bei großem  $y$  und bei zu kleinem  $Nes$  bei der Division

$$bs := \frac{y}{Nes}$$

wieder ein Überlauf auftreten kann! Bei der Rückskalierung mit  $2^{2s+2}$  berücksichtigt der Summand  $+2$  den Faktor 4 in (C.171):

$$2^{2s+2} \cdot bs = \frac{4 \cdot y}{x^2 + (1 - y)^2}$$

Für  $x \rightarrow +\infty$  oder  $y \rightarrow +\infty$  kann so das Argument der `lnp1`-Intervallfunktion ohne Überlauf berechnet werden.

Im Restbereich erhält man dann für  $Q_1([x], [y])$  durch Intervallauswertung von

$$(C.173) \quad Q_1([x], [y]) \subseteq \ln \left[ 1 + \frac{4[y]}{[x]^2 + (1 - [y])^2} \right]$$

eine garantierte Einschließung ohne zwischenzeitlichen Overflow. Für  $y \rightarrow 0$  und  $x \rightarrow \infty$  muss in (C.173) die Funktion `lnp1` mit dem obigen Bruch als Argument zur Anwendung kommen. Weitere Einzelheiten findet man im Quelltext der folgenden Funktion `Q_atan_UPSIGN(...)`, die bei der Auswertung des inversen Tangens zur Anwendung kommt.

Das folgende Beispiel benutzt mit der Current-Präzision `prec = 200` das Argumentintervall  $Z = [0, 0] + i \cdot [\text{pred}(1), \text{pred}(1)]$ , wobei nach (3.1) gilt:  $\text{pred}(1) = 1 - 2^{-\text{prec}} = 1 - 2^{-200}$ . Man erhält die Einschließung mit 60 Dezimalstellen:

$$\begin{aligned} \text{atan}(Z) &= ([0.0, 0.0], \\ &\quad [6.96612916462745035964318282065467450915877635032056530391283e1, \\ &\quad 6.96612916462745035964318282065467450915877635032056530391284e1]), \end{aligned}$$

Erhöht man mit  $Z = [0, 0] + i \cdot [\text{pred}(1), \text{pred}(1)]$  die Current-Präzision z.B. auf `prec = 2000000`, so erhält man nach einigen Sekunden die Einschließung:

$$\begin{aligned} \text{atan}(Z) &= ([0.0, 0.0], \\ &\quad [6.93147527133535589389886830074237297163784172110322434248307e5, \\ &\quad 6.93147527133535589389886830074237297163784172110322434248308e5, ]), \end{aligned}$$

wobei hier bei der Ausgabe die gleiche Dezimalstellenzahl wie im ersten Beispiel gewählt wurde.



Im Algorithmus von  $\text{atan}(Z)$  ist noch folgender Intervallausdruck auszuwerten:

$$\mathbf{D}(s) := 2^{2s} - y^2 - x^2, \quad s \in \mathbb{Z}, \quad x, y \text{ vom Typ } \text{MpfiClass},$$

wobei nur  $x$  ein Punktintervall ist. Zusätzlich wird vorausgesetzt, dass bei der Auswertung von  $2^{2s}$ ,  $y^2$  und  $x^2$  kein Überlauf entsteht. Eine fast optimale Einschließung von  $\mathbf{D}(s)$  wird berechnet mithilfe der Funktion

```
void TwoPow2s_y2_x2(const long s, const MpfiClass& y,
                    const MpfiClass& x, MpfiClass& D),
```

die in `mpfci.cpp` implementiert ist. Es gilt dann  $\mathbf{D}(s) \subseteq D$ .

Um die Problematik einer optimalen Einschließung von  $\mathbf{D}(s)$  zu erläutern, betrachten wir den Spezialfall  $s = 0$ , d.h. einzuschließen ist

$$(C.174) \quad \mathbf{D}_0 := 1 - y^2 - x^2, \quad x, y \text{ vom Typ } \text{MpfiClass}, \text{ wobei nur } x \text{ ein Punktintervall ist.}$$

Prinzipiell ist  $\mathbf{D}_0$  in (C.174) schon optimal, da nach [6, Seite 32] in einem Intervallterm zur optimalen Einschließung die Intervallvariablen jeweils nur einmal auftreten dürfen. In einigen Sonderfällen kann es jedoch zu Auslöschungseffekten kommen, die dann eine optimale Einschließung von  $\mathbf{D}_0$  dennoch verhindern.

Ein **erster Sonderfall** liegt vor, wenn z.B. das Punktintervall  $x$  sehr dicht bei 1 liegt. In diesem Fall wird  $\mathbf{D}_0$  intervallmäßig fast optimal ausgewertet mit

$$D := (1 \diamond x) \diamond (1 \diamond x) - \text{sqr}(y) \supseteq \mathbf{D}_0.$$

Beachten Sie, dass im Intervallausdruck rechts jetzt die Intervallvariable  $x$  zweimal vorkommt. Die damit verbundene Überschätzung ist jedoch ganz minimal, da  $x$  ein Punktintervall ist.

Ein **zweiter Sonderfall** liegt vor, wenn z.B. das echte Intervall  $y$  sehr schmal ist und sehr dicht bei 1 liegt. In diesem Fall wird  $\mathbf{D}_0$  intervallmäßig fast optimal ausgewertet mit

$$D := (1 \diamond y) \diamond (1 \diamond y) - \text{sqr}(x) \supseteq \mathbf{D}_0.$$

In der obigen Funktion `TwoPow2s_y2_x2(...)` findet man weitere Einzelheiten. Die folgenden Beispiele zeigen die gewonnenen Verbesserungen einiger Einschließungen.

Mit  $Z = [-\text{minfloat}(), -\text{minfloat}()] + i \cdot [1 - 2^{-70}, 1 - 2^{-70}]$  erhält man mit der Current-Präzision `prec = 70` die Einschließung

$$\text{atan}(Z) = ([-1.4064180...553497e - 323228476, -1.4064180...553495e - 323228476], \\ [2.46067249098780584842e1, 2.46067249098780584844e1]).$$

Mit  $Z = [-\text{minfloat}(), -\text{minfloat}()] + i \cdot [1 - 2^{-70000}, 1 - 2^{-70000}]$  erhält man mit der Current-Präzision `prec = 70000` die Einschließung

$$\text{atan}(Z) = ([-1.4986879...456363e - 323207425, -1.4986879...456362e - 323207425], \\ [2.42604978931883658022e4, 2.42604978931883658023e4]).$$

`minfloat()` :=  $2^{-1073741824}$  ist die, von der Current-Präzision unabhängige, kleinste positive Maschinenzahl.

### C.2.23. $\operatorname{arccot}(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - i \cdot [-1, +1]$  liefert die Funktion

```
MpfciClass acot(const MpfciClass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\operatorname{acot}(Z)$  für die komplexen Funktionswerte  $\operatorname{arccot}(z)$ , mit  $z \in Z$ .

$$\{\operatorname{arccot}(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \operatorname{acot}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der imaginären Achse von  $-i$  bis  $+i$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben.

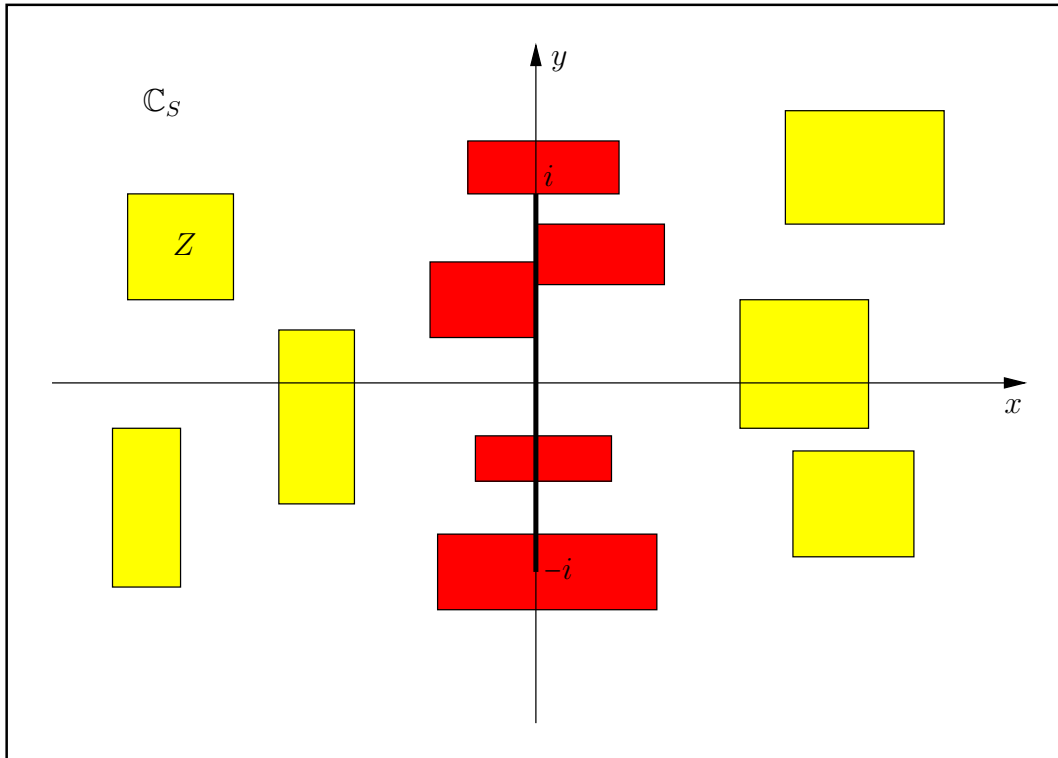


Abbildung C.42.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\operatorname{acot}(Z)$ .

Nach Abbildung C.42 darf also kein Punkt des Verzweigungsschnitts ein Element eines erlaubten, achsenparallelen Intervalls  $Z$  sein.

Mit  $Z = [2^{-1073741824}, 2^{-1073741824}] + i \cdot [1, 1]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 70$  die Einschließung<sup>19</sup>

$$\operatorname{acot}(Z) = ([1.57079632679489661922, 1.57079632679489661924], [-2.46067249098780584844e1, -2.46067249098780584842e1]).$$

Mit  $Z = [0.5, 1.5] + i \cdot [-1, 1]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 80$  die Einschließung

$$\operatorname{acot}(Z) = ([4.63647609000806116214256e - 1, 1.10714871779409050301707], [-7.08303336014054020062385e - 1, 7.08303336014054020062385e - 1]).$$

<sup>19</sup>Die Maschinenzahl  $2^{-1073741824} = \operatorname{minfloat}()$  ist präzisionsunabhängig die kleinste positive Maschinenzahl.

### C.2.24. $\operatorname{arsinh}(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - \{(-\infty, -1) \cup (+1, +\infty)\}$  liefert die Funktion

```
Mpfciclass asinh(const Mpfciclass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\operatorname{arsinh}(Z)$  für die komplexen Funktionswerte  $\operatorname{arsinh}(z)$ , mit  $z \in Z$ .

$$\{\operatorname{arsinh}(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \operatorname{arsinh}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der reellen Achse von  $-1$  bis  $-\infty$  bzw. von  $+1$  bis  $+\infty$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben. Die Implementierung erfolgt mit

$$\operatorname{arsinh}(Z) = i \cdot \operatorname{asin}(-i \cdot Z);$$

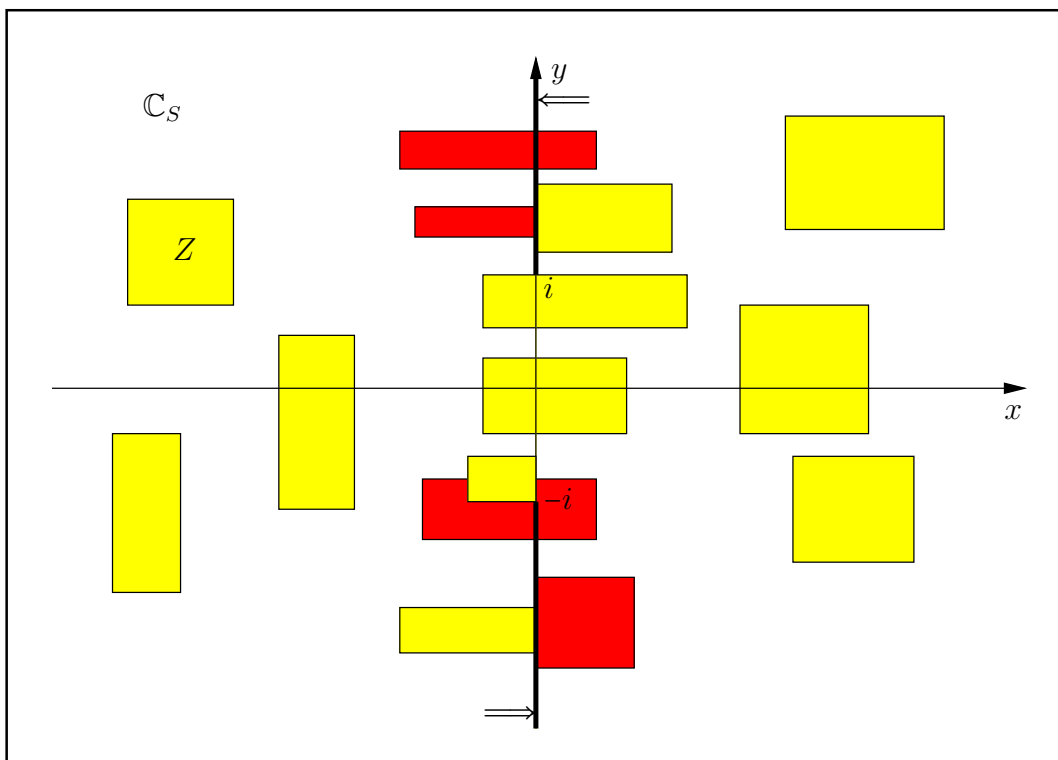


Abbildung C.43.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\operatorname{arsinh}(Z)$ .

Mit  $Z = [2^{1073741821}, 2^{1073741821}] + i \cdot [2^{1073741821}, 2^{1073741821}]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 80$  die Einschließung

$$\operatorname{arsinh}(Z) = ([7.44261116915172247033984e8, 7.44261116915172247033986e8], [7.85398163397448309615658e-1, 7.85398163397448309615664e-1]).$$

Mit  $Z = [-0.5, 0.5] + i \cdot [-1, -0.5]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 80$  die Einschließung

$$\operatorname{arsinh}(Z) = ([-7.32857675973645260888675e-1, 7.32857675973645260888675e-1], [-1.57079632679489661923133, -4.52278447151190682063657e-1]).$$

### C.2.25. $\operatorname{arcosh}(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - \{(-\infty, +1)\}$  liefert die Funktion

```
MpfciClass acosh(const MpfciClass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\operatorname{acosh}(Z)$  für die komplexen Funktionswerte  $\operatorname{arcosh}(z)$ , mit  $z \in Z$ .

$$\{\operatorname{arcosh}(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \operatorname{acosh}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der reellen Achse von  $-\infty$  bis  $+1$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben. Die Implementierung erfolgt mit

$$\operatorname{acosh}(Z) = i \cdot \operatorname{acos}(Z) = \pm i \cdot (\pi/2 - \operatorname{asin}(Z));$$

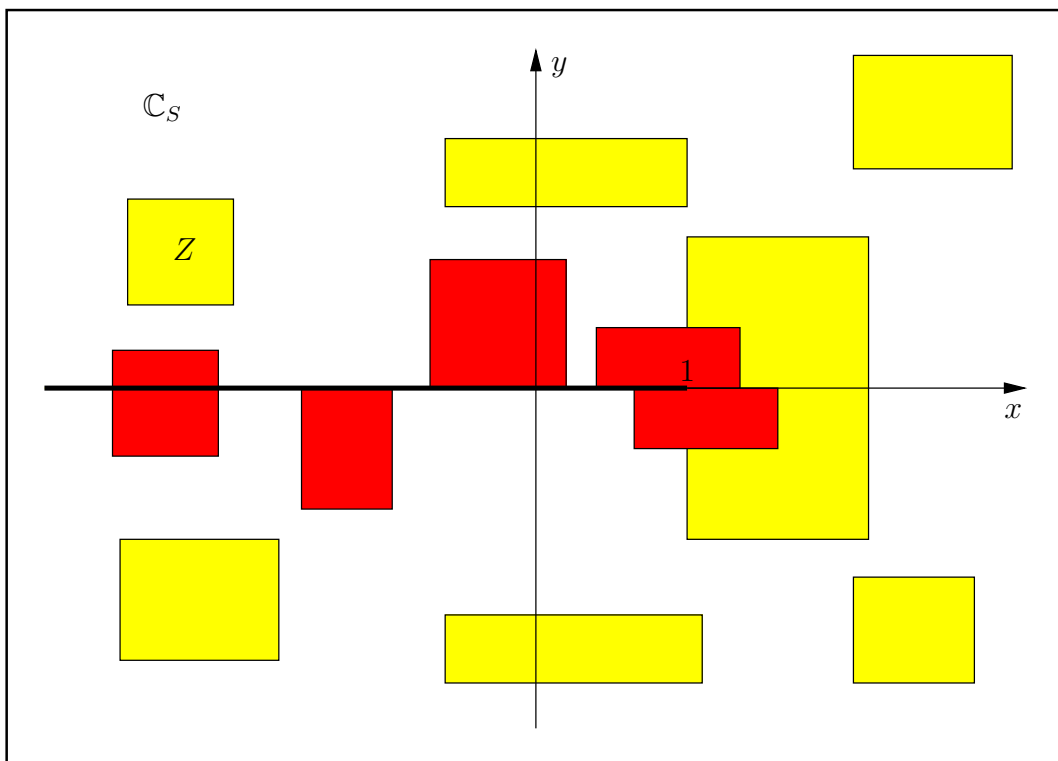


Abbildung C.44.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\operatorname{acosh}(Z)$ .

Mit  $Z = [1 - 2^{-70}, 1 - 2^{-70}] + i \cdot [2^{-1073741824}, 2^{-1073741824}]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 70$  die Einschließung

$$\operatorname{acosh}(Z) = ([5.78868064589607735286e - 323228487, 5.78868064589607735289e - 323228487], [4.11590317489199529168e - 11, 4.11590317489199529171e - 11]).$$

Mit  $Z = [1, 2] + i \cdot [-1, 0]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 150$  die Einschließung

$$\operatorname{acosh}(Z) = ([0, 1.46935174436818527325584431736164761678780335733478818], [-9.04556894302381364127316795661958721431094560961605070e - 1, 0]).$$

### C.2.26. $\operatorname{artanh}(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - \{(-\infty, -1) \cup (+1, +\infty)\}$  liefert die Funktion

```
MpfciClass atanh(const MpfciClass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\operatorname{atanh}(Z)$  für die komplexen Funktionswerte  $\operatorname{artanh}(z)$ , mit  $z \in Z$ .

$$\{\operatorname{artanh}(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \operatorname{atanh}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der reellen Achse von  $-1$  bis  $-\infty$  bzw. von  $+1$  bis  $+\infty$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben.

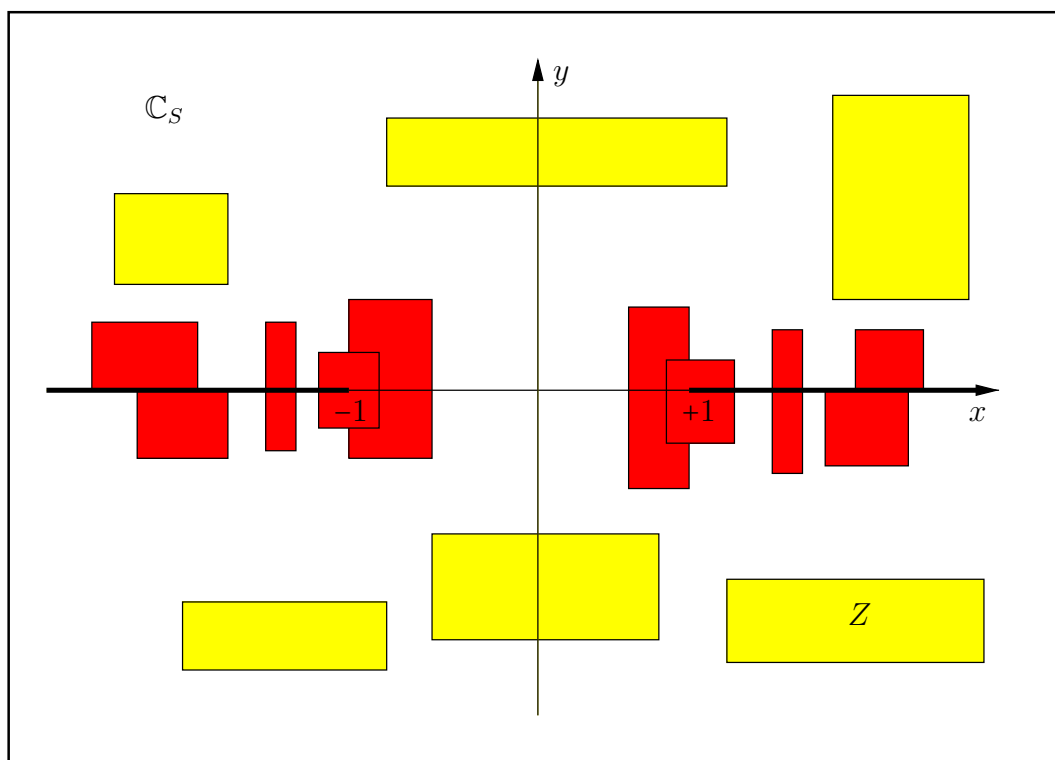


Abbildung C.45.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\operatorname{atanh}(Z)$ .

Mit  $Z = [1 - 2^{-70}, 1 - 2^{-70}] + i \cdot [2^{-1073741824}, 2^{-1073741824}]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 70$  die Einschließung

$$\operatorname{atanh}(Z) = ([2.46067249098780584842e1, 2.46067249098780584844e1], [1.4064180...553495e - 323228476, 1.4064180...553497e - 323228476]).$$

Mit  $Z = [1 - 2^{-70}, 1 - 2^{-70}] + i \cdot [2^{3000}, 2^{3000}]$  erhält man mit der Current-Präzision  $\operatorname{prec} = 150$  die Einschließung

$$\operatorname{atanh}(Z) = ([6.60733027580565499208e - 1807, 6.60733027580565499209e - 1807], [1.57079632679489661923, 1.57079632679489661924]),$$

wobei die Ausgabe wie im ersten Beispiel auf 21 Dezimalstellen begrenzt wurde.

### C.2.27. $\operatorname{arcoth}(z)$

Mit dem achsenparallelen Rechteckintervall  $Z \subset \mathbb{C}_S := \mathbb{C} - i \cdot [-1, +1]$  liefert die Funktion

```
Mpfciclass acoth(const Mpfciclass& Z);
```

die folgende achsenparallele Rechteckeinschließung  $\operatorname{acoth}(Z)$  für die komplexen Funktionswerte  $\operatorname{arcoth}(z)$ , mit  $z \in Z$ .

$$\{\operatorname{arcoth}(z) \mid z \in Z \subset \mathbb{C}_S\} \subseteq \operatorname{acoth}(Z).$$

$\mathbb{C}_S$  ist dabei die längs der imaginären Achse von  $-i$  bis  $+i$  aufgeschnittene komplexe Ebene. In der folgenden Abbildung sind einige erlaubte und **nicht erlaubte** Rechteckintervalle  $Z$  angegeben.

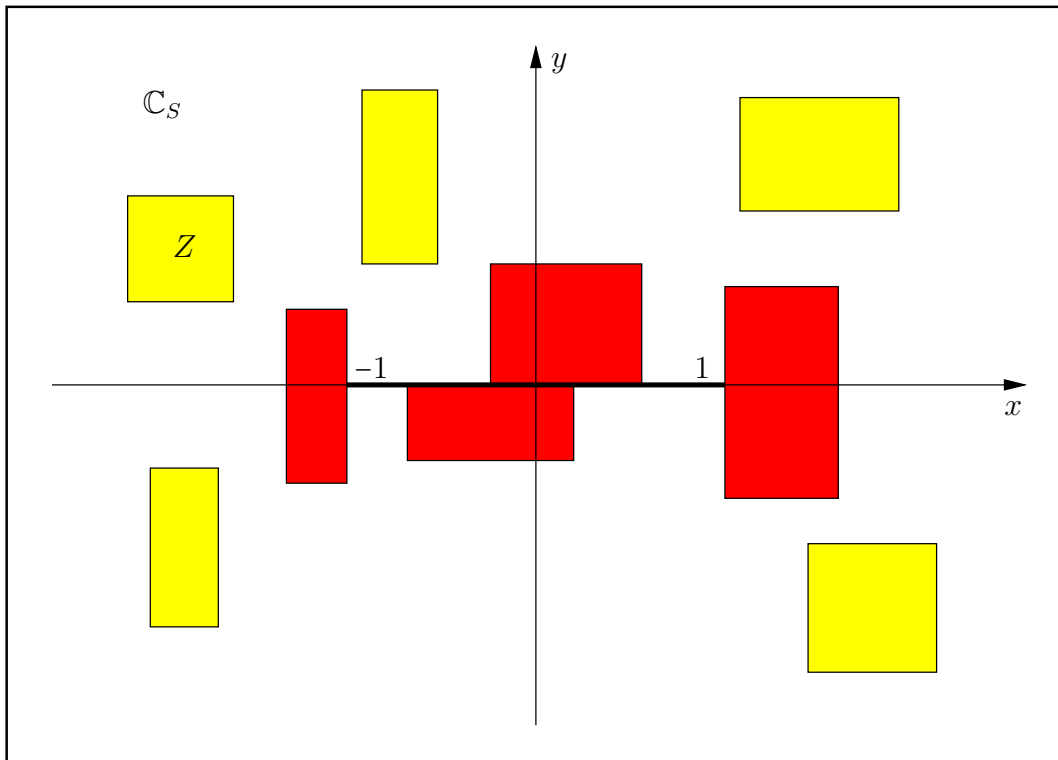


Abbildung C.46.: Erlaubte und **nicht erlaubte** Intervalle  $Z$  von  $\operatorname{acoth}(Z)$ .

Mit  $Z = [1, 1] + i \cdot [2^{-1073741824}, 2^{-1073741824}]$  erhält man mit der Current-Präzision  $\text{prec} = 70$  die Einschließung

$$\operatorname{acoth}(Z) = ([3.72130559324020099216e8, 3.72130559324020099218e8], [-7.85398163397448309617e-1, -7.85398163397448309615e-1]).$$

Mit  $Z = [1, 1] + i \cdot [2^{3000}, 2^{3000}]$  erhält man mit der Current-Präzision  $\text{prec} = 70$  die Einschließung

$$\operatorname{acoth}(Z) = ([6.60733027580565499207e-1807, 6.60733027580565499209e-1807], [-8.12854862555773544048e-904, -8.12854862555773544046e-904]),$$

wobei die Ausgabe wie im ersten Beispiel auf 21 Dezimalstellen begrenzt wurde.

### C.2.28. $z^p, p \in \mathbb{P}:\text{MpfiClass}$

Für  $z = |z| \cdot e^{i\varphi}$ ,  $-\pi/2 < \varphi < +3\pi/2$  definieren wir mit  $z \in \mathbb{Z}:\text{MpfiClass}$  und  $p \in \mathbb{P}:\text{MpfiClass}$  die Potenz

$$(C.175) \quad z^p := e^{p \cdot \ln(z)} = e^{p \cdot \ln|z|} \cdot e^{i \cdot p(\varphi + 2\pi k)}, \quad k \in \mathbb{Z}.$$

Wegen der Mehrdeutigkeit des komplexen Logarithmus, d.h. wegen  $k \in \mathbb{Z}$ , gibt es damit für festes  $z \neq 0$  und  $p \in \mathbb{R} - \mathbb{Q}$  beliebig viele Potenzen  $z_k^p$ , die alle auf einem Kreis um den Ursprung mit dem Radius  $R = e^{p \cdot \ln|z|}$  liegen. Die Aufgabe besteht nun darin, alle diese Potenzen für alle  $z \in \mathbb{Z}$  und für alle  $p \in \mathbb{P}$  möglichst optimal einzuschließen, d.h. einzuschließen ist die Menge

$$T := \{y \in \mathbb{C} \mid y = e^{p \cdot \ln|z|} \cdot e^{i \cdot p(\varphi + 2\pi k)}, \quad k \in \mathbb{Z}, z \in \mathbb{Z}, p \in \mathbb{P} = [p_1, p_2]\}.$$

Zur Einschließung von  $T$  sind drei Fälle zu unterscheiden:

**Fall 1:**  $0 \notin \mathbb{Z}$ .

$T$  ist jetzt ein Kreisring, der durch folgende Radien bestimmt ist.

$$r_1 = e^{\text{Inf}(\mathbb{P} \ln|z|)}, \quad r_2 = e^{\text{Sup}(\mathbb{P} \ln|z|)}.$$

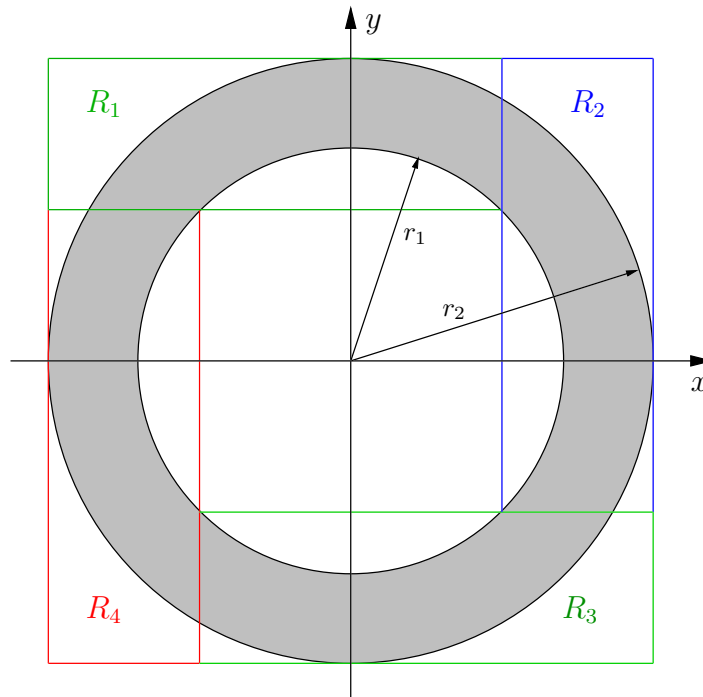


Abbildung C.47.: Den Kreisring einschließende Rechtecke  $R_\nu$ ,  $\nu = 1, 2, 3, 4$ .

Die vier Rechtecke  $R_\nu$ ,  $\nu = 1, 2, 3, 4$ , schließen den Kreisring  $T$  ein und werden von der Funktion

```
std::list<MpfiClass> pow_all( const MpfiClass& Z, const MpfiClass& P )
```

in einer Liste bereitgestellt. Beachten Sie bitte, dass bei kleinen Differenzen  $r_2 - r_1$  die Einschließung von  $T$  mit einem einzelnen Quadrat der Seitenlänge  $2r_2$  zu großen Überschätzungen führen würde. Darüber hinaus wäre die Null ein Element der Einschließung, womit weitere Komplikationen in nachfolgenden Rechnungen verbunden sein könnten.

Das folgende Programm `MPFR-09.cpp` berechnet mit  $Z = [1, 1.125] + i[1, 1.25]$  und mit dem Exponentenintervall  $P = [1, 1.5]$  die vier einschließenden Intervalle  $R_\nu$ ,  $\nu = 1, 2, 3, 4$ .

```

1 // MPFR-09.cpp
2 #include "mpfciclass.hpp"
3
4 using namespace MPFR;
5 using namespace MPFI;
6 using namespace cxsc;
7 using namespace std;
8
9 int main(void)
10 {
11     MPFI::MpfiClass::SetCurrPrecision(40);
12     cout << "GetCurrPrecision() = " << MPFI::MpfiClass::GetCurrPrecision() << endl;
13     MpfiClass Z(interval(1,1.125), interval(1,1.25), 53);
14     MpfiClass P(interval(1,1.5), 53);
15     cout.precision(40/3.321928095); // Ausgabe mit 9 Dez.-Stellen
16     cout << "Z = " << Z << endl;
17     cout << "Z.GetPrecision() = " << Z.GetPrecision() << endl;
18     cout << "P = " << P << endl;
19     cout << "P.GetPrecision() = " << P.GetPrecision() << endl;
20     cout << "Einschliessung aller Potenzen:" << endl;
21
22     list<MpfiClass> res;
23     res = pow_all(Z, P);
24
25     list<MpfiClass>::iterator pos;
26     // Ausgabe der n Potenzen:
27     int k(0);
28     for (pos = res.begin(); pos != res.end(); ++pos )
29     {
30         k++;
31         cout << "R" << k << " = " << *pos << endl; // Jedes Rechteck R in neue Zeile
32         cout << "Praezision = " << (*pos).GetPrecision() << endl;
33     }
34
35     return 0;
36 }

```

Das Programm liefert die Ausgabe

```

GetCurrPrecision() = 40
Z = ([1.00000000000,1.1250000000], [1.0000000000,1.2500000000])
Z.GetPrecision() = 53
P = [1.0000000000,1.5000000000]
P.GetPrecision() = 53

Einschliessung aller Potenzen:

R1 = ([9.9999999999e-1,2.18084073549], [-1.0000000000,2.18084073549])
Praezision = 40
R2 = ([-2.18084073549,1.0000000000], [9.9999999999e-1,2.18084073549])
Praezision = 40
R3 = ([-2.18084073549,-9.9999999999e-1], [-2.18084073549,1.0000000000])
Praezision = 40
R4 = ([-1.0000000000,2.18084073549], [-2.18084073549,-9.9999999999e-1])
Praezision = 40

```

mit den vier Rechteckintervallen  $R_\nu$ ,  $\nu = 1, 2, 3, 4$ .

Beachten Sie, dass sich nach (C.175) im Fall  $P = [p, p]$ , mit  $p \in \mathbb{Q}$  und bei Punktintervallen  $Z$  nur endlich viele Potenzwerte  $z_k^p$  auf dem Kreis mit dem Radius  $e^{p \ln |z|}$  befinden. Beispielsweise gilt  $1^2 = 1$ , aber `pow_all(1,2)` liefert wieder die vier einschließenden Rechtecke. Diese Spezialfälle werden also in `pow_all` nicht exktra behandelt!



**Fall 2:**  $0 \in Z$ ,  $p_1 > 0$ .

$T$  ist jetzt eine Kreisscheibe, die mit  $Z \rightarrow 0$  in den Ursprung übergeht. Diese Kreisscheibe wird jetzt mit `pow_all()` durch ein einziges Quadrat der Ergebnisliste optimal eingeschlossen.

**Fall 3:**  $0 \in Z$ ,  $p_1 \leq 0$ .

Da  $0^p$  für  $p \leq 0$  undefiniert ist, liefert `pow_all` jetzt eine Fehlermeldung mit Programmabbruch.

### Anmerkungen:

1. `pow_all()` berechnet eine Einschließung aller Potenzen  $z^p$ , mit  $z \in Z$  und  $p \in P$ , wobei  $P$  ein reelles Intervall sein muss. Würde man im Gegensatz dazu für den Exponenten einen von Null verschiedenen Imaginärteil zulassen, so wäre in (C.175) mit  $p = p_1 + i \cdot p_2$  der erste Faktor rechts gegeben durch

$$e^{p_1 \cdot \ln|z| - p_2(\varphi + 2\pi k)}, \quad k \in \mathbb{Z},$$

und wegen  $k \in \mathbb{Z}$  würden damit die Potenzen  $z_k^p$  über alle Grenzen wachsen und damit eine sinnvolle Einschließung verhindern. Bei der Funktion `pow(Z,W)` wählt man daher  $k = 0$  und berechnet damit nur Einschließungen des Hauptwertes.

## C.2.29. $e^z - 1$

Wir betrachten die Aufgabe, zu einem vorgegebenen komplexen Intervall

$$Z = X + i \cdot Y, \quad X = [x_1, x_2], \quad Y = [y_1, y_2]$$

eine möglichst optimale Einschließung der Menge  $W := \{w \in \mathbb{C} \mid w = e^z - 1, \quad z \in Z\}$  durch ein komplexes Rechteck-Intervall  $F \supseteq W$  mit der Präzision  $\text{prec} \geq 2$  zu berechnen. Beachten Sie dabei jedoch, dass nach Seite 243 mit einem solchen einschließenden Rechteck  $F = U + i \cdot V$  oft deutliche Überschätzungen nicht zu vermeiden sind.

### C.2.29.1. Realteil

Das gesuchte Einschließungsintervall  $U$  für den Realteil von  $e^Z - 1$  wird mit dem Maschinenintervall  $\cos(Y) = [c_1, c_2]$  berechnet durch:

$$(C.176) \quad U := e^X \diamond \cos(Y) \diamond 1 = e^{[x_1, x_2]} \diamond [c_1, c_2] \diamond 1.$$

Da die Subtraktion der Eins zu Auslöschungsproblemen führen kann, lässt sich in Abhängigkeit der Vorzeichen von  $c_1, c_2$  das Intervall  $U$  wie folgt berechnen:

$$(C.177) \quad U = \begin{cases} [e^{x_1 + \ln(c_1)} - 1, e^{x_2 + \ln(c_2)} - 1], & 0 \leq c_1 \leq c_2, \\ [c_1 \cdot e^{x_2} - 1, c_2 \cdot e^{x_1} - 1], & c_1 \leq c_2 \leq 0, \\ [c_1 \cdot e^{x_2} - 1, e^{x_2 + \ln(c_2)} - 1], & c_1 \leq 0 \leq c_2. \end{cases}$$

Dabei ist klar, dass eine Auslöschung nur in den drei Fällen auftreten kann, in denen die Summanden  $\ln(c_1)$  oder  $\ln(c_2)$  im Exponenten vorkommen. Bei der Berechnung dieser **Summanden** tritt jedoch noch folgendes Problem auf, das im Fall  $\ln(c_1)$  näher zu betrachten ist. Es gilt zunächst nach Definition:

$$\ln(c_1) := \ln(\text{Inf}(\cos(Y))),$$

und da die Logarithmusfunktion monoton wachsend ist, folgt zusätzlich

$$(C.178) \quad \ln(c_1) := \ln(\text{Inf}(\cos(Y))) = \text{Inf}(\text{ln\_cos}(Y)).$$

Wertet man jetzt die beiden Ausdrücke rechts in (C.178) auf der Maschine aus mit der Präzision  $\text{prec} = 180$  und dem Punktintervall<sup>20</sup>  $Y = [10^{-300}, 10^{-300}]$ , so erhält man die deutlich verschiedenen Ergebnisse

$$\begin{aligned} \widetilde{\text{Inf}}(\text{ln\_cos}(Y)) &= -5.00000e - 601, & \text{prec} &= 180, \\ \widetilde{\text{ln}}(\text{Inf}(\cos(Y))) &= -6.52530e - 55, & \text{prec} &= 180, \\ \widetilde{\text{ln}}(\text{Inf}(\cos(Y))) &= -5.00000e - 601, & \text{prec} &= 2160. \end{aligned}$$

Erst bei 12-facher Current-Präzision, d.h. mit  $\text{prec} = 2160$ , erhält man für den optimalen Wert  $\ln(c_1)$  mit dem Ausdruck  $\widetilde{\text{ln}}(\text{Inf}(\cos(Y)))$  das exakte Ergebnis  $-5 \cdot 10^{-601}$ , so dass aus Laufzeitgründen  $\alpha := \text{Inf}(\text{ln\_cos}(Y))$  auszuwerten ist! Entsprechendes gilt für<sup>21</sup>  $\beta := \text{Sup}(\text{ln\_cos}(Y))$ . Damit wird der Realteil viel besser eingeschlossen durch die Maschinenauswertung von

$$(C.179) \quad U = \begin{cases} [e^{x_1 + \alpha} - 1, e^{x_2 + \beta} - 1], & 0 \leq c_1 \leq c_2, \\ [c_1 \cdot e^{x_2} - 1, c_2 \cdot e^{x_1} - 1], & c_1 \leq c_2 \leq 0, \\ [c_1 \cdot e^{x_2} - 1, e^{x_2 + \beta} - 1], & c_1 \leq 0 \leq c_2. \end{cases}$$

<sup>20</sup>Mit  $10^{-300}$  ist hier die zu  $10^{-300}$  nächstgelegene IEEE-Maschinenzahl gemeint.

<sup>21</sup> $\text{ln\_cos}(\dots)$  ist die bereits in `mpfi.cpp` definierte Intervallfunktion zur Einschließung von  $\ln(\cos(Y))$ .

Bei der Auswertung der beiden Summen  $x_1 + \alpha$  und  $x_2 + \beta$  kann bei jeweils entgegengesetzten Vorzeichen beider Summanden im Exponenten erhebliche Auslöschung auftreten. In diesen Fällen ist die Präzision so lange zu verdoppeln, bis für die jeweilige Summe, also z.B. für  $s := x_1 + \alpha$ , die durch `prec` vorgegebene Genauigkeit erreicht wird. Um für  $e^{x_1 + \alpha} - 1$  den optimal abgerundeten Funktionswert zu erhalten, muss vorher natürlich, wie oben beschrieben, die **abgerundete** Summe  $x_1 + \alpha$  berechnet worden sein. Anschließend wird dann mit der schon implementierten Punktfunktion `expm1(s, RoundDown)` der gesuchte, optimal abgerundete Wert von  $e^{x_1 + \alpha} - 1$  berechnet. Ganz entsprechend wird der aufzurundende Ausdruck  $e^{x_2 + \beta} - 1$  ausgewertet. Weitere Einzelheiten findet man in den Funktionen

```
MPFR::MpfrClass expm1_H1(const MPFR::MpfrClass& x, const Mpfciclass& z,
                        bool& Ausl, RoundingMode rnd);
MPFR::MpfrClass expm1_H(const MPFR::MpfrClass& x, const Mpfciclass& z,
                        RoundingMode rnd);
```

die in `mpfciclass.cpp` definiert sind.

### C.2.29.2. Imaginärteil

Da bei der Berechnung des Imaginärteils von  $e^Z - 1$  keine Auslöschung auftritt, kann das Einschließungsintervall  $V$  intervallmäßig direkt und optimal berechnet werden:

$$(C.180) \quad V := e^X \diamond \sin(Y),$$

wobei jetzt  $e^X$  und  $\sin(Y)$  die intervallmäßig berechneten Funktionswerte bedeuten.

### C.2.29.3. Numerische Ergebnisse

Wir beginnen mit zwei Beispielen, in denen bei der Berechnung der Summen  $x_1 + \alpha$  und  $x_2 + \beta$  keine Auslöschung auftreten kann.

Im **1. Beispiel** wählen wir  $x_1 = x_2 = y_1 = y_2 = 2^{-1000000000}$ . Damit haben die Summanden  $x_1, \alpha$  bzw.  $x_2, \beta$  zwar unterschiedliche Vorzeichen, ihre Beträge sind jedoch so unterschiedlich, dass eine Auslöschung nicht eintreten kann. Beachten Sie jedoch, dass bei einer naiven Berechnung von z.B.  $e^{x_1 + \alpha} - 1$  wegen Auslöschung bei der Differenzberechnung sehr starke Überschätzungen auftreten würden, die jedoch durch Anwendung der `expm1`-Funktion bei dieser Implementierung vermieden wird. Mit `prec = 180` und  $Z = [x_1, x_2] + i \cdot [y_1, y_2]$  erhalten wir die Einschließung

$$e^Z - 1 \subseteq ([2.1677979676169...8690 \cdot 10^{-301029996}, 2.1677979676169...8692 \cdot 10^{-301029996}], \\ [2.1677979676169...8690 \cdot 10^{-301029996}, 2.1677979676169...8692 \cdot 10^{-301029996}]).$$

Beachten Sie, wie dicht diese fast optimale Einschließung von  $e^Z - 1$  am Ursprung liegt.

Das **2. Beispiel** zeigt, dass  $Z$  kein Punktintervall sein muss. Wir wählen  $x_1 = y_1 = 2^{-1000000000}$  und  $x_2 = y_2 = 2 \cdot x_1$ . Mit `prec = 180` und  $Z = [x_1, x_2] + i \cdot [y_1, y_2]$  erhalten wir die Einschließung

$$e^Z - 1 \subseteq ([2.1677979676169...8690 \cdot 10^{-301029996}, 4.3355959352338...7383 \cdot 10^{-301029996}], \\ [2.1677979676169...8690 \cdot 10^{-301029996}, 4.3355959352338...7383 \cdot 10^{-301029996}]).$$

Auch jetzt wird das Einschließungsintervall fast optimal berechnet und liegt ebenfalls sehr dicht am Ursprung.

Im **3. Beispiel** wählen wir zunächst  $y_1 = y_2 = 0.75$ . Um starke Auslöschung bei der Berechnung von z.B.  $x_1 + \alpha$  zu konstruieren, berechnen wir mit  $Y = [y_1, y_2]$  und  $\text{prec} = 1800$  zunächst eine Einschließung  $T$  von  $-\ln(\cos(Y)) \subseteq T$  und wählen mit  $x_1 = x_2 = \text{mid}(T)$  das Punktintervall  $X = [x_1, x_2]$ . Wegen  $\text{mid}(T) \approx -\alpha := \text{Inf}(\ln_{\cos}(Y))$  ist daher bei der Berechnung der Summe  $x_1 + \alpha$  erhebliche Auslöschung zu erwarten, vgl. (C.179) auf Seite 370. Für  $e^Z - 1$  erhalten wir die Einschließung

$$e^Z - 1 \subseteq ([-2.39983652734278649...3587 \cdot 10^{-544}, -2.39983652734278649...3586 \cdot 10^{-544}], \\ [9.31596459944072461165...0253 \cdot 10^{-1}, 9.31596459944072461165...0255 \cdot 10^{-1}]).$$

Auch diese Einschließung ist mit 541 korrekten Dezimalziffern nahezu optimal.

Im **4. Beispiel** konstruieren wir ganz analog zum 3. Beispiel noch stärkere Auslöschung und wählen  $y_1 = y_2 = 2^{-40000}$ . Mit  $\text{prec} = 1800$  erhalten wir dann für  $e^Z - 1$  nach etwa einer Minute die Einschließung

$$e^Z - 1 \subseteq ([-1.322855833975488484...2062 \cdot 10^{-48166}, -1.322855833975488484...2061 \cdot 10^{-48166}], \\ [6.3120937524672703508...5709 \cdot 10^{-12042}, 6.3120937524672703508...5710 \cdot 10^{-12042}]).$$

Auch diese Einschließung ist mit 541 korrekten Dezimalziffern nahezu optimal.

Im **5. Beispiel** bestätigen wir die bekannte Gleichung

$$(C.181) \quad e^{i \cdot 2\pi} - 1 = 0$$

und wählen dazu  $x_1 = x_2 = 0$ . Mit der Anweisung

```
MpfiClass Y = 2*MPFI::MpfiClass::Pi();
```

und  $\text{prec} = 1800$  schließt  $Y = [y_1, y_2]$  den Wert  $2\pi$  optimal ein, und für  $e^{i \cdot 2\pi} - 1$  erhalten wir die Einschließung

$$e^{i \cdot 2\pi} - 1 \in ([-4.2241520191547511029602402287327...0180 \cdot 10^{-1083}, 0], \\ [-9.191465627586006593...6038 \cdot 10^{-542}, 2.005434173965873156...9529 \cdot 10^{-542}]).$$

Mit wachsender Präzision  $\text{prec}$  wird  $e^{i \cdot 2\pi} - 1 = 0$  immer enger eingeschlossen, wobei zu beachten ist, dass die optimale Einschließung  $([0, 0], [0, 0])$  natürlich nicht erreicht werden kann, da  $2\pi$  auf der Maschine **nicht** darstellbar ist. Mit  $\text{prec} = 5\,000\,000$  erhält man nach einigen Minuten die sehr enge Einschließung

$$e^{i \cdot 2\pi} - 1 \in ([-2.93976...10^{-3010299}, 0], [-7.41700...10^{-1505151}, 7.66780...10^{-1505150}]).$$

Die größere Laufzeit wird verursacht durch die optimale Einschließung von  $2\pi$  mit ca. 1 505 149 Dezimalstellen. Beachten Sie, dass das Infimum des Realteils im Vergleich zum Imaginärteil fast den doppelten negativen Zehnerexponenten erhält und dass das Supremum **0** sogar exakt berechnet wird.

# D. Einschließung reeller und komplexer Ausdrücke

## D.1. Einschließung reeller arithmetischer Ausdrücke

### D.1.1. Problemstellung und Definitionen

In einem *reellen arithmetischen Ausdruck* sind per Definition reelle Konstanten und Funktionen aus Tabelle 3.11.3 mit den vier arithmetischen Grundoperationen verknüpft. Ist dann die reelle Funktion  $f : D_f \rightarrow \mathbb{R}$  ein solcher reeller arithmetischer Ausdruck und ist  $X = [x_1, x_2] \subseteq D_f \subseteq \mathbb{R}$  ein vorgegebenes Intervall, so ist

$$(D.1) \quad W_{f,X} := \{f(x) \mid x \in X \subset D_f\}$$

die Wertemenge  $W_{f,X}$  von  $f$  über dem reellen Intervall  $X \subset D_f$ . Ist  $f$  über dem abgeschlossenen Intervall  $X \subset D_f$  stetig, so gibt es Zahlen  $\xi_1, \xi_2 \in X$ , so dass gilt  $W_{f,X} = f(X) := [f(\xi_1), f(\xi_2)]$ , d.h. die Wertemenge ist ein endliches Intervall. Beachten Sie, dass die Stetigkeit von  $f(x)$  über  $X \subset D_f$  für arithmetische Ausdrücke mit Funktionen aus Tabelle 3.11.3 keine einschränkende Bedingung ist.

Bei der Intervallrechnung verfolgt man u.a. das Ziel, zu einer Funktion  $f(x)$  und zu einem gegebenen Maschinenintervall  $X \subset D_f$  mit einem möglichst einfachen Algorithmus die Wertemenge  $W_{f,X}$  durch ein Maschinenintervall  $I(X)$  garantiert einzuschließen, d.h. es soll gelten  $W_{f,X} \subseteq I(X)$ . Nach [6, Seite 31] erhält man eine solche Maschineneinschließung  $I(X)$ , wenn der Funktionsterm  $f(x)$  intervallmäßig ausgewertet wird, d.h. wenn die auftretenden Variablen  $x$  jeweils durch  $X$  und die Grundoperationen durch ihre entsprechenden Intervalloperationen ersetzt werden. Wir betrachten dazu das folgende Beispiel:

$$f(x) = 3x \cdot e^{-2x+1}, \quad X = [0, 2], \quad \text{mit: } W_{f,X} = [0, 3/2], \quad \text{und die intervallmäßige Auswertung liefert} \\ I(X) = 3 \diamond X \diamond \exp(-2 \diamond X \diamond 1) = [0, 16.309690 \dots] \supset W_{f,X}.$$

Die intervallmäßige Auswertung liefert bei diesem Beispiel mit  $I(X)$  eine erhebliche Überschätzung des Wertebereichs  $W_{f,X} = [0, 3/2]$ , da der Funktionsterm  $f(x)$  die Variable  $x$  zweimal enthält und weil die Intervallrechnung grundsätzlich annimmt, dass, im Gegensatz zu unserem Beispiel, die auftretenden Intervalle  $X$  beide voneinander unabhängig sind. Das Ziel ist also,  $I(X) = [I_1, I_2]$  so zu verbessern, dass im Idealfall  $I(X) = f(X) := [f(\xi_1), f(\xi_2)]$  erfüllt ist. Beachten Sie aber, dass  $I(X) = f(X)$  nur gelten kann, wenn  $f(\xi_1)$  und  $f(\xi_2)$  beides Maschinenzahlen sind. Nach Definition verlangen wir daher für eine **optimale** Einschließung des Wertebereichs  $W_{f,X} = f(X) := [f(\xi_1), f(\xi_2)]$  durch das Maschinenintervall  $I(X) = [I_1, I_2]$  die beiden folgenden Doppelungleichungen:

$$(D.2) \quad I_1 \leq f(\xi_1) < \text{succ}(I_1), \quad \text{pred}(I_2) < f(\xi_2) \leq I_2, \quad \xi_1, \xi_2 \in X.$$

Im folgenden Abschnitt wird gezeigt, wie man die durch intervallmäßige Auswertung berechnete Einschließung  $I(X)$  eines Wertebereichs weiter so verbessern kann, dass die verbesserte Einschließung dann **optimal** oder nahezu **optimal** ist. Als nahezu optimal bezeichnen wir eine Einschließung, wenn z.B. gilt  $I_1 \leq f(\xi_1) < \text{succ}(\text{succ}(I_1))$ .

### D.1.2. Verbesserung einer Einschließung

Eine erste Methode zur Verbesserung der Einschließung des Wertebereichs  $W_{f,X}$  wird realisiert durch Anwendung der *Zentrischen Form* oder der *Mittelwertform*, siehe z.B. [5], [56], [58]. Da diese Methoden jedoch nur bei hinreichend schmalen Intervallen  $X$  wirklich hilfreich sind, wird hier nicht näher darauf eingegangen. Eine zweite Methode basiert auf dem folgenden Satz<sup>1</sup>, [6, Seite 32]

**Satz D.1** (Eine Variable tritt explizit nur einmal auf).

Ist  $f(x)$  ein reeller, arithmetischer Ausdruck und tritt in ihm die Variable  $x$  explizit nur einmal auf, so liefert die intervallmäßige Auswertung von  $f(x)$  über dem Intervall  $X \subset D_f$  eine optimale Einschließung des Wertebereichs  $W_{f,X}$ , d.h. es gilt:  $W_{f,X} = [f(\xi_1), f(\xi_2)]$ ,  $\xi_1, \xi_2 \in X$ .

Eine unmittelbare Anwendung dieses Satzes besteht darin, in einem arithmetischen Ausdruck  $f(x)$  die Anzahl der auftretenden Variablen  $x$  möglichst zu minimieren. Als Beispiel betrachten wir die beiden äquivalenten Ausdrücke

$$f(x) = \frac{1+x}{x}, \quad g(x) = 1 + \frac{1}{x}, \quad \text{mit } X = [1, 4].$$

Durch intervallmäßige Auswertung erhalten wir die Ergebnisse

$$I_f(X) := (1 \diamond X) \diamond X = [0.5, 5]; \quad I_g(X) := 1 \diamond 1 \diamond X = [1.25, 2] = W_{g,X} \subset I_f(X),$$

wobei  $I_g(X)$  sogar die optimale Einschließung des Wertebereichs von  $f(x) \equiv g(x)$  über  $X$  ist, da  $g(x)$  die Variable  $x$  nur einmal enthält.

Eine zweite Anwendung von Satz D.1 ist ein quadratischer Ausdruck der Form

$$f(x) := x^2 + a \cdot x + b = (x + a/2)^2 - a^2/4 + b =: g(x), \quad x, a, b \in \mathbb{R}.$$

Wertet man jetzt  $g(x)$  intervallmäßig aus und benutzt dabei z.B.  $\text{sqr}(X \diamond a \diamond 2)$  für den ersten Summanden  $(x + a/2)^2$ , so erhält man mit  $I_g(X)$  eine optimale Einschließung des Wertebereichs  $W_{f,X}$  über einem Intervall  $X$ , wenn man zur Vermeidung einer möglichen Auslöschung  $g(x)$  insgesamt in hinreichend hoher Präzision auswertet. Beachten Sie dabei, dass man bei der intervallmäßigen Auswertung von  $g(x)$  für  $b$  ein beliebiges aber endliches Intervall zulassen darf, während man für eine optimale Einschließung  $a$  nur durch ein Punktintervall ersetzen darf. Als Beispiel wählen wir

$$f(x) := x^2 - 4 \cdot x - 4 = (x - 2)^2 - 8 =: g(x), \quad X = [1, 4].$$

Durch intervallmäßige Auswertung erhalten wir die Ergebnisse

$$I_f(X) := \text{sqr}(X) \diamond 4 \diamond X \diamond 4 = [-19, +8]; \quad I_g(X) := \text{sqr}(X \diamond 2) \diamond 8 = [-8, -4] \subset I_f(X),$$

wobei  $I_g(X)$  sogar die optimale Einschließung des Wertebereichs von  $f(x) \equiv g(x)$  über  $X$  ist, da  $g(x)$  die Variable  $x$  nur einmal enthält. Vergleichen Sie auch den Abschnitt B.5 auf Seite 226.

### D.1.3. Optimale Einschließungen

Mit  $f(x) = 3x \cdot e^{-2x+1}$  haben wir schon ein Beispiel kennengelernt, bei dem die intervallmäßige Auswertung über  $X = [0, 2]$  den Wertebereich  $W_{f,X} = [0, 3/2]$  erheblich überschätzt. Bei einmal stetig differenzierbaren Funktionen  $f(x)$  kann man diese Überschätzungen jedoch vermeiden, wenn man das Intervall  $X = \bigcup_i X_i$  in Maschinenintervalle  $X_i$  unterteilt, in denen  $f(x)$  jeweils monoton ist. Bedeutet dann  $I_f(X_i)$  die Einschließung aller Funktionswerte  $f(x)$ , mit  $x \in X_i$ , so

<sup>1</sup>Der Satz D.1 gilt sinngemäß auch dann, wenn neben  $x$  auch noch weitere Variablen, z.B.  $y, z, \dots$ , auftreten.

ist  $I_{f,X} := \cup_i I_f(X_i)$  eine deutlich bessere Einschließung von  $W_{f,X} \subseteq I_{f,X}$ . Mit dem Teilintervall  $X_i = [x_{i,1}, x_{i,2}]$  ist dann  $I_f(X_i)$  wie folgt zu berechnen:

$$(D.3) \quad I_f(X_i) = [\text{Inf}(f([x_{i,1}, x_{i,1}])), \text{Sup}(f([x_{i,2}, x_{i,2}]))].$$

Die Funktion  $f(x)$  ist also jeweils intervallmäßig auf den Intervallrändern von  $X_i = [x_{i,1}, x_{i,2}]$  auszuwerten. Um bei  $f(x) = 3x \cdot e^{-2x+1}$  die Monotonieintervalle  $X_i$  zu bestimmen, benötigen wir die Nullstellen von  $f'(x) = 3e^{-2x+1}(1-2x)$  und erhalten damit  $X_1 = [0, 0.5]$ ,  $X_2 = [0.5, 2]$  und mit diesen Intervallen nach (D.3) die Einschließungen  $I_f(X_1) = [0, 1.5]$ ,  $I_f(X_2) = [6 \cdot e^{-3}, 1.5]$ , d.h.  $I_{f,X} = I_f(X_1) \cup I_f(X_2) = [0, 3/2] = W_{f,X}$ . Wesentlich komplizierter wird das Verfahren jedoch, wenn die Nullstellen von  $f'(x)$  **keine** Maschinenzahlen sind. Ist dann z.B.  $Z_1 = [\zeta_1, \zeta_2]$  eine Maschineneinschließung einer solchen Nullstelle  $z_1$ , so muss  $f(x)$  wieder intervallmäßig über diesem Intervall  $Z$  ausgewertet werden, wodurch die schon bekannten Überschätzungen wieder auftreten können, und diese Überschätzungen machen sich besonders dann negativ bemerkbar, wenn die Nullstelle  $z_1$  von  $f'(x)$  entweder ganz in der Nähe einer Nullstelle von  $f(x)$  liegt oder sogar mit dieser Nullstelle übereinstimmt, wie es z.B. schon bei  $g(x) := (x - \sqrt{33})^2$  realisiert ist. Besonders aufwendig wird das Verfahren auch dann, wenn  $f'(x)$  im Intervall  $X$  sehr viele Nullstellen besitzt, die natürlich alle möglichst eng eingeschlossen werden müssen. Das beschriebene Verfahren liefert also nur bei verhältnismäßig einfachen Funktionen eine wirklich optimale Einschließung ihres Wertebereichs.

Mit Hilfe der *Globalen Optimierung* kann der Wertebereich einer zweimal stetig differenzierbaren Funktion  $f(x)$  ebenfalls garantiert eingeschlossen werden, [27], [28], [29], [31], [59]. Dabei wird eine Unter- bzw. Oberschranke des Wertebereichs  $W_{f,X}$  einer Funktion  $f(x)$  über einem vorgegebenen Intervall  $X$  berechnet. Vergleichen Sie dazu den Abschnitt 7.4 auf Seite 131.

Wir stellen uns jetzt dem Problem, wie man ausgehend von einer Funktion  $f : D_f \rightarrow \mathbb{R}$ , deren Wertebereich durch intervallmäßige Auswertung optimal eingeschlossen wird, neue Funktionen  $g : D_g \rightarrow \mathbb{R}$  konstruieren kann, deren Wertebereiche  $W_{g,X}$  über dem Intervall  $X$  dann ebenfalls optimal eingeschlossen werden. Wir formulieren dazu den recht einfachen

**Satz D.2** (Multiplikation mit einer reellen Konstanten  $a \in \mathbb{R}$ ).

Liefert  $f : D_f \rightarrow \mathbb{R}$  bei intervallmäßiger Auswertung über dem Intervall  $X$  eine optimale Einschließung ihres Wertebereichs  $W_{f,X}$ , so liefert auch die Funktion  $g(x) := a \cdot f(x)$  eine nahezu optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

*Beweis.* Nach Voraussetzung gibt es nach Seite 373 die optimale Einschließung  $I(X)$  des Wertebereichs  $W_{f,X}$ , d.h. es gilt  $W_{f,X} \subseteq I(X)$  mit den Ungleichungen aus (D.2). Mit  $g(x) := a \cdot f(x)$  gilt zusätzlich  $W_{g,X} = a \cdot W_{f,X}$  und damit  $a \cdot W_{f,X} \subseteq a \diamond I(X)$ , wobei  $a \diamond I(X)$  das Ergebnis der intervallmäßigen Auswertung von  $g(x)$  über dem Intervall  $X$  ist. Da die Multiplikation von  $a$  mit dem Intervall  $I(X)$  auf der Maschine optimal implementiert ist, wird  $W_{f,X}$  durch  $a \diamond I(X)$  optimal oder nahezu optimal eingeschlossen, wobei *nahezu optimal* bedeutet, dass bei der Maschinenmultiplikation  $a \diamond I(X)$  nach außen gerundet wird, wobei im ungünstigsten Fall nach rechts und links je ein Bit Genauigkeit verloren geht ■

Neben der Multiplikation kann man eine neue Funktion  $g(x)$  auch durch die Addition einer reellen Zahl  $b$  gewinnen, dies beschreibt der

**Satz D.3** (Addition einer reellen Konstanten  $b \in \mathbb{R}$ ).

Liefert  $f : D_f \rightarrow \mathbb{R}$  bei intervallmäßiger Auswertung über dem Intervall  $X$  eine optimale Einschließung ihres Wertebereichs  $W_{f,X}$ , so liefert auch die Funktion  $g(x) := f(x) + b$  eine nahezu optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

*Beweis.* Der einfache Beweis bleibt dem Leser überlassen.

Die beiden letzten Sätze lassen sich zusammenfassen zum

**Satz D.4** (Multiplikation oder Addition einer reellen Konstanten).

Liefert  $f : D_f \rightarrow \mathbb{R}$  bei intervallmäßiger Auswertung über dem Intervall  $X$  eine optimale Einschließung ihres Wertebereichs  $W_{f,X}$ , so liefert auch die Funktion  $g(x) := a \cdot f(x) + b$  eine nahezu optimale Einschließung ihres Wertebereichs  $W_{g,X}$ .

Als Beispiel betrachten wir die Exponentialfunktion  $f(x) = e^x$ , deren Wertebereich mit Hilfe der Intervallfunktion  $\exp(X)$  über dem Maschinenintervall  $X$  optimal eingeschlossen wird. Nach Satz D.4 wird dann auch der Wertebereich  $W_{g,X}$  der Funktion  $g(x) := \sqrt{2} \cdot e^x + \sqrt{3}$  mit Hilfe des Intervallausdrucks

$$\text{sqrt}(2) \diamond \exp(X) \oplus \text{sqrt}(3) \supseteq W_{g,X}$$

nahezu optimal eingeschlossen, wenn mit  $\text{sqrt}(2)$  und  $\text{sqrt}(3)$  die optimalen Einschließungen von  $\sqrt{2}$  bzw.  $\sqrt{3}$  bezeichnet werden, wobei  $\text{sqrt}(2)$  und  $\text{sqrt}(3)$  natürlich keine Punktintervalle sein können. Deshalb werden außer durch die Intervall-Multiplikation und Addition auch noch zusätzliche aber nur minimale Überschätzungen durch die Intervalle  $\text{sqrt}(2)$  und  $\text{sqrt}(3)$  verursacht, so dass die Einschließung von  $W_{g,X}$  nur *nahezu optimal* sein kann.

Ausgehend von einer Funktion  $f : D_f \rightarrow \mathbb{R}$ , die bei intervallmäßiger Auswertung über einem Maschinenintervall  $T \subset D_f$  eine optimale Einschließung ihres Wertebereichs  $W_{f,T}$  liefert, kann man mit einer Transformationsfunktion  $t : D_t \rightarrow \mathbb{R}$  neue Funktionen  $g(x) := f(t(x))$  gewinnen, wenn der Wertebereich  $W_{t,X}$  bez. des Maschinenintervalls  $X \subset D_t$  selbst wieder ein Intervall  $T$  ist, so dass gilt  $W_{t,X} = T \subset D_f$ . Diese neue Funktion  $g(x) := f(t(x))$  liefert dann bei intervallmäßiger Auswertung über dem Intervall  $X$  eine nahezu optimale Einschließung ihres Wertebereichs  $W_{g,X}$ , wenn auch die intervallmäßige Auswertung von  $t(x)$  über dem Intervall  $X$  den Wertebereich  $W_{t,X}$  optimal einschließt. Wir formulieren dies im

**Satz D.5** (Transformationsfunktion).

Die reelle Funktion  $f : D_f \rightarrow \mathbb{R}$  liefert nach Voraussetzung bei intervallmäßiger Auswertung über dem Intervall  $T \subset D_f$  eine optimale Einschließung ihres Wertebereichs  $W_{f,T}$  und  $t : D_t \rightarrow \mathbb{R}$  besitzt einen Wertebereich  $W_{t,X}$  über dem Intervall  $X \subset D_t$ , der selbst wieder ein Intervall  $T := W_{t,X} \subset D_f$  ist. Dann wird bei intervallmäßiger Auswertung der neuen Funktion  $g(x) := f(t(x))$  ihr Wertebereich  $W_{g,X}$  auch wieder nahezu optimal eingeschlossen, wenn bei intervallmäßiger Auswertung auch der Wertebereich  $T := W_{t,X}$  optimal oder nahezu optimal eingeschlossen wird.

*Beweis.* Nach Voraussetzung ist  $W_{t,X}$  selbst ein **Intervall**, das durch das Intervall  $T \subset D_f$  optimal oder nahezu optimal eingeschlossen wird und nach Voraussetzung liefert auch die intervallmäßige Auswertung von  $f(x)$  über dem Intervall  $T$  selbst wieder eine optimale Einschließung ihres Wertebereichs  $W_{f,T}$ , so dass  $W_{g,X}$  durch die intervallmäßige Auswertung von  $g(x) = f(t(x))$  nahezu optimal ausgewertet wird.

Als Beispiel wählen wir  $f(x) = e^x$  und als Transformationsfunktion  $t(x) = \sin(x) + 2$  mit dem beliebigen Maschinenintervall  $X$ . Damit sind die Voraussetzungen des Satzes D.5 erfüllt, und mit der intervallmäßigen Auswertung von  $g(x) = e^{\sin(x)+2}$  erhalten wir über jedem Intervall  $X$  eine nahezu optimale Einschließung des Wertebereichs  $W_{g,X}$ .

Beachten Sie, dass man mit den Intervallfunktionen aus Tabelle 4.14.3 auf Seite 68 eine Vielzahl von Transformationsfunktion  $t(x)$  konstruieren kann. Beachten Sie aber auch, dass z.B.  $t(x) = \sin(x) + x$  die Voraussetzungen von Satz D.5 nicht erfüllt, da  $t(x)$  die Variable  $x$  zweimal enthält, so dass  $W_{t,X}$  jetzt **keinesfalls** optimal eingeschlossen werden kann.



## D.2. Einschließung komplexer arithmetischer Ausdrücke

### D.2.1. Problemstellung und Definitionen

In einem *komplexen arithmetischen Ausdruck* sind per Definition konstante komplexe Zahlen und Funktionen aus der Tabelle 6.12.3 mit den vier arithmetischen Grundoperationen verknüpft. Die komplexe Intervall-Variable  $Z = X + i \cdot Y$  ist in der komplexen Ebene ein Rechteck, wobei  $X = [x_1, x_2]$  und  $Y = [y_1, y_2]$  reelle Intervalle sind, vgl. Seite 243.

Ist  $f : \mathbb{C} \rightarrow \mathbb{C}$  eine komplexwertige Funktion, mit z.B.  $f(z) = \sin(z)$ ,  $z = x + i \cdot y \in Z \subset \mathbb{C}$ , so ist

$$(D.4) \quad W_{f,Z} := \{f(z) \mid z \in Z \subset \mathbb{C}\}$$

die Wertemenge  $W_{f,Z}$  von  $f$  über dem komplexen Intervall  $Z$ .

Bei der komplexen Intervallrechnung ist es nun das Ziel, zu einer Funktion  $f(z)$  und zu einem vorgegebenen Intervall  $Z$  mit einem geeigneten Algorithmus die Wertemenge  $W_{f,Z}$  durch ein Rechteck  $F(Z)$  wenigstens **optimal** einzuschließen, d.h. es muss gelten  $W_{f,Z} \subseteq F(Z)$ , und jede Seite von  $F(Z)$  muss mit wenigstens einem Punkt der Menge  $W_{f,Z}$  übereinstimmen, vgl. Seite 243. Wir **definieren**:

Ist  $f : \mathbb{C} \rightarrow \mathbb{C}$ , mit  $f(z) = u(x, y) + i \cdot v(x, y)$ , über dem Rechteck  $Z = X + i \cdot Y$  eine holomorphe Funktion, dann ist  $F(Z)$  eine optimale Einschließung des Wertebereichs  $W_{f,Z}$   $\iff$

$$\begin{aligned} \exists x_{u,m} \in X, \exists y_{u,m} \in Y, u_m := u(x_{u,m}, y_{u,m}); \quad & \exists x_{u,M} \in X, \exists y_{u,M} \in Y, u_M := u(x_{u,M}, y_{u,M}); \\ \exists x_{v,m} \in X, \exists y_{v,m} \in Y, v_m := v(x_{v,m}, y_{v,m}); \quad & \exists x_{v,M} \in X, \exists y_{v,M} \in Y, v_M := v(x_{v,M}, y_{v,M}); \\ \implies \quad \forall x \in X \wedge \forall y \in Y \quad \text{gilt:} \quad & u_m \leq u(x, y) \leq u_M, \quad v_m \leq v(x, y) \leq v_M. \end{aligned}$$

**Anmerkung:** Die Punkte  $(x_{u,m}, y_{u,m})$  und  $(x_{u,M}, y_{u,M})$  sind die Extrempunkte  $m, M$  von  $u(x, y)$  auf dem Rand von  $Z$ . Entsprechendes gilt für die Punkte  $(x_{v,m}, y_{v,m})$  und  $(x_{v,M}, y_{v,M})$ .

Mit Ausnahme der Funktionen `power_fast` und `pow_all` liefern alle Funktionen aus Tabelle 6.12.3 eine **optimale** Einschließung  $F(Z) \supseteq W_{f,Z}$ . Gilt sogar  $W_{f,Z} = F(Z)$ , so nennt man die Einschließung  $F(Z)$  **exakt**, wobei die exakten Einschließungen nur bei einfachen Funktionen, wie z.B.  $f(z) = a \cdot z + c$ ,  $a \in \mathbb{R}$ ,  $c \in \mathbb{C}$  auftreten.

Wertet man einen arithmetischen Ausdruck, z.B.  $f(z) = z/(1+z)$ , intervallmäßig aus, indem man die Variable  $z$  durch ein vorgegebenes Intervall  $Z$  ersetzt, so erhält man mit dem Ergebnis

$$(D.5) \quad f(Z) := Z \diamond (1 \oplus Z)$$

eine Einschließung von  $W_{f,Z}$ , d.h. es gilt:  $W_{f,Z} \subset f(Z)$ , aber es ist keinesfalls sicher, dass  $f(Z)$  auch eine optimale Einschließung von  $W_{f,Z}$  ist. Im folgenden Abschnitt soll untersucht werden, wie z.B. der Ausdruck  $f(z) = z/(1+z)$  umgeformt werden kann, um die Einschließung weiter zu verbessern.

Fast alle Funktionen aus Tabelle 6.12.3 liefern bei intervallmäßiger Auswertung eine **optimale** Einschließung  $F(Z)$  der Wertemenge  $W_{f,Z} := \{f(z) \mid z \in Z \subset \mathbb{C}\}$ . Für die praktische Anwendung wäre es interessant zu wissen, mit welchen Transformationen, z.B. mit  $z \rightarrow 1/z$ , man ausgehend von  $f(z)$  eine neue Funktion  $g(z) := f(1/z)$  erhält, die ebenfalls mit  $g(Z)$  eine optimale Einschließung des neuen Wertebereichs  $W_{g,Z}$  liefert. Als Beispiel betrachten wir  $f(z) = z$ , und mit  $z \rightarrow 1/z$  erhält man  $g(z) = 1/z$ , so dass sowohl  $f(Z) = Z$  und  $g(Z) = 1 \diamond Z = \text{reci}(Z)$  jeweils optimale Einschließungen liefern. Bei  $f(Z) = Z$  ist dies wegen der identischen Abbildung klar, und bei  $g(Z) = \text{reci}(Z)$  liefert `reci(Z)` nach Tabelle 6.12.3 automatisch eine optimale Einschließung, falls  $0 \notin Z$ . Beachten Sie, dass auch  $g(Z) = 1 \diamond Z$  eine optimale Einschließung liefert, da die direkte Intervalldivision als eine der vier Grundoperationen eine optimale Einschließung berechnet, vergleichen Sie dazu die Anmerkungen auf Seite 287.

Im nächsten Beispiel betrachten wir  $f(z) = z^2$ , und mit der gleichen Transformation  $z \rightarrow 1/z$  erhalten wir  $g(z) = f(1/z) = 1/z^2$ . Mit  $f(Z) = \text{sqr}(Z)$  erhält man eine optimale Einschließung,

aber  $g(Z) = 1 \diamond \text{sqr}(Z)$  oder  $g(Z) = \text{sqr}(1 \diamond Z)$  liefern nach Seite 288 i.a. **keine** optimale Einschließung. Erst mit  $G(Z) := \text{reci\_z2}(Z)$  wird der Wertebereich  $W_{g,Z}$  optimal eingeschlossen. Die Transformation  $z \rightarrow 1/z$  liefert also i.a. **keine** optimale Einschließung!

Im Abschnitt D.2.3 werden die Bedingungen für optimale Einschließungen genauer untersucht.

## D.2.2. Verbesserung einer Einschließung

In diesem Abschnitt wird untersucht, unter welchen Gesichtspunkten die durch intervallmäßige Auswertung gewonnene garantierte Einschließung  $f(Z)$  eines gegebenen arithmetischen Ausdrucks  $f(z)$ ,  $z \in Z \subset \mathbb{C}$ , noch weiter verbessert werden kann. In nachfolgender Tabelle sind dazu für  $Z = [2, 5] + i \cdot [1, 3]$  die durch intervallmäßige Auswertung berechneten Einschließungen  $g(Z) \subset f(Z)$  der jeweils identischen Ausdrücke  $g(z) \equiv f(z)$ ,  $z \in Z$ , zusammengestellt.

$Z = X + i \cdot Y = [2, 5] + i \cdot [1, 3];$ Präzision: <code>prec = 15.</code>	
$f(z) = z/(1+z)$	$f(Z) = ([3.333282e-1, 1.800049], [-6.666871e-1, 7.000123e-1])$
$g(z) = 1/(1+1/z)$	$g(Z) = ([6.921691e-1, 8.708802e-1], [1.960563e-2, 1.814042e-1])$
$f(z) = (z+1)/z$	$f(Z) = ([5.293884e-1, 3.000000], [-1.270813, 6.213379e-1])$
$g(z) = 1+1/z$	$g(Z) = ([1.147033, 1.400025], [-2.500000e-1, -3.845977e-2])$
$f(z) = 1/(1/z)$	$f(Z) = ([1.747924, 6.364991], [2.381515e-1, 3.400147])$
$g(z) = z$	$g(Z) = ([2.000000, 5.000000], [1.000000, 3.000000])$
$f(z) = \frac{z^2-1}{z+1}$	$f(Z) = ([-1.400025, 1.013477e1], [-3.166749, 9.600098])$
$g(z) = z-1$	$g(Z) = ([1.000000, 4.000000], [1.000000, 3.000000])$

Die oberen vier Beispiele zeigen wegen  $g(Z) \subset f(Z)$ , dass ein komplexer, arithmetischer Ausdruck so umgeformt werden sollte, dass die Variable  $z$  möglichst nur einmal vorkommt und die Zahl der arithmetische Operatoren minimal ist.

$Z = X + i \cdot Y = [2, 5] + i \cdot [1, 3];$ $a = [-2, -2] + i \cdot [2, 2],$ $b = [2, 2] - i \cdot [1, 1];$	
$f(z) = \text{sqr}(z) \cdot z + (\text{sqr}(z) + az + b)$	$f(Z) = ([-2.500e1, 1.100e1], [-2.700e1, 9.000])$
$g(z) = \text{sqr}(z) \cdot z + \text{poly2}(z, a, b)$	$g(Z) = ([-2.100e1, 7.000], [-2.300e1, 5.000])$

Das letzte Beispiel zeigt, dass ein Teilausdruck, der nach Tabelle 6.12.3 optimal eingeschlossen wird, durch die jeweilige Funktion dieser Tabelle, hier `poly2(z, a, b)`, ersetzt werden sollte. Wir fassen zusammen:

Zur Verbesserung der Einschließung eines komplexen, arithmetischen Ausdrucks sollte dieser so äquivalent umgeformt werden, dass die Variable  $z$  möglichst nur einmal vorkommt und die Zahl der arithmetischen Operatoren minimal ist. Falls möglich sollten Teilausdrücke durch ihre gleichwertigen Funktionen aus Tabelle 6.12.3 mit optimaler Einschließung ersetzt werden. Die verbesserte Einschließung ist dann aber nicht notwendig optimal!

## D.2.3. Optimale Einschließungen

In diesem Abschnitt betrachten wir die in einem Gebiet  $G \subset \mathbb{C}$  holomorphen Funktionen  $f(z)$  der Tabelle 6.12.3, die mit Ausnahme der beiden Funktionen `power_fast` und `pow_all` zu einem gegebenen Intervallargument  $Z = X + i \cdot Y \subset G$  bei intervallmäßiger Auswertung mit  $f(Z)$  eine

optimale Einschließung des Wertebereichs  $W_{f,Z} := \{f(z) \mid z \in Z \subset G\}$  liefern und stellen jetzt die Frage, mit welchen Manipulationen man, ausgehend von  $f(z)$ , eine neue holomorphe Funktion  $g(z)$  konstruieren kann, die bei intervallmäßiger Auswertung mit  $g(Z)$  ebenfalls eine optimale Einschließung ihres Wertebereichs  $W_{g,Z} := \{g(z) \mid z \in Z \subset G\}$  liefert.

Bevor wir jedoch zur Beantwortung dieser Frage kommen, soll zunächst noch geklärt werden, unter welchen Voraussetzungen garantiert ist, dass eine gegebene holomorphe Funktion  $f(z) = u(x, y) + i \cdot v(x, y)$  mit  $f(Z)$  eine optimale Einschließung von  $W_{f,Z}$  liefert. Mit  $z = x + i \cdot y \in Z \subset G$  gilt:

1. Wenn in den Funktionstermen<sup>2</sup> für  $u(x, y)$  und  $v(x, y)$  die Variablen  $x$  und  $y$  jeweils nur einmal vorkommen, so liefert  $f(Z)$  eine optimale Einschließung von  $W_{f,Z}$ , [6, Seite 32]. Ein Beispiel ist die Exponentialfunktion

$$e^z = e^x \cdot \cos(y) + i \cdot e^x \cdot \sin(y), \quad x \in X, \quad y \in Y.$$

2. Wenn man in der  $x, y$ -Ebene auf dem Rand eines beliebigen Rechtecks  $Z \subset G$  für  $u(x, y)$  und  $v(x, y)$  jeweils die Maximum- und Minimum-Punkte  $M, m$  bestimmt hat, so kann man mithilfe der Funktionswerte von  $u(x, y)$  und  $v(x, y)$  an diesen Punkten für  $f(z)$  eine optimale Einschließung von  $W_{f,Z}$  berechnen, vgl. für  $f(z) = 1/z^2$  die Seiten 210 und 216.

Wir kommen jetzt zur Beantwortung der schon angesprochenen Frage, wie man eine holomorphe Funktion  $f(z)$  so manipulieren kann, dass die neue, ebenfalls holomorphe Funktion,  $g(z)$  bei intervallmäßiger Auswertung auch eine optimale Einschließung von  $W_{g,Z}$  liefert. Wir formulieren dazu den ersten, recht einfachen

**Satz D.6** (Addition einer Konstanten  $c = c_1 + i \cdot c_2 \in \mathbb{C}$ ).

Ist  $f(z)$  holomorph in  $Z = X + i \cdot Y \subset \mathbb{C}$  und liefert die intervallmäßige Auswertung von  $f$  eine optimale Einschließung des Wertebereichs  $W_{f,Z}$ , so liefert auch  $g(z) := f(z) + c$  eine optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

*Beweis.* Mit  $f(z) = u(x, y) + i \cdot v(x, y)$  gilt:  $g(z) = (u(x, y) + c_1) + i \cdot (v(x, y) + c_2)$ , und nach Voraussetzung gibt es auf dem Rand von  $Z \subset \mathbb{C}$  für  $u(x, y)$  und  $v(x, y)$  die entsprechenden Extrempunkte  $m, M$ , so dass  $\forall x \in X \wedge \forall y \in Y$  gilt:

$$u_m \leq u(x, y) \leq u_M, \quad v_m \leq v(x, y) \leq v_M, \quad \text{und daraus folgt für die gleichen Extrempunkte} \\ u_m + c_1 \leq u(x, y) + c_1 \leq u_M + c_1, \quad v_m + c_2 \leq v(x, y) + c_2 \leq v_M + c_2,$$

wobei zu beachten ist, dass in den obigen Ungleichungen das Gleichheitszeichen gilt, wenn die jeweiligen Extrempunkte  $m, M$  auf dem Rand von  $Z$  gewählt werden, und nach Definition sind die beiden letzten Doppelungleichungen gerade die Bedingung für eine optimale Einschließung des Wertebereichs  $W_{g,Z}$  ■

Als **Anwendung** betrachten wir mit  $c = 1 + i$  die Funktionen  $f(z) = 1/z$  und  $g(z) = 1/z + c$ , die nach Satz D.6 zu einem beliebigen Rechteck  $Z \not\ni 0$  bei intervallmäßiger Auswertung jeweils **optimale** Einschließungen ihrer verschiedenen Wertebereiche  $W_{f,Z}$  bzw.  $W_{g,Z}$  liefern. Beachten Sie jedoch, dass die äquivalente Funktion  $h(z) = (1 + c \cdot z)/z \equiv g(z) = 1/z + c$  **keine optimale** Einschließung von  $W_{h,Z} = W_{g,Z}$  mehr liefert!

**Satz D.7** (Multiplikation mit einer reellen Konstanten  $r \in \mathbb{R}$ ).

Ist  $f(z)$  holomorph in  $Z = X + i \cdot Y \subset \mathbb{C}$  und liefert die intervallmäßige Auswertung von  $f$  eine optimale Einschließung des Wertebereichs  $W_{f,Z}$ , so liefert auch  $g(z) := r \cdot f(z)$  eine optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

<sup>2</sup>Die Funktionsterme für  $u(x, y)$  und  $v(x, y)$  können z.B. mit *Mathematica* und `ComplexExpand[]` berechnet werden.

*Beweis.* Mit  $f(z) = u(x, y) + i \cdot v(x, y)$  gilt  $g(z) = r \cdot u(x, y) + i \cdot r \cdot v(x, y)$ . Nach Voraussetzung gibt es auf dem Rand von  $Z \subset \mathbb{C}$  für  $u(x, y)$  und  $v(x, y)$  die entsprechenden Extrempunkte  $m, M$ , so dass  $\forall x \in X \wedge \forall y \in Y$  gilt:

$$u_m \leq u(x, y) \leq u_M, \quad v_m \leq v(x, y) \leq v_M, \quad \text{und daraus folgt für die gleichen Extrempunkte} \\ r \cdot u_m \leq r \cdot u(x, y) \leq r \cdot u_M, \quad r \cdot v_m \leq r \cdot v(x, y) \leq r \cdot v_M, \quad \text{falls } r \geq 0,$$

wobei zu beachten ist, dass in den obigen Ungleichungen das Gleichheitszeichen gilt, wenn für  $x, y$  die Koordinaten der jeweiligen Extrempunkte  $m, M$  auf dem Rand von  $Z$  gewählt werden, und nach Definition sind die beiden letzten Doppelungleichungen gerade die Bedingung für eine optimale Einschließung des Wertebereichs  $W_{g,Z}$ . Im Fall  $r < 0$  erhält man ebenfalls mit den gleichen Extrempunkten die beiden Doppelungleichungen

$$r \cdot u_M \leq r \cdot u(x, y) \leq r \cdot u_m, \quad r \cdot v_M \leq r \cdot v(x, y) \leq r \cdot v_m, \quad \text{falls } r \geq 0 \quad \blacksquare$$

**Satz D.8** (Multiplikation mit der imaginären Einheit  $i = \sqrt{-1}$ ).

Ist  $f(z)$  holomorph in  $Z = X + i \cdot Y \subset \mathbb{C}$  und liefert die intervallmäßige Auswertung von  $f$  eine optimale Einschließung des Wertebereichs  $W_{f,Z}$ , so liefert auch  $g(z) := i \cdot f(z)$  eine optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

*Beweis.* Mit  $f(z) = u(x, y) + i \cdot v(x, y)$  gilt:  $g(z) = -v(x, y) + i \cdot u(x, y)$ , und nach Voraussetzung gibt es auf dem Rand von  $Z \subset \mathbb{C}$  für  $u(x, y)$  und  $v(x, y)$  die entsprechenden Extrempunkte  $m, M$ , so dass  $\forall x \in X \wedge \forall y \in Y$  gilt:

$$u_m \leq u(x, y) \leq u_M, \quad v_m \leq v(x, y) \leq v_M, \quad \text{und daraus folgt für die gleichen Extrempunkte} \\ u_m \leq u(x, y) \leq u_M, \quad -v_M \leq -v(x, y) \leq -v_m, \quad \text{falls } r \geq 0,$$

wobei zu beachten ist, dass in den obigen Ungleichungen das Gleichheitszeichen gilt, wenn die jeweiligen Extrempunkte  $m, M$  auf dem Rand von  $Z$  gewählt werden, und nach Definition sind die beiden letzten Doppelungleichungen gerade die Bedingung für eine optimale Einschließung des Wertebereichs  $W_{g,Z}$  ■.

Man kann die ersten drei Sätze zusammenfassen:

**Satz D.9** (Multiplikation mit  $r \in \mathbb{R}$  oder mit  $i \cdot r$  und anschließender Translation mit  $c \in \mathbb{C}$ ).

Ist  $f(z)$  holomorph in  $Z = X + i \cdot Y \subset \mathbb{C}$  und liefert die intervallmäßige Auswertung von  $f(z)$  eine optimale Einschließung des Wertebereichs  $W_{f,Z}$ , so liefern sowohl  $g(z) := r \cdot f(z) + c$  als auch  $h(z) := i \cdot r \cdot f(z) + c$  jeweils eine optimale Einschließung ihrer Wertebereiche  $W_{g,Z}$  bzw.  $W_{h,Z}$ .

Mit  $a \in \mathbb{C}$  und  $b \in \mathbb{IC}$  liefern daher nach Satz D.9  $f(z) = z^2 + az + b$ ,  $g(z) = r \cdot (z^2 + az + b)$  und  $h(z) = ir \cdot (z^2 + az + b)$  bei intervallmäßiger Auswertung jeweils eine **optimale** Einschließung ihrer Wertebereiche, vgl. die Tabelle auf Seite 115. Beachten Sie, dass die beiden Klammern () jeweils mit der Funktion `poly2(z, a, b)` auszuwerten sind und dass z.B.  $r \cdot (z^2 + az + b)$  nicht ausmultipliziert werden darf. Ist beispielsweise der Ausdruck  $rz^2 + \alpha z + \beta$  gegeben, so erhält man eine optimale Einschließung über einem beliebigen Rechteck  $Z \subset \mathbb{C}$  nur dann, wenn man für  $r \neq 0$  den äquivalenten Term  $r \cdot (z^2 + \alpha/r + \beta/r)$  intervallmäßig auswertet,  $\alpha \in \mathbb{C}$ ,  $\beta \in \mathbb{IC}$ .

Bisher hatten wir durch Addition oder durch geeignete Multiplikation aus einer holomorphen Funktion  $f(z)$ , die bei intervallmäßiger Auswertung eine optimale Einschließung liefert, eine neue holomorphe Funktion  $g(z)$  mit der gleichen Eigenschaft gewonnen.

Jetzt suchen wir zu einer holomorphen Funktion  $f : D_f \rightarrow \mathbb{C}$  eine holomorphe Transformationsfunktion  $t(z)$ , mit der man eine neue Funktion  $g(z) = f(t(z))$  gewinnt, die ebenfalls bei intervallmäßiger Auswertung eine optimale Einschließung liefert. Eine solche Funktion  $t(z)$  müsste die Eigenschaft haben, ein achsenparalleles Rechteck  $Z \subset D_f$  wieder in ein achsenparalleles Rechteck  $Z_t \subset D_f$  abzubilden. Mit dem neuen Rechteck  $Z_t$  liefert dann  $f(z)$  nach Voraussetzung bei intervallmäßiger Auswertung eine optimale Einschließung des Wertebereichs  $W_{f,Z_t}$ .

Eine geeignetes  $t(z)$  ist die lineare Funktion

$$(D.6) \quad t(z) = \alpha \cdot z + c, \quad \text{mit } \alpha = r \in \mathbb{R} \text{ oder } \alpha = r \cdot i \text{ und } c \in \mathbb{C},$$

und die neue Funktion ist gegeben durch  $g(z) = f(t(z))$ . Daraus ergibt sich der

### Satz D.10.

Die Funktion  $f : D_f \rightarrow \mathbb{C}$  sei holomorph in  $D_f$  und liefere bei intervallmäßige Auswertung über einem beliebigen Rechteck  $Z \subset D_f$  eine optimale Einschließung ihres Wertebereichs  $W_{f,Z}$ . Ist dann  $t : D_f \rightarrow D_f \subset \mathbb{C}$  definiert durch  $t(z) = \alpha \cdot z + c$ , mit  $\alpha = r \in \mathbb{R}$  oder  $\alpha = r \cdot i$  und  $c \in \mathbb{C}$ , so liefert auch  $g(z) := f(t(z))$  bei intervallmäßiger Auswertung eine optimale Einschließung ihres Wertebereichs  $W_{g,Z}$ .

### Anmerkungen:

1.  $t(z)$  ist entweder eine Streckung mit anschließender Translation oder eine Drehstreckung um  $\pm\pi$  mit anschließender Translation, so dass ein achsenparalleles Rechteck wieder in ein solches abgebildet wird. Die intervallmäßige Auswertung von  $t(z)$  liefert daher mit  $t(Z)$  nicht nur eine optimale sondern auch eine **exakte** Einschließung von  $W_{t,Z}$ .  
Mit  $t(z) = 2iz + 1 + i$  und  $f(z) = \sin(z)$  erhält man mit  $g(z) = \sin(2iz + 1 + i)$  eine in ganz  $\mathbb{C}$  holomorphe Funktion, die über einem beliebigen Rechteck  $Z \subset \mathbb{C}$  bei intervallmäßiger Auswertung mit  $g(Z)$  eine optimale Einschließung des Wertebereichs  $W_{g,Z}$  liefert.
2. Definiert man  $t(z) = \alpha \cdot z + c$  jetzt mit beliebigem  $\alpha \in \mathbb{C}$ , so ist die Transformation  $t(z)$  wieder eine Drehstreckung mit anschließender Translation, aber der Wertebereich  $W_{t,Z}$  von  $t(z)$  ist jetzt ein beliebig gedrehtes und damit i.a. kein achsenparalleles Rechteck, so dass die intervallmäßige Auswertung  $t(Z)$  zwar eine optimale aber **keine exakte** Einschließung von  $W_{t,Z}$  liefert. Die intervallmäßige Auswertung von  $g(z) := f(t(z))$  ergibt damit i.a. **keine optimale** Einschließung von  $W_{g,Z}$ .
3. Spiegelt man beispielsweise ein achsenparalleles Rechteck  $Z = X + i \cdot Y$  an der 1. Winkelhalbierenden (Vertauschung von  $X, Y$ ), so erhält man wieder ein achsenparalleles Rechteck  $Z_S$ , und wenn dies wieder in  $D_f$  liegt, so ist die Einschließung von  $W_{f,Z_S}$  auch über diesem Rechteck  $Z_S$  wieder optimal. Da aber  $t(z)$  durch das Vertauschen von  $x$  und  $y$  keine holomorphe Funktion mehr ist, so ist auch die neue Funktion  $g(z) := f(t(z))$  nicht mehr holomorph, und die Extrema von  $g(z)$  müssen nicht mehr nur auf dem Rand von  $Z$  liegen, so dass Satz D.10 nicht mehr anwendbar ist.

Eine interessante **Anwendung** zur optimalen Einschließung ist die **Möbiustransformation**:

$$(D.7) \quad T(z) := \frac{az + b}{cz + d}; \quad a, b, c, d, z \in \mathbb{C}, \quad ad - bc \neq 0.$$

$$(D.8) \quad T(z) = \begin{cases} \frac{a}{d} \cdot z + \frac{b}{d}, & c = 0 \\ \frac{az + b}{z + d}, & c \neq 0, \quad a, b, d, z \in \mathbb{C}, \quad ad - b \neq 0. \end{cases}$$

Den letzten Ausdruck in (D.8) findet man im Fall  $c \neq 0$ , wenn man in (D.7) Zähler und Nenner durch  $c$  dividiert und dann z.B.  $a/c$  wieder mit  $a$  bezeichnet. Der letzte Ausdruck in (D.8) hat bei intervallmäßiger Auswertung über  $Z$  für eine optimale Einschließung des Wertebereichs  $W_{T,Z}$  jedoch den Nachteil, dass  $z$  im Zähler **und** im Nenner auftaucht. Dies wird jedoch vermieden, wenn man mithilfe der *Polynomdivision* wie folgt umformt

$$(D.9) \quad T(z) \equiv T_1(z) := \frac{b - ad}{z + d} + a, \quad a, b, d, z \in \mathbb{C}, \quad ad - b \neq 0.$$

Ist nun  $Z$  in der komplexen Ebene ein achsenparalleles Rechteck, welches den Punkt  $z_0 = -d$  nicht enthält, so erhält man bei intervallmäßiger Auswertung des Bruchs in (D.9) eine **optimale** Einschließung seines Wertebereichs, da der Nenner **exakt** eingeschlossen wird und weil die komplexe Intervalldivision eine optimale Einschließung liefert, wenn Zähler und Nenner voneinander unabhängige Intervalle sind. Dies ist in unserem Fall gewährleistet, da  $a, b, d$  durch Punktintervalle einzuschließen sind. Die anschließende Addition von  $a \in \mathbb{C}$  ändert dann nach **Satz D.6** nichts an der optimalen Einschließung des Wertebereichs  $W_{T_1,Z}$ .

Mit  $a = 1 + 2i, b = -1 + i, d = -2 + 2i$  und  $Z = [1, 2] + i \cdot [-1, 0]$  vergleichen wir die Einschließungen, die mit den Ausdrücken  $T_1$  und  $T$  mit der Präzision `prec = 40` berechnet werden, man erhält:

$Z = X + i \cdot Y = [1, 2] + i \cdot [-1, 0]; \quad \text{Präzision: } \text{prec} = 40.$	
$T_1(z) = \frac{b - ad}{z + d} + a$	$T_1(Z) = ([0, 4.000000000000], [-3.41547594745, -5.000000000000e - 1])$
$T(z) = \frac{az + b}{z + d}$	$T(Z) = ([-5.000000000000e - 1, 5.000000000000], [-4.41547594744, -0])$

Man erkennt im Vergleich zu  $T(Z)$  deutlich die viel bessere Einschließung  $T_1(Z) \subset T(Z)$ , wobei  $T_1(Z)$  sogar die **optimale** Einschließung von  $W_{T,Z} = \{w \in \mathbb{C} \mid w = (az + b)/(z + d) \wedge z \in Z\}$  ist.

# E. Tabellen mathematischer Funktionen

## E.1. Reelle Punkt-Funktionen

Tabelle E.1.: Elementarfunktionen mit  $x, y, a, b$  vom Typ `MpfrClass`,  $n$ : `unsigned long int`;

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\ln(1+x)$	<code>lnp1(x)</code>
$x^2$	<code>sqr(x)</code>	$\log_2(x)$	<code>log2(x)</code>
$x^2 + y^2$	<code>x2py2(x,y)</code>	$\log_{10}(x)$	<code>log10(x)</code>
$x^2 - y^2$	<code>x2my2(x,y)</code>	$\ln(\sin(x))$	<code>ln_sin(x)</code>
$1/x$	<code>reci(x)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$x/(x^2 + y^2)$	<code>x_div_x2py2(x,y)</code>	$\ln(\sqrt{x^2 + y^2})$	<code>ln_sqrtx2y2(x,y)</code>
$(x^2 - y^2)/(x^2 + y^2)^2$	<code>Re_rz2(x,y)</code>	$\ln(\sqrt{(1+x)^2 + y^2})$	<code>ln_sqrtxp1_2y2(x,y)</code>
$2xy/(x^2 + y^2)^2$	<code>mIm_rz2(x,y)</code>	$x^k, k \in \mathbb{Z}$	<code>power(x,k)</code>
$\frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Re_r1pz2(x,y)</code>	$x^y$	<code>pow(x,y)</code>
$\frac{2xy}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Im_r1pz2(x,y)</code>	$\sin(x)$	<code>sin(x)</code>
$\sqrt{x}$	<code>sqr(x)</code>	$1/\sin(x)$	<code>csc(x)</code>
$\sqrt{n}, n \in \mathbb{N}_0$	<code>sqr_n(n)</code>	$\cos(x)$	<code>cos(x)</code>
$1/\sqrt{x}$	<code>sqr_r(x)</code>	$1/\cos(x)$	<code>sec(x)</code>
$\sqrt[3]{x}, x \in \mathbb{R}$	<code>cbrt(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqr(x,n)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt{x+1} - 1$	<code>sqrtp1m1(x)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt{1+x^2}$	<code>sqrtp1px2(x)</code>	$\arccos(x)$	<code>acos(x)</code>
$\sqrt{1-x^2}$	<code>sqrtp1mx2(x)</code>	$\arctan(x)$	<code>atan(x)</code>

Fortsetzung auf der nächsten Seite

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\sqrt{x^2 - 1}$	<code>sqrtx2m1(x)</code>	$\arctan(y/x)$	<code>atan2(y,x)</code>
$x/\sqrt{1+x^2}$	<code>xdsqrt1px2(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$x/\sqrt{1-x^2}$	<code>xdsqrt1mx2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$x/\sqrt{x^2-1}$	<code>xdsqrtx2m1(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$\sqrt{x^2+y^2}$	<code>sqrtx2y2(x,y)</code>	$\cosh(x)$	<code>cosh(x)</code>
$e^x$	<code>exp(x)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$2^x$	<code>exp2(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$\ln(x)$	<code>ln(x)</code>	$\operatorname{AGM}(x,y)$	<code>agm(x,y)</code>

Weitere Anmerkungen findet man auf Seite 38.



## E.2. Reelle Punkt-Funktionen der Mathematischen Physik

Tabelle E.2.: Funktionen der Mathematischen Physik mit  $x$  vom Typ `MpfrClass`

Funktionsterm	Aufruf	Anmerkung
$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>	monoton wachsend
$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>	monoton fallend
$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt, x > 0$	<code>gamma(x)</code>	Pole: $x=0, -1, -2, \dots$
$\Gamma'(x)$	<code>gamma_D(x)</code>	Pole: $x=0, -1, -2, \dots$
$\frac{1}{\Gamma(x)}$	<code>gamma_reci(x)</code>	überall differenzierbar
$\left(\frac{1}{\Gamma(x)}\right)'$	<code>gamma_reci_D(x)</code>	überall differenzierbar
$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$	<code>digamma(x)</code>	Pole: $x=0, -1, -2, \dots$
$\ln(\Gamma(x))$	<code>lngamma(x)</code>	$2k-1 \leq x \leq -2k \rightsquigarrow \text{NaN}, k=0,1,2,\dots$
$\ln( \Gamma(x) )$	<code>int k; lgamma(x,k);</code>	$k = \begin{cases} +1, & \Gamma(x) > 0 \\ -1, & \Gamma(x) < 0 \end{cases}$
$k!$	<code>factorial(k)</code>	unsigned long int k
$\zeta(x) = \sum_{k=1}^\infty k^{-x}, x > 1;$	<code>zeta(x)</code>	$x \neq +1$
$\zeta(k), k = 0, 2, 3, 4, \dots$	<code>zeta(k)</code>	unsigned long int k
$\operatorname{Ei}(x) = \gamma + \ln(x) + \sum_{k=1}^\infty \frac{x^k}{k \cdot k!}, x > 0;$	<code>Ei(x)</code>	$x=0 \rightsquigarrow -\text{Inf}; x < 0 \rightsquigarrow \text{NaN};$
$\operatorname{Li}_2(x) = -\int_0^x \frac{\ln(1-t)}{t} dt, x < 1;$	<code>Li2(x)</code>	$x > 1 \rightsquigarrow$ Nur Realteil!
$J_n(x) = \sum_{k=0}^\infty \frac{(-1)^k}{k! \Gamma(k+n+1)} \left(\frac{x}{2}\right)^{2k+n};$	<code>Jn(n,x)</code>	Bessel-Fkt. 1. Art; $n \in \mathbb{Z};$
$J_0(x);$	<code>J0(x)</code>	Bessel-Fkt. 1. Art;
$J_1(x);$	<code>J1(x)</code>	Bessel-Fkt. 1. Art;
$Y_n(x);$	<code>Yn(n,x)</code>	Bessel-Fkt. 2. Art; $n \in \mathbb{Z};$
$Y_0(x);$	<code>Y0(x)</code>	Bessel-Fkt. 2. Art;
$Y_1(x);$	<code>Y1(x)</code>	Bessel-Fkt. 2. Art;

Weitere Anmerkungen findet man auf Seite 45.

### E.3. Reelle Intervall-Funktionen

Tabelle E.3.: Elementarfunktionen mit  $x, y$  vom Typ `MpfiClass`,  $a$  vom Typ `MpfrClass`;

Funktion	Aufruf	Funktion	Aufruf
$ x $	<code>abs(x)</code>	$\ln(x)$	<code>ln(x)</code>
$x^2$	<code>sqr(x)</code>	$\ln(1+x)$	<code>lnp1(x)</code>
$x^2 + y^2$	<code>x2py2(x,y)</code>	$\log_2(x)$	<code>log2(x)</code>
$x^2 - y^2$	<code>x2my2(x,y)</code>	$\log_{10}(x)$	<code>log10(x)</code>
$1/x$	<code>reci(x)</code>	$\ln(\sin(x))$	<code>ln_sin(x)</code>
$x/(x^2 + y^2)$	<code>x_div_x2py2(x,y)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$(x^2 - y^2)/(x^2 + y^2)^2$	<code>Re_rz2(x,y)</code>	$\ln(\sqrt{x^2 + y^2})$	<code>ln_sqrtx2y2(x,y)</code>
$2xy/(x^2 + y^2)^2$	<code>mIm_rz2(x,y)</code>	$\ln(\sqrt{(1+x)^2 + y^2})$	<code>ln_sqrtxp1_2y2(x,y)</code>
$\frac{1+x^2-y^2}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Re_r1pz2(x,y)</code>	$x^k, k \in \mathbb{Z}$	<code>power(x,k)</code>
$\frac{2xy}{4x^2y^2+(1+x^2-y^2)^2}$	<code>Im_r1pz2(x,y)</code>	$x^y$	<code>pow(x,y)</code>
$\sqrt{x}$	<code>sqr(x)</code>	$\sin(x)$	<code>sin(x)</code>
$\sqrt{n}, n \in \mathbb{N}_0$	<code>sqrtn(n)</code>	$1/\sin(x)$	<code>csc(x)</code>
$a/\sqrt{n}, n \in \mathbb{N}$	<code>xdivsqrtn(a,n)</code>	$\cos(x)$	<code>cos(x)</code>
$x/\sqrt{n}, n \in \mathbb{N}$	<code>xdivsqrtn(x,n)</code>	$1/\cos(x)$	<code>sec(x)</code>
$1/\sqrt{x}$	<code>sqrtr(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt[3]{x}$	<code>cbrt(x)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqrtn(x,n)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt{x+1} - 1$	<code>sqrtp1m1(x)</code>	$\arccos(x)$	<code>acos(x)</code>
$\sqrt{1+x^2}$	<code>sqrtp1x2(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{1-x^2}$	<code>sqrtp1mx2(x)</code>	$\arg(x + i \cdot y)$	<code>atan2(y,x)</code>
$\sqrt{x^2-1}$	<code>sqrtp1x2m1(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>

Fortsetzung auf der nächsten Seite

Fortsetzung der vorherigen Seite			
Funktion	Aufruf	Funktion	Aufruf
$x/\sqrt{1+x^2}$	<code>xdsqrt1px2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$x/\sqrt{1-x^2}$	<code>xdsqrt1mx2(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$x/\sqrt{x^2-1}$	<code>xdsqrtx2m1(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\sqrt{x^2+y^2}$	<code>sqrtox2y2(x,y)</code>	$1/\cosh(x)$	<code>sech(x)</code>
$e^x$	<code>exp(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$2^x$	<code>exp2(x)</code>	$\coth(x)$	<code>coth(x)</code>
$10^x$	<code>exp10(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{AGM}(x,y)$	<code>agm(x,y)</code>

Weitere Anmerkungen findet man auf Seite 69.

## E.4. Reelle Intervall-Funktionen der Mathematischen Physik

Tabelle E.4.: Funktionen der Mathematischen Physik mit  $x$  vom Typ `MpfiClass`

Funktionsterm	Aufruf	Anmerkung
$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>	monoton wachsend
$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>	monoton fallend
$\Gamma'(x)$	* <code>gamma_D(x)</code>	Pole: $x=0, -1, -2, \dots$
$\psi(x) = \frac{\Gamma'(x)}{\Gamma(x)}$	<code>digamma(x)</code>	Pole: $x=0, -1, -2, \dots$
$k!$	<code>ifactorial(k)</code>	unsigned long int k
$\zeta(k), k = 0, 2, 3, 4, \dots$	<code>izeta(k)</code>	unsigned long int k
$\operatorname{Ei}(x) = \gamma + \ln(x) + \sum_{k=1}^{\infty} \frac{x^k}{k \cdot k!}, x > 0;$	<code>Ei(x)</code>	$\operatorname{Inf}(x)=0 \rightsquigarrow -\operatorname{Inf}; \operatorname{Inf}(x)<0 \rightsquigarrow \operatorname{NaN};$

Weitere Anmerkungen findet man auf Seite 71.

## E.5. Reelle Funktionen zur Automatischen Differentiation

Tabelle E.5.: Funktionen vom Typ MPDerivType zur Automatischen Differentiation

Funktion	Aufruf	Funktion	Aufruf
$x^2$	<code>sqr(x)</code>	$\ln(\cos(x))$	<code>ln_cos(x)</code>
$x^n, n \in \mathbb{Z}$	<code>power(x,n)</code>	$\sin(x)$	<code>sin(x)</code>
$x^y$	<code>pow(x, y)</code>	$\cos(x)$	<code>cos(x)</code>
$1/x$	<code>reci(x)</code>	$\tan(x)$	<code>tan(x)</code>
$\sqrt{x}$	<code>sqrt(x)</code>	$\cot(x)$	<code>cot(x)</code>
$\sqrt[3]{x}, x \in \mathbb{R}$	<code>cbrt(x)</code>	$\arcsin(x)$	<code>asin(x)</code>
$\sqrt[n]{x}, n = 2, 3, \dots$	<code>sqrt(x,n)</code>	$\arccos(x)$	<code>acos(x)</code>
$1/\sqrt{x}$	<code>sqrt_r(x)</code>	$\arctan(x)$	<code>atan(x)</code>
$\sqrt{1+x^2}$	<code>sqrt1px2(x)</code>	$\operatorname{arccot}(x)$	<code>acot(x)</code>
$\sqrt{1-x^2}$	<code>sqrt1mx2(x)</code>	$\sinh(x)$	<code>sinh(x)</code>
$\sqrt{x^2-1}$	<code>sqrtx2m1(x)</code>	$\cosh(x)$	<code>cosh(x)</code>
$\sqrt{x+1}-1$	<code>sqrtp1m1(x)</code>	$\tanh(x)$	<code>tanh(x)</code>
$e^x$	<code>exp(x)</code>	$\operatorname{coth}(x)$	<code>coth(x)</code>
$e^x - 1$	<code>expm1(x)</code>	$\operatorname{arsinh}(x)$	<code>asinh(x)</code>
$e^{x^2}$	<code>expx2(x)</code>	$\operatorname{arcosh}(x)$	<code>acosh(x)</code>
$e^{x^2} - 1$	<code>expx2m1(x)</code>	$\operatorname{arcosh}(1+x)$	<code>acoshp1(x)</code>
$e^{-x^2}$	<code>expmx2(x)</code>	$\operatorname{artanh}(x)$	<code>atanh(x)</code>
$e^{-x^2} - 1$	<code>expmx2m1(x)</code>	$\operatorname{arcoth}(x)$	<code>acoth(x)</code>
$2^x$	<code>exp2(x)</code>	$1/\sin(x)$	<code>csc(x)</code>
$10^x$	<code>exp10(x)</code>	$1/\cos(x)$	<code>sec(x)</code>
$\ln(x)$	<code>ln(x)</code>	$1/\sinh(x)$	<code>csch(x)</code>
$\ln(1+x)$	<code>lnp1(x)</code>	$1/\cosh(x)$	<code>sech(x)</code>

*Fortsetzung auf der nächsten Seite*

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\log_2(x)$	<code>log2(x)</code>	$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	<code>erf(x)</code>
$\log_{10}(x)$	<code>log10(x)</code>	$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>
$\ln(\sin(x))$	<code>ln_sin(x)</code>		

Weitere Anmerkungen findet man auf Seite 128.

## E.6. Reelle Funktionen zur Taylorarithmetik

Tabelle E.6.: Funktionen vom Typ MPitaylor zur Taylor-Arithmetik

Funktion	Aufruf	Funktion	Aufruf
$x^2$	sqr(x)	$\cos(x)$	cos(x)
$\sqrt{x}$	sqrt(x)	$\tan(x)$	tan(x)
$\sqrt[n]{x}$ , $n = 2, 3, \dots$	sqrt(x, n)	$\cot(x)$	cot(x)
$\sqrt{1+x^2}$	sqrt1px2(x)	$\arcsin(x)$	asin(x)
$\sqrt{1-x^2}$	sqrt1mx2(x)	$\arccos(x)$	acos(x)
$\sqrt{x^2-1}$	sqrtx2m1(x)	$\arctan(x)$	atan(x)
$\sqrt{x+1}-1$	sqrtp1m1(x)	$\operatorname{arccot}(x)$	acot(x)
$x^y$	pow(x, y)	$\sinh(x)$	sinh(x)
$x^n$ , $n \in \mathbb{Z}$	power(x, n)	$\cosh(x)$	cosh(x)
$e^x$	exp(x)	$\tanh(x)$	tanh(x)
$2^x$	exp2(x)	$\operatorname{coth}(x)$	coth(x)
$10^x$	exp10(x)	$\operatorname{arsinh}(x)$	asinh(x)
$e^x - 1$	expm1(x)	$\operatorname{arcosh}(x)$	acosh(x)
$e^{x^2} - 1$	expx2m1(x)	$\operatorname{artanh}(x)$	atanh(x)
$e^{-x^2} - 1$	expmx2m1(x)	$\operatorname{arcoth}(x)$	acoth(x)
$\ln(x)$	ln(x)	$\operatorname{arcosh}(1+x)$	acoshp1(x)
$\log_2(x)$	log2(x)	$1/\sin(x)$	csc(x)
$\log_{10}(x)$	log10(x)	$1/\cos(x)$	sec(x)
$\ln(1+x)$	lnp1(x)	$1/\sinh(x)$	csch(x)
$\ln(\sin(x))$	ln_sin(x)	$1/\cosh(x)$	sech(x)
$\ln(\cos(x))$	ln_cos(x)	$\operatorname{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-t^2} dt$	erf(x)
<i>Fortsetzung auf der nächsten Seite</i>			

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\sin(x)$	<code>sin(x)</code>	$\operatorname{erfc}(x) = \frac{2}{\sqrt{\pi}} \int_x^\infty e^{-t^2} dt$	<code>erfc(x)</code>

Weitere Anmerkungen findet man auf Seite 139.



## E.7. Komplexe Punkt-Funktionen

Tabelle E.7.: Elementarfunktionen mit  $z = x + i \cdot y$ ,  $a, b$  vom Typ `MpfcClass`

Funktion	Aufruf	Funktion	Aufruf
$\text{conj}(z) = x - i \cdot y$	<code>conj(z)</code>	$e^z - 1$	<code>expm1(z)</code>
$\text{Re}(z) = x$	<code>Re(z)</code>	$e^{z^2}$	<code>expx2(z)</code>
$\text{Im}(z) = y$	<code>Im(z)</code>	$e^{z^2} - 1$	<code>expx2m1(z)</code>
$ z $	<code>abs(z)</code>	$e^{-z^2}$	<code>expmx2(z)</code>
$\text{arg}(z)$	<code>arg(z)</code>	$e^{-z^2} - 1$	<code>expmx2m1(z)</code>
$\text{arg}(1 + z)$	<code>argp1(z)</code>	$2^z$	<code>exp2(z)</code>
$z^2$	<code>sqr(z)</code>	$10^z$	<code>exp10(z)</code>
$z^2 + a \cdot z + b$	<code>poly2(z,a,b)</code>	$z^n, n \in \mathbb{Z}$	<code>power(z,n)</code>
$1/z$	<code>reci(z)</code>	$z^r, r \in \mathbb{R}$	<code>pow(z,r)</code>
$1/z^2$	<code>reci_z2(z)</code>	$z^w, w \in \mathbb{C}$	<code>pow(z,w)</code>
$1/(1 + z^2)$	<code>reci_1pz2(z)</code>	$\sin(z)$	<code>sin(z)</code>
$1/(1 - z^2)$	<code>reci_1mz2(z)</code>	$\cos(z)$	<code>cos(z)</code>
$\sqrt{z}$	<code>sqr(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z}$	<code>sqr_all(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{z+1} - 1$	<code>sqrtp1m1(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$\sqrt{1+z^2}$	<code>sqrtp1px2(z)</code>	$\arccos(z)$	<code>acos(z)</code>
$\sqrt{1-z^2}$	<code>sqrtp1mx2(z)</code>	$\arctan(z)$	<code>atan(z)</code>
$\sqrt{z^2-1}$	<code>sqrtpx2m1(z)</code>	$\text{arccot}(z)$	<code>acot(z)</code>
$\sqrt{-z^2-1}$	<code>sqrtpmx2m1(z)</code>	$\sinh(z)$	<code>sinh(z)</code>
$\sqrt[n]{z}, n \in \mathbb{Z}$	<code>sqr(z,n)</code>	$\cosh(z)$	<code>cosh(z)</code>
$\sqrt[n]{z}, n \in \mathbb{Z}$	<code>sqr_all(z,n)</code>	$\tanh(z)$	<code>tanh(z)</code>
$\log(z)$	<code>ln(z)</code>	$\coth(z)$	<code>coth(z)</code>

Fortsetzung auf der nächsten Seite

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\log(1 + z)$	<code>lnp1(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$\log_2(z)$	<code>log2(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$\log_{10}(z)$	<code>log10(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$e^z$	<code>exp(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>

Weitere Anmerkungen findet man auf Seite 93.

## E.8. Komplexe Intervall-Funktionen

Tabelle E.8.: Funktionen mit  $z = x + i \cdot y$ ,  $a, b$  vom Typ `MpfiClass`;  $x, y$  vom Typ `MpfiClass`

Funktion	Aufruf	Funktion	Aufruf
$\text{conj}(z) = x - i \cdot y$	<code>conj(z)</code>	$2^z$	<code>exp2(z)</code>
$\text{Re}(z) = x$	<code>Re(z)</code>	$10^z$	<code>exp10(z)</code>
$\text{Im}(z) = y$	<code>Im(z)</code>	$e^z$	<code>exp(z)</code>
$ z $	<code>abs(z)</code>	$e^z - 1$	<code>expm1(z)</code>
$\text{Arg}(z)$	<code>Arg(z)</code>	$e^{z^2}$	<code>expx2(z)</code>
$\text{arg}(z)$	<code>arg(z)</code>	$e^{z^2} - 1$	<code>expx2m1(z)</code>
$\text{arg}(1 + z)$	<code>argp1(z)</code>	$e^{-z^2}$	<code>expmx2(z)</code>
$z^2$	<code>sqr(z)</code>	$e^{-z^2} - 1$	<code>expmx2m1(z)</code>
$z^2 + a \cdot z + b$	<code>poly2(z, a, b)</code>	$z^n, n \in \mathbb{Z}$	<code>power_fast(z, n)</code>
$1/z$	<code>reci(z)</code>	$z^n, n \in \mathbb{Z}$	<code>power(z, n)</code>
$1/z^2$	<code>reci_z2(z)</code>	$z^p, p: \text{MpfiClass}$	<code>pow(z, p)</code>
$1/(1 + z^2)$	<code>reci_1pz2(z)</code>	$z^p, p: \text{MpfiClass}$	<code>pow_all(z, p)</code>
$1/(1 - z^2)$	<code>reci_1mz2(z)</code>	$z^w, w: \text{MpfiClass}$	<code>pow(z, w)</code>
$\sqrt{z}$	<code>sqr(z)</code>	$\sin(z)$	<code>sin(z)</code>
$\sqrt{z}$	<code>sqr_all(z)</code>	$\cos(z)$	<code>cos(z)</code>
$1/\sqrt{z}$	<code>sqr_r(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z+1} - 1$	<code>sqrtp1m1(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{1+z^2}$	<code>sqrtp1px2(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$\sqrt{1-z^2}$	<code>sqrtp1mx2(z)</code>	$\arccos(z)$	<code>acos(z)</code>
$\sqrt{z^2-1}$	<code>sqrtpx2m1(z)</code>	$\arctan(z)$	<code>atan(z)</code>
$\sqrt{-z^2-1}$	<code>sqrtpmx2m1(z)</code>	$\text{arccot}(z)$	<code>acot(z)</code>
$\sqrt[n]{z}$	<code>sqr(z, n)</code>	$\sinh(z)$	<code>sinh(z)</code>
$\sqrt[n]{z}$	<code>sqr_all(z, n)</code>	$\cosh(z)$	<code>cosh(z)</code>
$\log(z)$	<code>Ln(z)</code>	$\tanh(z)$	<code>tanh(z)</code>
$\log(1 + z)$	<code>Ln1(z)</code>	$\text{coth}(z)$	<code>coth(z)</code>

Fortsetzung auf der nächsten Seite

<i>Fortsetzung der vorherigen Seite</i>			
Funktion	Aufruf	Funktion	Aufruf
$\log(z)$	<code>ln(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$\log(1+z)$	<code>lnp1(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$\log_2(z)$	<code>log2(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$\log_{10}(z)$	<code>log10(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>

Weitere Anmerkungen findet man auf Seite 116.

## E.9. Komplexe Funktionen zur Taylorarithmetik

Tabelle E.9.: Funktionen vom Typ `MPcitaylor` zur Taylor-Arithmetik,  $z$  vom Typ `MPcitaylor`,  $w$  vom Typ `Mpfciclass`;

Funktion	Aufruf	Funktion	Aufruf
$z^2$	<code>sqr(z)</code>	$\log_{10}(z)$	<code>log10(z)</code>
$\sqrt{z}$	<code>sqrt(z)</code>	$\ln(1+z)$	<code>lnp1(z)</code>
$\sqrt[n]{z}, n=2,3,\dots$	<code>sqrt(z,n)</code>	$\sin(z)$	<code>sin(z)</code>
$\sqrt{1+z^2}$	<code>sqrt1px2(z)</code>	$\cos(z)$	<code>cos(z)</code>
$\sqrt{1-z^2}$	<code>sqrt1mx2(z)</code>	$\tan(z)$	<code>tan(z)</code>
$\sqrt{z^2-1}$	<code>sqrtox2m1(z)</code>	$\cot(z)$	<code>cot(z)</code>
$\sqrt{z+1}-1$	<code>sqrtp1m1(z)</code>	$\arcsin(z)$	<code>asin(z)</code>
$z^w$	<code>pow(z,w)</code>	$\arccos(z)$	<code>acos(z)</code>
$z^n, n \in \mathbb{Z}$	<code>power(z,n)</code>	$\arctan(z)$	<code>atan(z)</code>
$e^z$	<code>exp(z)</code>	$\operatorname{arccot}(z)$	<code>acot(z)</code>
$2^z$	<code>exp2(z)</code>	$\sinh(z)$	<code>sinh(z)</code>
$10^z$	<code>exp10(z)</code>	$\cosh(z)$	<code>cosh(z)</code>
$e^z - 1$	<code>expm1(z)</code>	$\tanh(z)$	<code>tanh(z)</code>
$e^{z^2}$	<code>expx2(z)</code>	$\coth(z)$	<code>coth(z)</code>
$e^{z^2} - 1$	<code>expx2m1(z)</code>	$\operatorname{arsinh}(z)$	<code>asinh(z)</code>
$e^{-z^2}$	<code>expmx2(z)</code>	$\operatorname{arcosh}(z)$	<code>acosh(z)</code>
$e^{-z^2} - 1$	<code>expmx2m1(z)</code>	$\operatorname{artanh}(z)$	<code>atanh(z)</code>
$\ln(z)$	<code>ln(z)</code>	$\operatorname{arcoth}(z)$	<code>acoth(z)</code>
$\log_2(z)$	<code>log2(z)</code>		

Weitere Anmerkungen findet man auf Seite 147.



## F. Laufzeitvergleiche

Mit der Current-Präzision  $\text{prec} = 53$  werden die Laufzeiten der Langzahl-Bibliotheken MPFR, MPFI, MPFC und MPFCI jeweils verglichen mit den Laufzeiten der Arithmetiken mit den Datentypen `real`, `interval`, `complex` und `cinterval`. Dabei zeigt sich, dass die Laufzeiten mit den letzten vier Datentypen etwa zehnmal günstiger sind als mit den entsprechenden Langzahl-Bibliotheken mit der Präzision  $\text{prec} = 53$ .

Ein zweiter Laufzeitvergleich mit den C-XSC-Datentypen `l_real`, `l_interval`, `l_complex` und `l_cinterval` zeigt die große Überlegenheit der MPF\*-Bibliotheken, wenn die Präzision größer als  $\text{prec} = 53$  Bits, d.h. größer als 16 Dezimalstellen gewählt wird. Bei den genannten Staggered Correction Arithmetiken wird die Präzision mit der Variablen `stagprec` festgelegt, wobei z.B. `stagprec = 15` eine Präzision von etwa  $15 \cdot 16 = 240$  Dezimalstellen definiert.

$x = 200000000, \quad f(x) = \sin(x);$				
real	MPFR	stagprec = 2	stagprec = 5	stagprec = 15
1/10	1	177	305	844, 4.89

$x = [0.25, 0.5], \quad f(x) = \arcsin(x);$				
interval	MPFI	stagprec = 2	stagprec = 5	stagprec = 15
1/23	1	39	38	76

$z = 2 + i \cdot 3, \quad f(z) = e^z, \quad z \in \mathbb{C}$				
complex	MPFC	stagprec = 2	stagprec = 5	stagprec = 15
1/9	1	125	309	838, 2.75

$Z = [1, 2] + i \cdot [2, 3], \quad f(z) = \arctan(z), \quad z \in Z$				
cinterval	MPFCI	stagprec = 2	stagprec = 5	stagprec = 15
1/13	1	39	56	154

### Anmerkungen:

- In der letzten Zeile bedeutet 1/13, dass die Laufzeit mit `cinterval` 13-mal kleiner ist als die Laufzeit mit MPFCI und der entsprechenden Präzision von  $\text{prec} = 53$  Bits.
- In der letzten Zeile bedeutet 154, dass die Laufzeit mit `l_cinterval` und `stagprec = 15`  $\approx 240$  Dezimalstellen ganze 154-mal größer ist als die Laufzeit mit der MPFCI-Bibliothek mit der Präzision von  $\text{prec} = 798$  Bits, die etwa 240 Dezimalstellen entspricht.
- Der große Laufzeitvorteil der vier MPF\*-Bibliotheken beruht auf deren direkten Hardware-Zugriff, während die Staggered-Arithmetiken den langen Akkumulator benutzen, der leider nur software-mäßig simuliert vorhanden ist.

- Die grün unterlegten Werte beziehen sich auf den Laufzeitvergleich zwischen *Mathematica* und MPFR bzw. MPFC, wobei in *Mathematica* jeweils die zu `stagprec = 15` entsprechende Dezimalstellenzahl 240 gewählt wurde.



# Literaturverzeichnis

- [1] Abramowitz M. and Stegun I. *Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables*. National Bureau of Standards, Washington, 1964.
- [2] Auzinger, W. and Stetter H.J. *Accurate Arithmetic Results for Decimal Data on Non-Decimal Computers*. Computing 35, 141-151, 1985.
- [3] Bohlender, G.: What do we need beyond IEEE arithmetic? pp. 1-32 in: Ch. Ullrich: Contributions to Computer Arithmetic and Self-Validating Numerical Methods. J.C. Baltzer AG, Scientific Publishing Co., Basel, 1990.
- [4] Alefeld G. and Grigorieff, R. D. *Fundamentals of Numerical Computation (Computer-Oriented Numerical Analysis)*. Computing Supplementum 2, Springer-Verlag, Wien / New York, 1980.
- [5] Alefeld G. and Herzberger J. *Introduction to Interval Computations*. Academic Press, New York, 1983.
- [6] Alefeld G. and Herzberger J. *Einführung in die Intervallrechnung*. BI, Reihe Informatik/12, 1983.
- [7] Behnke H., Sommer F. *Theorie der analytischen Funktionen einer komplexen Veränderlichen*. Springer, Berlin, 1962.
- [8] American National Standards Institute/Institute of Electrical and Electronics Engineers: "IEEE Standard for Binary Floating-Point Arithmetic"; ANSI/IEEE Std 754-1985, New York, 1985.
- [9] Blomquist, F.; Hofschuster, W.; Krämer, W.: Realisierung der hyperbolischen Cotangensfunktion in einer Staggered-Correction-Intervallarithmetic in C-XSC. Preprint 2004/3, Wissenschaftliches Rechnen / Softwaretechnologie, Universität Wuppertal, 2004.
- [10] Blomquist, F.: Automatische a priori Fehlerabschätzungen zur Entwicklung optimaler Algorithmen und Intervallfunktionen in C-XSC. Universität Wuppertal, 425 Seiten, Nov. 2005. [http://www2.math.uni-wuppertal.de/~xsc/literatur/a\\_priori.pdf](http://www2.math.uni-wuppertal.de/~xsc/literatur/a_priori.pdf)
- [11] Blomquist, F.; Hofschuster, W.; Krämer, W.: Real and Complex Taylor Arithmetic in C-XSC. Universität Wuppertal, 57 Seiten, 2005. [http://www2.math.uni-wuppertal.de/wrswt/preprints/prep\\_05\\_4.pdf](http://www2.math.uni-wuppertal.de/wrswt/preprints/prep_05_4.pdf)
- [12] Blomquist F.: Verbesserungen der komplexen Standardfunktionen von Markus Neher (see [54]). Interne Mitteilung, Bergische Universität Wuppertal, 2005.
- [13] Blomquist, F.; Hofschuster, W.; Krämer, W.; Neher, M.: Complex Interval Functions in C-XSC. Preprint BUW-WRSWT 2005/2, Bergische Universität Wuppertal, pp 1-48, 2005.
- [14] Blomquist, F., Hofschuster, W. and Krämer, W.: *A Modified Staggered Correction Arithmetic with Enhanced Accuracy and Very Wide Exponent Range*. Lecture Notes in Computer Science LNCS 5492, Springer-Verlag Berlin Heidelberg, 41-67, 2009.

- [15] Blomquist, F.: Staggered Correction Computations with Enhanced Accuracy and Extremely Wide Exponent Range. *Reliable Computing*, Vol. 15, pp. 26-35, May, 2011.
- [16] Brand, Hans-Stephan: Integration und Test einer Langzahlintervallbibliothek in C-XSC. Bachelor-Arbeit, Universität Wuppertal, 2010.
- [17] Braune, K., Krämer, W.: *High Accuracy Standard Functions for Real and Complex Intervals*. In Kaucher E., Kulisch U. and Ullrich Ch., editors, *Computerarithmetic: Scientific Computation and Programming Languages*, pages 81-114. Teubner, Stuttgart, 1987.
- [18] Braune, K.: Standard Functions for Real and Complex Point and Interval Arguments with Dynamic Accuracy. *Computing Supplementum*, 6:159-184, 1988.
- [19] Corliss, G. F.: *Automatic Differentiation Bibliography*. In [26], pp 331-353, 1991.
- [20] Hammer, R.; Hocks, M.; Kulisch, U.; Ratz, D.: *C++ Toolbox for Verified Computing: Basic Numerical Problems*. Springer-Verlag, Berlin / Heidelberg / New York, 1995.
- [21] Laurent Fousse, Guillaume Hanrot, Vincent Lefèvre, Patrick Pélissier, Paul Zimmermann: MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Transactions on Mathematical Software (TOMS)*, Vol. 33 Issue 2, pp. June 2007.
- [22] Fritzsche, A.: *Grundkurs Funktionentheorie*. Spektrum Akademischer Verlag Heidelberg, 2009.
- [23] GNU MP LIBRARY, <http://gmplib.org/>
- [24] GNU MPFR LIBRARY, <http://www.mpfr.org/mpfr-current/mpfr.html>
- [25] Griewank, A.: *On Automatic Differentiation*. In Iri, M., Tanabe, K.: *Mathematical Programming: Recent Developments and Applications*, pp. 83-108, Kluwer Academic Publishers, 1989.
- [26] Griewank, A., Corliss, G.: *Automatic Differentiation of Algorithms: Theory, Implementation, and Applications*. Proceedings of Workshop on Automatic Differentiation at Breckenridge, SIAM, Philadelphia, 1991.
- [27] Hansen, E.: *Global Optimization Using Interval Analysis -The One-Dimensional Case*. *Journal of Optimization Theory and Applications* 29, pp 331-344, 1979.
- [28] Hansen, E.: *Global Optimization Using Interval Analysis -The Multi-Dimensional Case*. *Numerische Mathematik* 34, pp 247-270, 1980.
- [29] Hansen, E.: *Global Optimization Using Interval Analysis*. Marcel Dekker, New York, 1992.
- [30] Ratscheck, H.; Rokne, J.: *Computer Methods for the Range of Functions*. Ellis Horwood Limited, Chichester, 1984.
- [31] Ratscheck, H.; Rokne, J.: *New Computer Methods for Global Optimization*. Ellis Horwood Limited, Chichester, 1988.
- [32] Herzberger, J. (Ed): *Topics in Validated Computations*. Proceedings of IMACS-GAMM International Workshop on Validated Numerics, Oldenburg, 1993. North Holland, 1994.
- [33] IBM: *High-Accuracy Arithmetic Subroutine Library (ACRITH)*. IBM Deutschland GmbH, third edition, 1986.

- [34] Hofschuster, W., Krämer, W.: C-XSC – A C++ Class Library for Extended Scientific Computing. *Numerical Software with Result Verification*. R. Alt, A. Frommer, B. Kearfott, W. Luther (eds), Springer Lecture Notes in Computer Science, 2004.
- [35] Hofschuster, W.; Krämer, W.: FLIB, eine schnelle und portable Funktionsbibliothek für reelle Argumente und reelle Intervalle im IEEE-*double*-Format. Preprint 98/7 des IWRMM, Universität Karlsruhe, 227 Seiten, 1998.
- [36] Hofschuster, W., Krämer, W., and Neher, M.: *C-XSC and Closely Related Software Packages*. Lecture Notes in Computer Science LNCS 5492, Springer-Verlag Berlin Heidelberg, 68-102, 2009.
- [37] Klatte, R.; Kulisch, U.; Lawo, C.; Rauch, M.; Wiethoff, A.: C-XSC, A C++ Class Library for Extended Scientific Computing. Springer-Verlag, Berlin / Heidelberg / New York, 1993.
- [38] Krämer, W.: Inverse Standardfunktionen für reelle und komplexe Intervallargumente mit a priori Fehlerabschätzungen für beliebige Datenformate, Dissertation, Universität Karlsruhe, 1987.
- [39] Krämer W.: Inverse Standard Functions for Real and Complex Point and Interval Arguments with Dynamic Accuracy. *Computing Supplementum*, 6:185-212, 1988.
- [40] Krämer, W.: A priori Worst Case Error Bounds for Floating-Point Computations, IEEE Transactions on Computers, Vol. 47, No. 7, July 1998.
- [41] Krämer, W., Wolff von Gudenberg, J. (eds): *Scientific Computing, Validated Numerics, Interval Methods*, Kluwer Academic Publishers Boston/Dordrecht/London, 398 pages, 2001.
- [42] Krämer, W.: Mehrfachgenaue reelle und intervallmäßige Staggered-Correction Arithmetik mit zugehörigen Standardfunktionen, Bericht des Instituts für Angewandte Mathematik, Universität Karlsruhe, S. 1-80, 1988.
- [43] Krämer, W.: Die Berechnung von Funktionen und Konstanten in Rechenanlagen. Habilitationsschrift, Universität Karlsruhe 1993.
- [44] Krämer, W.: Multiple-Precision Computations with Result Verification, in: Adams, E., Kulisch, U.(editors): *Scientific Computing with Automatic Result Verification*. Academic Press, pp. 325-356, 1993.
- [45] Krämer, W.; Kulisch, U.; Lohner, R.: Numerical Toolbox for Verified Computing II. Springer Verlag. Draft version, Karlsruhe, 1994.
- [46] Kulisch, U.: *Computer Arithmetic and Validity – Theory, Implementation and Applications*. de Gruyter, Berlin, 2008.
- [47] Lerch, M., Tischler, G., Wolff von Gudenberg, J., Hofschuster, W., and Krämer, W.: *fi-lib++*, a Fast Interval Library Supporting Containment Computations. ACM Transactions on Mathematical Software Vol 32, Number 2, pp. 299-324, 2006.
- [48] Lohner, R., Wolff von Gudenberg, J.: Complex interval division with maximum accuracy. Proc. of the 7th IEEE Symposium on Computer Arithmetic in Urbana (Illinois), pp 332-336, IEEE Comp. Soc., 1985.
- [49] Lohner, R.: Interval arithmetic in staggered correction format. In: Adams, E., Kulisch, U.(Eds): *Scientific Computing with Automatic Result Verification*. Academic Press, San Diego, pp 301-321, 1993.
- [50] Moore, R.E.: *Interval Analysis*. Prentice-Hall, Englewood Cliffs, New Jersey, 1966

- [51] MPFR manual. April 2011.  
<http://www.mpfr.org/mpfr-current/mpfr.pdf>
- [52] MPFI library for arbitrary precision interval arithmetic.  
<http://cadadr.org/fm/package/mpfi.html>
- [53] Neher, M.: The mean value form for complex analytic functions. *Computing*, 67:255-268, 2001.
- [54] Neher M.: Complex Standard Functions and their Implementation in the CoStLy Library. Preprint Nr. 04/18, Universität Karlsruhe, 2004.
- [55] Neher, M.: Complex Standard Functions and Their Implementation in the CoStLy Library, ACM Transactions on Mathematical Software, Vol. 33, Number 1, 27 pages, 2007.
- [56] Neumaier, A.: Interval Methods for Systems of Equations. Cambridge University Press, Cambridge, 1990.
- [57] Rall, L. B.: *Applications of Software for Automatic Differentiation in Numerical Computation*. In [4], pp. 141-156, 1980.
- [58] Ratscheck, H.: Die Subdistributivität in der Intervallarithmetik. ZAMM 51, pp 189-192, 1971.
- [59] Ratz, D.: Automatische Ergebnisverifikation bei globalen Optimierungsproblemen. Dissertation, Universität Karlsruhe, 1992.
- [60] Revol, Nathalie and Rouillier, Fabrice: MPFI, a multiple precision interval arithmetic library based on MPFR. <http://perso.ens-lyon.fr/nathalie.revol/software.html>, 2001.
- [61] Rotmaier, B.: Die Berechnung der elementaren Funktionen mit beliebiger Genauigkeit. Dissertation, Universität Karlsruhe, 1971.
- [62] M. R. Spiegel: *KOMPLEXE VARIABLEN, Theorie und Anwendung*, Schaum's Überblicke, Aufgaben, 1991.
- [63] Stetter, H.J.: *Staggered Correction Representation, a Feasible Approach to Dynamic Precision*. Proceedings of the Symposium on Scientific Software, edited by Cai, Fosdick, Huang, China University of Science and Technology Press, Beijing, China, 1989.
- [64] Törn, A., Zilinskas, A.: *Global Optimization*. Lecture Notes in Computer Science, No. 350, Springer-Verlag, Berlin, 1989.
- [65] XSC website on programming languages for scientific computing with validation.  
<http://www.xsc.de>
- [66] [http://en.wikipedia.org/wiki/Wilkinson's\\_polynomial](http://en.wikipedia.org/wiki/Wilkinson's_polynomial)

# Stichwortverzeichnis

$\Gamma$ -Funktion, *siehe* gamma

$\arg(1 + z)$ , 251

$\psi$ -Funktion, *siehe* digamma

$\zeta$ -Funktion, *siehe* zeta

Abfragen

`isBounded`, 58, 103

`isEmpty`, 58

`isEven`, 25

`isInf`, 25, 58, 85, 103

`isInteger`, 25

`isNaN`, 25, 58, 85, 103

`isNeg`, 25, 58

`isNonNeg`, 58

`isNonPos`, 58

`isNumber`, 25, 85

`isOdd`, 25

`isPoint`, 58, 103

`isPos`, 25, 58

`isStrictlyNeg`, 58

`isStrictlyPos`, 58

`isZero`, 25, 58, 85, 103

Ableitungen, 125

`abs(X)`, 68

`abs(x)`, 33, 37

`abs(Z)`, 116

`abs(z)`, 90, 93, 252

`AbsMax`, 65

`AbsMin`, 65

Absolutbetrag

`MPFC`, 93

`MPFCI`, 116

`MPFI`, 68, 69

`MPFR`, 37

`acos(X)`, 68

`acos(x)`, 37

`acos(Z)`, 116, 339, 340

`acos(z)`, 93, 265

`acosh(X)`, 68

`acosh(x)`, 37

`acosh(Z)`, 116, 364

`acosh(z)`, 93

`acoshp1(X)`, 68

`acoshp1(x)`, 37, 202

`acot(X)`, 68

`acot(x)`, 37

`acot(Z)`, 116, 362

`acot(z)`, 93

`acoth(X)`, 68

`acoth(x)`, 37

`acoth(Z)`, 116, 366

`acoth(z)`, 93

`agm(X, Y)`, 69

`agm(x, y)`, 39

Akkumulator, 11

Algorithmen

Komplexe Intervalle

$\arccos(z)$ , 339

$\operatorname{arcosh}(z)$ , 364

$\arcsin(z)$ , 329

$\arctan(z)$ , 354

$\operatorname{arsinh}(z)$ , 363

$\cot(z)$ , 328

$\log(1 + z)$ , 296

$\log(z)$ , 293

$z^2$ , 287

$z^p$ , 367

Komplexe Nullstellen, 155

Komplexe Punktargumente

$1/\sqrt{z}$ , 258

$\arccos(z)$ , 265

$\arcsin(z)$ , 261

$\arg(z)$ , 250

$\cos(z)$ , 249

$\cosh(z)$ , 260

$\cot(z)$ , 249

$\coth(z)$ , 260

$\log(1 + z)$ , 266

$\log(z)$ , 252

$|z|$ , 252

$\sin(z)$ , 249

$\sinh(z)$ , 260

$\sqrt{z}$ , 256

$\tan(z)$ , 249

$\tanh(z)$ , 260

$e^z$ , 248

$z^2$ , 256  
 Reelle Punktargumente  
    $\operatorname{arcosh}(1+x)$ , 202  
    $\log(\cos(x))$ , 200  
    $\log(\sin(x))$ , 199  
    $\log(\sqrt{x^2+y^2})$ , 201  
    $\sqrt{x^2-1}$ , 197  
    $x^2+y^2$ , 184  
    $x^2-y^2$ , 184  
 Alle Potenzen, 367  
 Alle Wurzeln, 94, 301, 304  
 AllZeros, 158  
 Anwendungsprogramme, *siehe* Programme  
 Arg(Z), 116  
 arg(Z), 116  
 arg(z), 93, 250  
 argp1(Z), 116  
 argp1(z), 93, 251  
 Argumentfunktionen, *siehe* Arg, arg  
 Argumentintervall  
   komplex Z, 245  
   reell U, V, 243  
   reell X, Y, 243  
 Arithmetisch-Geometrisches Mittel, 39, 69  
 asin(X), 68  
 asin(x), 37  
 asin(Z), 116, 244, 329  
 asin(z), 93, 261  
 asinh(X), 68  
 asinh(x), 37  
 asinh(Z), 116, 363  
 asinh(z), 93  
 atan(X), 68  
 atan(x), 37  
 atan(Z), 116, 354  
 atan(z), 93  
 ATAN2(Y, X), 70  
 atan2(Y, X), 68  
 atan2(y, x), 37  
 atanh(X), 68  
 atanh(x), 37  
 atanh(Z), 116, 365  
 atanh(z), 93  
 Ausgabe, *siehe* Eingabe/Ausgabe  
 Auslöschung, 124, 268, 271  
 Automatische Differentiation, 125  
   Beispiel, 125  
  
 Basis, *siehe* Eingabe/Ausgabe  
 Beispielprogramme, *siehe* Programme  
 Besselfunktionen  
   Erster Art  $J_n(x)$ , 44  
   Nullstellen, 121  
   Zweiter Art  $Y_n(x)$ , 44, 45  
 Blow(X, eps), 65  
 Blow(Z, eps), 114  
  
 C-XSC, 11, 12, 15  
 cbrt(X), 68  
 cbrt(x), 37  
 Ceil(x), 34  
 common\_decimals, 66, 123  
 comp, 33  
 conj(Z), 114  
 conj(z), 90  
 cos(X), 68  
 cos(x), 37  
 cos(Z), 116  
 cos(z), 93, 249  
 cosh(X), 68  
 cosh(x), 37  
 cosh(Z), 116  
 cosh(z), 93, 260  
 cot(X), 68  
 cot(x), 37  
 cot(Z), 116  
 cot(z), 93, 249  
 coth(X), 68  
 coth(x), 37  
 coth(Z), 116  
 coth(z), 93, 260  
 csc(X), 68  
 csc(x), 37  
 csch(X), 68  
 csch(x), 37  
 Current-Präzision, 17, 51, 62, 66, 67, 73, 75,  
   142, 150, 268  
 Current-Precision, *siehe* Current-Präzision  
 Current-Rundungsmodus, 19, 20, 23, 25, 30,  
   34, 266  
  
 Default-Präzision, 20  
 Default-Rundungsmodus, 16  
 Destruktor, *siehe* Konstruktor  
 Dezimalstellen, 11, 20, 23, 30  
   Ausgabe, 54  
   Gemeinsame beim Intervall, 66  
 diam, 65  
 Differentiation  
   Automatisch, 125  
   Numerisch, 125  
   Symbolisch, 125  
 digamma(X), 71  
 digamma(x), 44

Disjoint, 58  
 Disjoint( $X, Y$ ), 58  
 Drehstreckung, 381  
 Durchmesser  
   Absolut, 65  
   Relativ, 65  
 Durchschnitt, 62, 107  
   leer, 58  
  
 $E_i(X)$ , 71  
 $E_i(x)$ , 44  
 Eingabe/Ausgabe, 21, 30, 35, 66, 73, 81, 94,  
   301, 304, 367  
   cin, 19, 54, 81, 100  
   cout, 19, 54, 81, 100  
   GetBase, 20, 55, 77, 96, 100  
   SetBase, 20, 55, 77, 96, 100  
 Einschließung  
   Erste Nullstelle von  $J_0(x)$ , 121  
   exakte, 377, 381  
   Komplexe arithm. Ausdrücke, 119, **138**,  
     **243**, 320, **379**  
   Komplexe Nullstellen, 155  
   nahezu optimale, 113, 157, 186, 223, 291,  
     **375**, **376**  
   optimale, 11, 73, 74, 119, 124, **135**, 157,  
     207, **243**, 269, 288, 289, 298, 320  
   Reelle arithm. Ausdrücke, **74**, 124  
 Elementarfunktionen  
   Komplexe Intervallargumente, 116  
   Komplexe Punktargumente, 93  
   Reelle Intervallargumente, 68  
   Reelle Punktargumente, 37  
 Elliptische Integrale, 39  
 EmptyIntVal, 49  
 erf( $X$ ), 71  
 erf( $x$ ), 44  
 erfc( $X$ ), 71  
 erfc( $x$ ), 44  
 exp( $X$ ), 68  
 exp( $x$ ), 37  
 exp( $Z$ ), 116  
 exp( $z$ ), 93, 248  
 exp10( $X$ ), 68  
 exp10( $x$ ), 37  
 exp10( $Z$ ), 116  
 exp10( $z$ ), 93  
 exp2( $X$ ), 68  
 exp2( $x$ ), 37  
 exp2( $Z$ ), 116  
 exp2( $z$ ), 93  
 expm1( $X$ ), 68  
 expm1( $x$ ), 37  
 expm1( $Z$ ), 116  
 expm1( $z$ ), 271  
 expmx2( $X$ ), 68  
 expmx2( $x$ ), 37  
 expmx2( $Z$ ), 116, 118  
 expmx2( $z$ ), 93  
 expmx2m1( $X$ ), 68  
 expmx2m1( $x$ ), 37  
 expmx2m1( $Z$ ), 116  
 expmx2m1( $z$ ), 93  
 expo, 17, 24  
 Exponent, *siehe* expo  
 Exponentialfunktion, 32, 248  
 expx2( $X$ ), 68  
 expx2( $x$ ), 37  
 expx2( $Z$ ), 116, 118, 135  
 expx2( $z$ ), 93  
 expx2m1( $X$ ), 68  
 expx2m1( $x$ ), 37  
 expx2m1( $Z$ ), 116  
 expx2m1( $z$ ), 93  
 Extremalkurve, 205, 210, 217, 227, 234, 247  
 Extremalpunkte  $m, M$ , 205, 210, 217, 227,  
   234, 246, 340, 352  
  
 factorial( $k$ ), 44  
 Floor( $x$ ), 34  
 Frac( $x$ ), 34  
 Funktionen  
    $(1 + X^2 - Y^2)/(4X^2Y^2 + (1 + X^2 - Y^2)^2)$ ,  
     227  
    $(1 + x^2 - y^2)/(4x^2y^2 + (1 + x^2 - y^2)^2)$ , 189  
    $(1/\Gamma(x))'$ , 44  
    $(X^2 - Y^2)/(X^2 + Y^2)^2$ , 68  
    $(x^2 - y^2)/(x^2 + y^2)^2$ , 38, 186, 210  
    $-2xy/(x^2 + y^2)^2$ , 187  
    $1/(1 + z^2)$ , 255, 289  
    $1/(1 - z^2)$ , 265, 290  
    $1/X$ , 68  
    $1/Z$ , 116  
    $1/\Gamma(x)$ , 44  
    $1/\cos(X)$ , 68  
    $1/\cos(x)$ , 37  
    $1/\cosh(X)$ , 68  
    $1/\cosh(x)$ , 37  
    $1/\sin(X)$ , 68  
    $1/\sin(x)$ , 37  
    $1/\sinh(X)$ , 68  
    $1/\sinh(x)$ , 37  
    $1/\sqrt{X}$ , 68  
    $1/\sqrt{Z}$ , 116, 321

$1/\sqrt{x}$ , 37  
 $1/\sqrt{z}$ , 258  
 $1/x$ , 37  
 $1/z$ , 93, 253, 287  
 $1/z^2$ , 254, 288  
 $10^X$ , 68  
 $10^Z$ , 116  
 $10^x$ , 37  
 $10^z$ , 93  
 $2XY/(X^2 + Y^2)^2$ , 68  
 $2\cos(x) \cdot \sinh(y)/(\cos(2x) - \cosh(2y))$ ,  
195  
 $2\sin(x) \cdot \cosh(y)/(\cosh(2y) - \cos(2x))$ ,  
192  
 $2^X$ , 68  
 $2^Z$ , 116  
 $2^x$ , 37  
 $2^z$ , 93  
 $2xy/(4x^2y^2 + (1 + x^2 - y^2)^2)$ , 188, 233  
 $2xy/(x^2 + y^2)^2$ , 38, 187, 216  
 $X/(X^2 + Y^2)$ , 68  
 $X/\sqrt{n}$ , 68  
 $X^2$ , 68  
 $X^2 + A \cdot X + B$ , 226  
 $X^2 + Y^2$ , 68  
 $X^2 - Y^2$ , 68  
 $X^Y$ , 68  
 $X^k$ , 68  
 $Z^2 + A \cdot Z + B$ , 291  
 $Z^W$ , 116  
 $Z^n$ , 116  
 $Z^p$ , 116  
 $\Gamma'(X)$ , 71  
 $\Gamma'(x)$ , 44  
 $\Gamma(x)$ , 44  
 $\arccos(X)$ , 68  
 $\arccos(Z)$ , 116  
 $\arccos(x)$ , 37  
 $\arccos(z)$ , 93, 265  
 $\text{arccot}(X)$ , 68  
 $\text{arccot}(Z)$ , 116  
 $\text{arccot}(x)$ , 37  
 $\text{arccot}(z)$ , 93  
 $\text{arcosh}(1+x)$ , 37, 68, 202  
 $\text{arcosh}(X)$ , 68  
 $\text{arcosh}(Z)$ , 116  
 $\text{arcosh}(x)$ , 37  
 $\text{arcosh}(z)$ , 93  
 $\text{arcoth}(X)$ , 68  
 $\text{arcoth}(Z)$ , 116  
 $\text{arcoth}(x)$ , 37  
 $\text{arcoth}(z)$ , 93  
 $\arcsin(X)$ , 68  
 $\arcsin(Z)$ , 116  
 $\arcsin(x)$ , 37  
 $\arcsin(z)$ , 93, 261  
 $\arctan(X)$ , 68  
 $\arctan(Y/X)$ , 68  
 $\arctan(Z)$ , 116  
 $\arctan(x)$ , 37  
 $\arctan(y/x)$ , 37  
 $\arctan(z)$ , 93  
 $\text{arsinh}(X)$ , 68  
 $\text{arsinh}(Z)$ , 116  
 $\text{arsinh}(x)$ , 37  
 $\text{arsinh}(z)$ , 93  
 $\text{artanh}(X)$ , 68  
 $\text{artanh}(Z)$ , 116  
 $\text{artanh}(x)$ , 37  
 $\text{artanh}(z)$ , 93  
 $\cos(X)$ , 68  
 $\cos(Z)$ , 116  
 $\cos(x)$ , 37  
 $\cos(z)$ , 93  
 $\cosh(X)$ , 68  
 $\cosh(Z)$ , 116  
 $\cosh(x)$ , 37  
 $\cosh(z)$ , 93  
 $\cot(X)$ , 68  
 $\cot(Z)$ , 116  
 $\cot(x)$ , 37  
 $\cot(z)$ , 93  
 $\coth(X)$ , 68  
 $\coth(Z)$ , 116  
 $\coth(x)$ , 37  
 $\coth(z)$ , 93  
 $\log(1+X)$ , 68  
 $\log(1+Z)$ , 116, 297, 299  
 $\log(1+x)$ , 37  
 $\log(1+z)$ , 266  
 $\log(1+z)$ , 93  
 $\log(X)$ , 68  
 $\log(Z)$ , 116  
 $\log(\Gamma(x))$ , 44  
 $\log(\cos(X))$ , 68  
 $\log(\cos(x))$ , 37, 200  
 $\log(|\Gamma(x)|)$ , 44  
 $\log(\sin(X))$ , 68  
 $\log(\sin(x))$ , 37, 199  
 $\log(\sqrt{(1+x)^2 + y^2})$ , 38, 69, 205, 267  
 $\log(\sqrt{X^2 + Y^2})$ , 68  
 $\log(\sqrt{x^2 + y^2})$ , 37, 201  
 $\log(x)$ , 37



$\log(z)$ , 252  
 $\log(z)$ , 93  
 $\log_2(Z)$ , 116  
 $\log_2(z)$ , 93  
 $\log_{10}(Z)$ , 116  
 $\log_{10}(z)$ , 93  
 $\text{AGM}(X, Y)$ , 69  
 $\text{AGM}(x, y)$ , 39  
 $\text{Arg}(Z)$ , 116  
 $\text{Ei}(X)$ , 71  
 $\text{Ei}(x)$ , 44  
 $J_0(x)$ , 44  
 $J_1(x)$ , 44  
 $J_n(x)$ , 44  
 $\text{Li}_2(x)$ , 44  
 $Y_0(x)$ , 44  
 $Y_1(x)$ , 44  
 $Y_n(x)$ , 44  
 $\arg(1 + Z)$ , 116  
 $\arg(1 + z)$ , 93  
 $\arg(Z)$ , 116  
 $\arg(z)$ , 93  
 $\text{erfc}(X)$ , 71  
 $\text{erfc}(x)$ , 44  
 $\text{erf}(X)$ , 71  
 $\text{erf}(x)$ , 44  
 $\psi(X) = \Gamma'(X)/\Gamma(X)$ , 71  
 $\psi(x) = \Gamma'(x)/\Gamma(x)$ , 44  
 $|X|$ , 68  
 $|Z|$ , 116  
 $|x|$ , 37  
 $|z|$ , 93  
 $\sin(X)$ , 68  
 $\sin(Z)$ , 116  
 $\sin(x)$ , 37  
 $\sin(z)$ , 93  
 $\sinh(X)$ , 68  
 $\sinh(Z)$ , 116  
 $\sinh(x)$ , 37  
 $\sinh(z)$ , 93  
 $\sqrt[3]{X}$ , 68  
 $\sqrt[3]{x}$ , 37  
 $\sqrt[n]{X}$ , 68  
 $\sqrt[n]{Z}$ , 116  
 $\sqrt[n]{x}$ , 37  
 $\sqrt[n]{z}$ , 93  
 $\sqrt{-Z^2 - 1}$ , 116, 319  
 $\sqrt{-z^2 - 1}$ , 93, 286  
 $\sqrt{1 + X^2}$ , 68  
 $\sqrt{1 + Z^2}$ , 116, 310  
 $\sqrt{1 + Z} - 1$ , 116  
 $\sqrt{1 + x^2}$ , 37  
 $\sqrt{1 + z^2}$ , 93, 279  
 $\sqrt{1 - X^2}$ , 68  
 $\sqrt{1 - Z^2}$ , 116, 314  
 $\sqrt{1 - x^2}$ , 37  
 $\sqrt{1 - z^2}$ , 93, 283  
 $\sqrt{X + 1} - 1$ , 68  
 $\sqrt{X^2 - 1}$ , 68  
 $\sqrt{Z^2 - 1}$ , 116, 315  
 $\sqrt{n}$ , 37, 68  
 $\sqrt{x + 1} - 1$ , 37  
 $\sqrt{x^2 - 1}$ , 37, 197  
 $\sqrt{z + 1} - 1$ , 93, 274  
 $\sqrt{z^2 - 1}$ , 93, 284  
 $\tan(X)$ , 68  
 $\tan(Z)$ , 116  
 $\tan(x)$ , 37  
 $\tan(z)$ , 93  
 $\tanh(X)$ , 68  
 $\tanh(Z)$ , 116  
 $\tanh(x)$ , 37  
 $\tanh(z)$ , 93  
 $\zeta(k)$ , 44, 71  
 $\zeta(x)$ , 44  
 $a/\sqrt{n}$ , 68  
 $e^X$ , 68  
 $e^X - 1$ , 68  
 $e^Z$ , 116  
 $e^Z - 1$ , 116  
 $e^x$ , 37  
 $e^x - 1$ , 37  
 $e^z$ , 93  
 $e^z$ , 248  
 $e^z - 1$ , 271  
 $e^{-X^2}$ , 68  
 $e^{-X^2} - 1$ , 68  
 $e^{-Z^2}$ , 116, 118  
 $e^{-Z^2} - 1$ , 116  
 $e^{-x^2}$ , 37  
 $e^{-x^2} - 1$ , 37  
 $e^{-z^2}$ , 93  
 $e^{-z^2} - 1$ , 93  
 $e^{X^2}$ , 68  
 $e^{X^2} - 1$ , 68  
 $e^{Z^2}$ , 116, 118, 135  
 $e^{Z^2} - 1$ , 116  
 $e^{x^2}$ , 37  
 $e^{x^2} - 1$ , 37  
 $e^{z^2}$ , 93  
 $e^{z^2} - 1$ , 93  
 $k!$ , 44, 71  
 $x/(x^2 + y^2)$ , 38, 185, 208  
 $x/\sqrt{1 + x^2}$ , 198, 241

$x/\sqrt{1-x^2}$ , 198, 241  
 $x/\sqrt{x^2-1}$ , 198, 241  
 $x^2$ , 37  
 $x^2 + a \cdot x + b$ , 191  
 $x^2 + y^2$ , 37, 184  
 $x^2 - y^2$ , 37, 184, 256  
 $x^k$ , 37  
 $x^y$ , 37  
 $z^2$ , 256, 287  
 $z^2 + a \cdot z + b$ , 257, 291  
 $z^n$ , 93  
 $z^r$ , 93  
 $z^w$ , 93  
 harmonisch, 205, 208, 210, 216, 227, 233, 244  
 holomorph, 205, 208, 210, 216, 227, 233, 244  
 mehrdeutig, 93, 245, 261, 265  
 separabel, 243  
      $\cos(z)$ , 244  
      $\cosh(z)$ , 244  
      $\sin(z)$ , 244  
      $\sinh(z)$ , 244  
      $e^z$ , 244  
 Funktionen der Mathematischen Physik  
     Reelle Intervallargumente, 71  
     Reelle Punktargumente, 45  
  
 $\gamma(x)$ , 44  
 Gamma-Funktion, *siehe*  $\gamma$   
 $\gamma_D(X)$ , 71  
 $\gamma_D(x)$ , 44  
 $\gamma_{\text{reci}}(x)$ , 44  
 $\gamma_{\text{reci}_D}(x)$ , 44  
 Gerichtete Rundung, *siehe* Rundung  
 GetBase, 20, 55, 77, 96, 100  
 GetCurrPrecision, 20, 55, 77, 96, 142, 151  
 GetCurrRndMode, 20, 77  
 GetPrecision, 20, 55, 77, 96, 143, 151  
 Globale Optimierung, 131, 135, 375  
 Grundoperationen  
     komplexe, 248, **287**  
     reelle, 11, 74  
  
 Höhenlinie, *siehe* Extremalkurve  
 harmonische Funktionen, 244, 311  
 Hauptwert, 245  
 Hauptzweig, 93  
 holomorphe Funktionen, 244, 311, 320  
 Horner-Schema, 160  
  
 IEEE, *siehe* Zahlenformat  
 ifactorial(k), 71  
  
 $\text{Im}(Z)$ , 113  
 $\text{Im}(z)$ , 91  
 $\text{Im\_csc}(x, y, \text{rnd})$ , 195  
 $\text{Im\_r1pz2}(X, Y)$ , 233  
 $\text{Im\_r1pz2}(x, y)$ , 188  
 $\text{Im\_rz2}(x, y)$ , 187  
 Imaginärteil, *siehe*  $\text{Im}$   
 $\text{in}(X, W)$ , 106  
 $\text{in}(x, W)$ , 106  
 $\text{in}(X, Y)$ , 61  
 $\text{in}(x, Y)$ , 61  
 $\text{in}(Z, W)$ , 106  
 $\text{in}(z, W)$ , 106  
 $\text{Inf}(X, \text{prec})$ , 65  
 $\text{Inf}(Z, \text{prec})$ , 113  
 Interval\_Scaling, 72  
 Intervall  
     common\_decimals, 66  
     Durchmesser  
         Absolut, 65  
         Relativ, 65  
     Gemeinsame Dezim.-Stellen, 66  
     Leeres, 49  
     Maximum der Absolutbeträge, 65  
     Minimum der Absolutbeträge, 65  
     Mittelpunkt, 65  
     Randpunkte vertauschen, 65  
     Skalarprodukt, 73  
     Skalierung, 72  
     Zwei Intervalle vertauschen, 66  
 Intervallauswertung  
     naiv, 271, 313  
     optimal, 288, 289, 313, 379  
 isBounded, 58, 103  
 isEmpty, 58  
 isEven, 25  
 isInf, 25, 58, 85, 103  
 isInteger, 25  
 isNaN, 25, 58, 85, 103  
 isNeg, 25, 58  
 isNonNeg, 58  
 isNonPos, 58  
 isNumber, 25, 85  
 isOdd, 25  
 isPoint, 58, 103  
 isPos, 25, 58  
 isStrictlyNeg, 58  
 isStrictlyPos, 58  
 isZero, 25, 58, 85, 103  
 izeta(k), 71  
  
 J0(x), 44

$J_1(x)$ , 44  
 $J_n(n, x)$ , 44  
 Konstanten  
      $\log(2)$ , 36, 67  
      $\pi$ , 36, 67  
      $e$ , 36, 67  
     Catalan, 36, 67  
 Konstruktor  
     MPcitaylor, 149  
     MPFC, 78  
     MPFCI, 97  
     MPFI, 49  
     MPFR, 18  
     MPitaylor, 141  
 Konvexe Hülle, 63, 108  
 Kreisring, 117, 367  
  
 Laufzeit, 12, 117, 203, 245, 355  
     Laufzeitvergleiche, 287, 399  
 $\lgamma(x, k)$ , 44  
 $Li_2(x)$ , 44  
 Lineare Funktion, 381  
 Liste  
     Auslesen, 94, 301, 304, 368  
 Listen, 94, 304, 367  
 $\ln(1+Z)$ , 299  
 $\ln(1+z)$ , 266  
 $\ln(X)$ , 68  
 $\ln(x)$ , 37  
 $\ln(Z)$ , 116, 294  
 $\ln(Z)$ , 116, 295  
 $\ln(z)$ , 93, 252  
 $\ln\_cos(X)$ , 68  
 $\ln\_cos(x)$ , 37, 200  
 $\ln\_sin(X)$ , 68  
 $\ln\_sin(x)$ , 37, 199  
 $\ln\_sqrtx_2y_2(X, Y)$ , 68  
 $\ln\_sqrtx_2y_2(x, y)$ , 37, 201  
 $\ln\_sqrtxp_1_2y_2(X, Y)$ , 205  
 $\ln\_sqrtxp_1_2y_2(x, y)$ , 38, 266  
 $\lngamma(x)$ , 44  
 $\lnp1(X)$ , 68  
 $\lnp1(x)$ , 37  
 $\lnp1(Z)$ , 116, 297  
 $\lnp1(Z)$ , 299  
 $\lnp1(z)$ , 93  
 $\log_{10}(Z)$ , 116  
 $\log_{10}(z)$ , 93  
 $\log_2(Z)$ , 116  
 $\log_2(z)$ , 93  
  
 Möbiustransformation, 382  
  
 $\text{mant}$ , 17, 24, 33  
 Maple, 333, 340  
 Maschinenzahlen, 191  
 Mathematica, 119, 333, 400  
     ComplexExpand[...], 379  
     FullSimplify[...], 231  
     Solve[...], 227, 234  
 $\text{max}$ , 33  
 $\text{MaxFloat}$ , 17, 24, 34, 281, 331  
 Maximum/Minimum  
     Harmonischer Funktionen, 244  
 $\text{MaxReal}$ , 17  
 $\text{mid}$ , 65  
 $\text{mid}(X)$ , 65  
 $\text{mid}(Z)$ , 113  
 $\text{mIm\_rz2}(X, Y)$ , 216  
 $\text{mIm\_rz2}(x, y)$ , 38, 187  
 $\text{min}$ , 33  
 $\text{minfloat}$ , 17, 24, 34, 329, 339, 354, 362  
 $\text{MinReal}$ , 17  
 $\text{minreal}$ , 17  
 $\text{MPFR\_PREC\_MAX}$ , 15  
 Multiplikation  
     Aufgerundet, 30  
     mit  $2^n$ , *siehe* times2pown  
  
 Nachfolger, *siehe* succ  
 Newton-Verfahren, 121, 155  
 Nullstellen  
     Besselfunktionen, 121  
     Eindeutigkeit, 121  
     Komplexer Ausdrücke, 155  
     Nichtlineare Funktionen, 158  
     Polynome, 170, 173  
 $\text{Number\_Scaling}(x, k, \text{rnd})$ , 34  
 $\text{Number\_Scaling\_S}(x, k, \text{rnd})$ , 34  
 Numerische Ergebnisse  
      $\arccos(Z)$ , 339  
      $\text{arccot}(Z)$ , 362  
      $\text{arcosh}(Z)$ , 364  
      $\text{arcoth}(Z)$ , 366  
      $\arcsin(Z)$ , 329, 338  
      $\arctan(Z)$ , 354, 360, 361  
      $\arg(1+z)$ , 251  
      $\arg(z)$ , 250  
      $\text{arsinh}(Z)$ , 363  
      $\log(1+Z)$ , 298, 299  
      $\log(1+z)$ , 270  
      $\log(Z)$ , 294, 295  
      $\log(\sqrt{(1+X)^2+Y^2})$ , 207  
      $\log(\sqrt{(1+x)^2+y^2})$ , 269  
      $\sqrt{Z}$ , 300, 303

$\sqrt{z+1} - 1$ , 277  
 $e^z - 1$ , 272

Operatoren

- Arithmetische
  - MPcitaylor, 150
  - MPFC, 87
  - MPFCI, 109
  - MPFI, 51
  - MPFR, 28
  - MPitaylor, 142
- Durchschnitt
  - MPFCI, 107
  - MPFI, 62
- Eingabe/Ausgabe
  - <<, 19, 54, 81, 100
  - >>, 19, 54, 81, 100
- Enthalten im Innern, *siehe in*
- Konvexe Hülle
  - MPFCI, 108
  - MPFI, 63
- Vergleiche
  - MPFC, 86
  - MPFCI, 104
  - MPFI, 59
  - MPFR, 26
- Zuweisungen
  - MPcitaylor, 150
  - MPFC, 80
  - MPFCI, 99
  - MPFI, 50
  - MPFR, 25
  - MPitaylor, 141

Optimale Intervallauswertung, *siehe* Intervallauswertung

Overflow, *siehe* Überlauf

poly2(X,A,B), 226  
 poly2(x,a,b), 191  
 poly2(Z,A,B), 291  
 poly2(z,a,b), 257  
 Polynom, 144, 153
 

- Komplexes
  - 2. Grades, 257, 291
- Wilkinson, 170

Potenzen, 37, 68, 92, 115, 367  
 pow(X,Y), 68  
 pow(x,y), 37  
 pow(Z,P), 116, 117  
 pow(z,r), 93  
 pow(Z,W), 116, 118  
 pow(z,w), 93

pow\_all(Z,P), 116, 367  
 power(X,k), 68  
 power(x,k), 37  
 power(Z,n), 116, 117  
 power(z,n), 93  
 power\_fast(Z,n), 116, 117

Präzision, 15, 21, 30, 35, 42, 66, 73, 94, 271, 301, 304, 367
 

- Default, 20
- GetCurrPrecision, 20, 55, 77, 96, 142, 151
- GetPrecision, 20, 55, 77, 96, 143, 151
- maximale, 15
- MPFR\_PREC\_MAX, 15
- RoundPrecision, 20, 55, 77, 96, 143, 151
- SetCurrPrecision, 20, 55, 77, 96, 142, 151
- SetPrecision, 20, 55, 77, 96, 143, 151
- Standard, 20

PrecisionType, 15, 48, 75, 95

pred, 17, 337  
 Prod\_H1, 40  
 prod\_H1, 40

Programme

- Alle Potenzen, 367
- Alle reelle Nullstellen, 161
- Alle Wurzeln, 94, 301, 304
- Automatische Differentiation, 125
- Eingabe/Ausgabe, 21, 30, 35, 66, 73, 94, 301, 304, 367
- Einschließung arithm. Ausdrücke, 124
- Erste Nullstelle von  $J_0(x)$ , 121
- Globale Optimierung, 132
- Komplexe Nullstellen, 155
- Komplexe Taylor-Arithmetik
  - Funktionsterme, 154
  - Polynom, 153
- Listen, 94, 304, 367
- MPFR-01, 21
- MPFR-02, 30
- MPFR-03, 42
- MPFR-04, 66
- MPFR-05, 35
- MPFR-06, 73
- MPFR-07, 301
- MPFR-08, 304
- MPFR-09, 367
- MPFR-10, 94
- MPFR-11, 155
- MPFR-12, 124
- MPFR-13, 121
- MPFR-14, 125

MPFR-15, 161  
 MPFR-16, 132  
 MPFR-17, 144  
 MPFR-18, 145  
 MPFR-19, 153  
 MPFR-20, 154  
 MPFR-21, 164  
 MPFR-22, 170  
 MPFR-23, 173  
 Präzision, 21, 30, 35, 42, 66, 73, 94, 301, 304, 367  
 Rundung, 21, 30, 35, 42, 66, 94  
 Skalarprodukt, 42, 73  
 Taylor-Arithmetik  
     Funktionsterme, 145  
     Polynom, 144  
 Zufallsintervalle, 66  
 Zufallszahlen, 35

random, 35, 66  
 Re( $Z$ ), 113  
 Re( $z$ ), 91  
 Re\_csc( $x, y, \text{rnd}$ ), 192  
 Re\_r1pz2( $X, Y$ ), 227  
 Re\_r1pz2( $x, y$ ), 189  
 Re\_rz2( $X, Y$ ), 210  
 Re\_rz2( $x, y$ ), 38, 186  
 Realteil, *siehe* Re  
 Rechteck, *siehe* Rechteckintervall  
 Rechteckintervall, 117, 205, 216, 223, 231, 233, 236, 243, 370  
 reci( $X$ ), 68  
 reci( $x$ ), 37  
 reci( $Z$ ), 116, 287  
 reci( $z$ ), 93, 253  
 reci\_1mz2( $Z$ ), 290  
 reci\_1mz2( $z$ ), 265  
 reci\_1pz2( $Z$ ), 289  
 reci\_1pz2( $z$ ), 255  
 reci\_z2( $Z$ ), 288  
 reci\_z2( $z$ ), 254  
 Reelle Taylor-Arithmetik,  
     *siehe* Taylor-Arithmetik  
 RelDiam, 65  
 Riemann'sche  $\zeta$ -Funktion, *siehe* zeta  
 Round( $x$ ), 34  
 RoundDown, 16, 20, 248  
 RoundFromZero, 16, 20  
 RoundingMode, 16, 75  
 RoundNearest, 16, 20, 179, 248  
 RoundPrecision, 20, 55, 77, 96, 143, 151  
 RoundToZero, 16, 20

RoundUp, 16, 20, 248  
 Rundung, 21, 30, 35, 42, 66, 94, 266  
     Arithm. Ausdruck, 179  
     Default, 16  
     GetCurrRndMode, 20, 77  
     nicht-optimale, 266  
     optimale, 266  
     SetCurrRndMode, 20, 77  
     Standard, 16  
 Rundungsparameter, 16, 20, 77, 248

scal\_prod, 41, 73  
 Scal\_prod\_k, 40  
 scal\_prod\_k, 40  
 sec( $X$ ), 68  
 sec( $x$ ), 37  
 sech( $X$ ), 68  
 sech( $x$ ), 37  
 set\_inf( $X$ ), 67  
 set\_inf( $x, k$ ), 35  
 set\_inf( $Z$ ), 114  
 set\_inf( $z, k$ ), 91  
 set\_Mpfc( $z$ ), 84  
 set\_Mpfc\_i( $Z$ ), 103  
 set\_Mpfi( $X$ ), 57  
 set\_Mpfr( $x$ ), 23  
 set\_nan( $X$ ), 67  
 set\_nan( $x$ ), 35  
 set\_nan( $Z$ ), 114  
 set\_nan( $z$ ), 91  
 set\_zero( $X$ ), 67  
 set\_zero( $x, k$ ), 35  
 set\_zero( $Z$ ), 114  
 set\_zero( $z$ ), 91  
 SetBase, 20, 55, 77, 96, 100  
 SetCurrPrecision, 20, 55, 77, 96, 142, 151  
 SetCurrRndMode, 20, 77  
 SetPrecision, 20, 55, 77, 96, 143, 151  
 sign, 24  
 sin( $X$ ), 68  
 sin( $x$ ), 37  
 sin( $Z$ ), 116  
 sin( $z$ ), 93, 249  
 sinh( $X$ ), 68  
 sinh( $x$ ), 37  
 sinh( $Z$ ), 116  
 sinh( $z$ ), 93, 260  
 Skalarprodukt, 11, 40–42, 72, 73  
     Maschinenintervalle, 72  
     Maschinenzahlen, 40  
 Skalierung, 72, 187–189, 223, 276, 291, 360  
     Intervall, 72

Zahl, 34  
 Spezielle Funktionen, *siehe* Funktionen der  
 Mathematischen Physik  
 Spiegelung  
 1. Winkelhalbierende, 216  
 Ursprung, 216, 233  
 y-Achse, 216, 233  
`sqr(X)`, 68  
`sqr(x)`, 37  
`sqr(Z)`, 287  
`sqr(z)`, 256  
`sqrt(X,n)`, 68  
`sqrt(x,n)`, 37  
`sqrt(Z)`, 300  
`sqrt(z)`, 256  
`sqrt(Z+1)-1`, 306  
`sqrt(Z,n)`, 302  
`sqrt(z,n)`, 93  
`sqrt1mx2(X)`, 68  
`sqrt1mx2(x)`, 37  
`sqrt1mx2(Z)`, 116, 314  
`sqrt1mx2(z)`, 93, 283  
`sqrt1px2(X)`, 68  
`sqrt1px2(x)`, 37  
`sqrt1px2(Z)`, 116, 310  
`sqrt1px2(z)`, 93, 279  
`sqrt_all(Z,n)`, 116  
`sqrt_all(z,n)`, 93  
`sqrt_all(Z)`, 301  
`sqrt_all(Z,n)`, 304  
`sqrt_N(n)`, 68  
`sqrt_n(n)`, 37  
`sqrt_r(X)`, 68  
`sqrt_r(x)`, 37  
`sqrt_r(Z)`, 116, 321  
`sqrt_r(z)`, 258  
`sqrtmx2m1(Z)`, 116, 319  
`sqrtmx2m1(z)`, 93, 286  
`sqrtp1m1(X)`, 68  
`sqrtp1m1(x)`, 37  
`sqrtp1m1(Z)`, 116  
`sqrtp1m1(z)`, 93, 274  
`sqrtx2m1(X)`, 68  
`sqrtx2m1(x)`, 37, 197  
`sqrtx2m1(Z)`, 116, 315  
`sqrtx2m1(z)`, 93, 284  
 Staggered-Arithmetik, 11  
 Standard-Präzision, 20  
 Standard-Rundungsmodus, 16  
 Streckung, 381  
 string, *siehe* Zeichenketten  
 succ, 17, 347

`sum_k_H1`, 40  
`Sup(X,prec)`, 65  
`Sup(Z,prec)`, 113  
 swap, 35, 65, 66  
 swap\_endpoints, 65  
`tan(X)`, 68  
`tan(x)`, 37  
`tan(Z)`, 116  
`tan(z)`, 93, 249  
`tanh(X)`, 68  
`tanh(x)`, 37  
`tanh(Z)`, 116  
`tanh(z)`, 93, 260  
`times2pown(X,n)`, 67  
`times2pown(x,n,rnd)`, 34  
`times2pown(Z,n)`, 114  
`times2pown(z,n,rnd)`, 91  
 Transformationsfunktion, 376, 381  
 Translation, 381  
`Trunc(x)`, 34  
 Tschebyscheffpolynom, 173  
 Typumwandlungen  
 MPFC, 82  
 MPFCI, 101  
 MPFI, 56  
 MPFR, 22  
 Überlauf, 11, 34, 41, 72, 184, 197, 201, 267  
 Überschätzungen, 243, 306, 310, 314, 370  
 Underflow, *siehe* Unterlauf  
 Unterlauf, 11, 17, 24, 43, 72, 91, 256, 335,  
 337  
 Variable  
 Initialisieren  
`set_inf(X)`, 67  
`set_inf(x,k)`, 35  
`set_inf(Z)`, 114  
`set_inf(z,k)`, 91  
`set_Mpfc(z)`, 84  
`set_MpfcI(Z)`, 103  
`set_Mpfi(X)`, 57  
`set_Mpfr(x)`, 23  
`set_nan(X)`, 67  
`set_nan(x)`, 35  
`set_nan(Z)`, 114  
`set_nan(z)`, 91  
`set_zero(X)`, 67  
`set_zero(x,k)`, 35  
`set_zero(Z)`, 114  
`set_zero(z)`, 91  
 Variablentyp

- PrecisionType, 15, 48, 75, 95
- RoundingMode, 16, 75
- Vergleichsoperatoren, 26, 59, 86, 104
- Verzweigungspunkt, 245, 261, 265, 293, 296
- Verzweigungsschnitt, 117, 118, 245, 250, 252,
  - 256, 258, 261, 265, 266, 284, 286,
  - 293, 295, 300, 303, 306, 329, 339,
  - 340, 354, 362–366
- $1/\sqrt{z}$ , 258
- $1/\sqrt{Z}$ , 321
- $\arccos(Z)$ , 339
- $\arccos(z)$ , 341
- $\operatorname{arccot}(Z)$ , 362
- $\operatorname{arcosh}(Z)$ , 364
- $\operatorname{arcoth}(Z)$ , 366
- $\arcsin(Z)$ , 329
- $\arcsin(z)$ , 261
- $\arctan(Z)$ , 354
- $\operatorname{arsinh}(Z)$ , 363
- $\operatorname{artanh}(Z)$ , 365
- $\log(1+z)$ , 266
- $\log(z)$ , 252
- $\operatorname{Ln}(Z)$ , 294, 295
- $\operatorname{Lnp1}(Z)$ , 297
- $\operatorname{lnp1}(Z)$ , 299
- $\sqrt[n]{Z}$ , 303
- $\sqrt{1-z^2}$ , 283
- $\sqrt{Z+1}-1$ , 306
- $\sqrt{Z}$ , 300
- $\sqrt{z}$ , 256
- Vorgänger, *siehe* pred
- Vorzeichen, *siehe* sign
- Wertebereich
  - komplex, 138, 308, 314, 318, 320, 379
  - reell, 74, 320
- Wilkinsonpolynom, 170
- Wurzelfunktionen
  - Alle dritte Wurzeln, 94
  - MPFC, 92
  - MPFCI, 115
  - MPFI, 68
  - MPFR, 37
- $x_2my_2$ , 37, 68, 256
- $x_2my_2(x, y)$ , 184
- $x_2py_2$ , 37, 68
- $x_2py_2(x, y)$ , 184
- $x\_div\_x_2py_2(x, y)$ , 185
- $x\_div\_x_2py_2(X, Y)$ , 208
- $\operatorname{xdsqrt1mx2}(X)$ , 241
- $\operatorname{xdsqrt1mx2}(x)$ , 198
- $\operatorname{xdsqrt1px2}(X)$ , 241
- $\operatorname{xdsqrt1px2}(x)$ , 198
- $\operatorname{xdsqrtx2m1}(X)$ , 241
- $\operatorname{xdsqrtx2m1}(x)$ , 198
- $Y_0(x)$ , 44
- $Y_1(x)$ , 44
- $Y_n(n, x)$ , 44
- Zahlen
  - Nachfolger, *siehe* succ
  - vertauschen, 35
  - Vorgänger, *siehe* pred
- Zahlenbereich
  - denormalisiert, 17
  - normalisiert, 17
- Zahlenformat
  - IEEE, 17
  - MPFR, 17
  - Staggered-Arithmetik, 11
- Zehnerpotenz
  - $\operatorname{exp10}$ , 37, 68, 92, 115
- Zeichenketten
  - $\operatorname{MpfcClass}(s, \operatorname{rnd}, \operatorname{prec})$ , 79
  - $\operatorname{MpfcIClass}(s, \operatorname{prec})$ , 98
  - $\operatorname{MpfiClass}(s, \operatorname{prec})$ , 49
  - $\operatorname{MpfrClass}(s, \operatorname{rnd}, \operatorname{prec})$ , 18
  - $\operatorname{string2Mpfc}(s, \operatorname{rnd}, \operatorname{prec})$ , 84
  - $\operatorname{string2MpfcI}(s, \operatorname{prec})$ , 102
  - $\operatorname{string2Mpfi}(s, \operatorname{prec})$ , 57
  - $\operatorname{string2Mpfr}(s, \operatorname{rnd}, \operatorname{prec})$ , 23
  - $\operatorname{to\_string}(X, \operatorname{prec})$ , 57
  - $\operatorname{to\_string}(x, \operatorname{rnd}, \operatorname{prec})$ , 23
  - $\operatorname{to\_string}(Z, \operatorname{prec})$ , 102
  - $\operatorname{to\_string}(z, \operatorname{rnd}, \operatorname{prec})$ , 83
- $\operatorname{zeta}(k)$ , 44
- $\operatorname{zeta}(x)$ , 44
- Zufallsintervalle, 66
- Zufallszahl
  - aus Intervall, 66
- Zufallszahlen, 35
- Zweierexponent, *siehe* expo
- Zweierpotenz
  - $\operatorname{exp2}$ , 37, 68, 92, 115