

Implementation of elementary functions of the *double*, *single* and *extended* class for MATLAB

Copyright © Marcel Leutenegger, 2003–2007
École Polytechnique Fédérale de Lausanne (EPFL)
Laboratoire d'Optique Biomédicale (LOB)
BM - Station 17
1015 Lausanne
Switzerland

June 27, 2007

This manual summarizes specific implementation details for the efficient computation of elementary MATLAB functions in case of complex floating point numbers. In general, the limited precision of the representation of floating point values leads to round-off errors in sums and differences involved in complex calculations, which potentially degrades the accuracy of the final result. This issue is briefly discussed and suitable implementations are proposed.

This library is free software; you can redistribute it and/or modify it under the terms of the GNU Lesser General Public License as published by the Free Software Foundation; version 2.1 of the License.

This library is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU Lesser General Public License for more details.

You should have received a copy of the GNU Lesser General Public License along with this library; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA.

Contents

1	Introduction	3
2	Basic functions	3
2.1	Exponential	3
2.2	Natural logarithm	4
2.3	Square root	4
3	Transcendental functions	5
3.1	Cosine	5
3.2	Hyperbolic cosine	5
3.3	Sine	6
3.4	Hyperbolic sine	6
3.5	Tangent	6
3.6	Hyperbolic tangent	6
4	Inverse transcendental functions	7
4.1	Inverse cosine	7
4.2	Inverse hyperbolic cosine	7
4.3	Inverse sine	7
4.4	Inverse hyperbolic sine	7
4.5	Inverse tangent	7
4.6	Inverse hyperbolic tangent	7

1 Introduction

Complex values z are represented as two floating point numbers representing the real and the imaginary part, respectively:

$$z = \Re(z) + i\Im(z) = a + ib \quad \text{where } i = \sqrt{-1} \text{ is the complex unit.}$$

The following abbreviations are used for function arguments and results:

f is a finite real value $\in (-\infty, +\infty)$,
 n is a finite, negative real value $\in (-\infty, -0]$,
 p is a finite, positive real value $\in [+0, +\infty)$,
 r is a real value $\in [-\infty, +\infty]$,
 u is an undefined value (NaN) and
 a is an affine value $\in \{\pm\infty, \pm\text{NaN}\}$.

The result of an undefined argument (NaN) is in general undefined. The rare exceptions are listed here-off. Although the sign does not matter, the MATLAB representation of +NaN is effectively the IEEE representation of -NaN.

2 Basic functions

2.1 Exponential

$$\exp(z) = e^a \cos(b) + ie^a \sin(b) \quad (1)$$

See also the comment on π in section 3 Transcendental functions.

Argument	Result
$-\infty$	0
f	r
$+\infty$	$+\infty$
$f+if$	$r+ir$
$f+ia$	$u+iu$
$-\infty+if$	0
$-\infty+ia$	0
$+\infty+if$	$a+ia$
$+\infty\pm i\infty$	$u+iu$

2.2 Natural logarithm

$$\log(z) = \frac{1}{2} \log(a^2 + b^2) + i \operatorname{atan} \frac{b}{a} \quad (2)$$

The logarithm "compresses" the dynamic range of the input value rendering it sensible to a misproportion between $|a|$ and $|b|$. Such a misproportion leads to a round-off error in the sum $a^2 + b^2$ affecting the result. Over the full complex plane $\mathbb{R} \times i\mathbb{R}$, the logarithm is computed as

$$\log(z) = \log|a| + \frac{1}{2} \log\left(1 + \frac{b^2}{a^2}\right) + i \operatorname{atan} \frac{b}{a} \quad \forall a^2 > 2^9 b^2 \quad (3)$$

$$\log(z) = \log|b| + \frac{1}{2} \log\left(1 + \frac{a^2}{b^2}\right) + i \operatorname{atan} \frac{b}{a} \quad \forall b^2 > 2^9 a^2 \quad (4)$$

The ratio is 2 to the power of the mantissa length difference (64bit – 53bit) less 2 rounding bits, hence $2^9 = 512$ at *double* precision. For *single* and *extended* precision the ratio is adapted accordingly.

Argument	Result
$-\infty$	$+\infty + i\pi$
n	$r + i\pi$
0	$-\infty$
p	r
$+\infty$	$+\infty$
$r + ir$	$r + if$

2.3 Square root

In principle, the square root could be computed as

$$\operatorname{sqrt}(z) = \sqrt{\frac{|z| + a}{2}} \pm i \sqrt{\frac{|z| - a}{2}} \quad \text{with} \quad |z| = \sqrt{a^2 + b^2} \quad (5)$$

But for $|a| \gg |b|$, the imaginary part would be truncated to zero. To maintain accuracy over the entire complex plane $\mathbb{R} \times i\mathbb{R}$, the square root is implemented as

$$\operatorname{sqrt}(z) = \sqrt{\frac{|z| + a}{2}} + i \frac{b}{2} \sqrt{\frac{2}{|z| + a}} \quad \forall a > 0 \quad (6)$$

$$\operatorname{sqrt}(z) = \frac{|b|}{2} \sqrt{\frac{2}{|z| - a}} \pm i \sqrt{\frac{|z| - a}{2}} \quad \forall a < 0 \quad (7)$$

The sign in equation (6) is automatically correct. In equation (7) it is adjusted accordingly.

Argument	Result
$-\infty$	$+i\infty$
n	ip
p	p
$+\infty$	$+\infty$
$f+if$	$f+if$
$f\pm i\infty$	$+\infty\pm i\infty$
$-\infty+if$	$+i\infty$
$+\infty+if$	$+\infty$
$\pm\infty\pm i\infty$	$u+iu$

3 Transcendental functions

Particular results of the transcendentals and their inverse functions will not be given explicitly hereoff. For affine arguments, they reply either by an undefined real or complex value or by

$$\text{function}(\infty + i\infty) = \lim_{z \rightarrow \infty + i\infty} \text{function}(z)$$

if this limit exists.

Comment on π The transcendental FPU instructions *fcos*, *fsin*, *fsincos* and *fpTan* reduce their arguments by computing the 2π -remainder. These FPU instructions use an internal representation of π with a 66bit mantissa regardless of the precision setting. Normal floating point operations never exceed the 64bit mantissa in *temporary real* format.

To minimize the round-off error, the arguments for the transcendental FPU functions are explicitly reduced by taking the 2π -remainder represented with a 53bit mantissa (*double* precision). In combination with a prescaling by a *temporary real* factor, the round-off error is reduced to less than 5.5E-20. This truncation and prescaling of the argument lowers the performance by 10% to 15%.

3.1 Cosine

$$\cos(z) = \frac{e^{iz} + e^{-iz}}{2} = \frac{1}{2}(e^{-b} + e^b)\cos(a) + \frac{i}{2}(e^{-b} - e^b)\sin(a) \quad (8)$$

3.2 Hyperbolic cosine

$$\cosh(z) = \cos(b - ia) = \frac{1}{2}(e^a + e^{-a})\cos(b) + \frac{i}{2}(e^a - e^{-a})\sin(b) \quad (9)$$

3.3 Sine

$$\sin(z) = \frac{e^{iz} - e^{-iz}}{2i} = \frac{1}{2}(e^b + e^{-b})\sin(a) + \frac{i}{2}(e^b - e^{-b})\cos(a) \quad (10)$$

3.4 Hyperbolic sine

$$\sinh(z) = i \sin(b - ia) = \frac{1}{2}(e^a - e^{-a})\cos(b) + \frac{i}{2}(e^a + e^{-a})\sin(b) \quad (11)$$

3.5 Tangent

$$\tan(z) = \frac{\sin(z)}{\cos(z)} = \frac{4e^{2b}\cos(a)\sin(a) + i(e^{2b} + 1)(e^{2b} - 1)}{((e^{2b} + 1)\cos(a))^2 + ((e^{2b} - 1)\sin(a))^2} \quad (12)$$

The tangent maintains accuracy up to an imaginary part $|b| \lesssim 373$. A larger $|b|$ causes either an overflow of e^{2b} to infinity or a truncation of the real part of the result to zero when saving as *double*. For a *single* (*extended*) result, the limit would be $|b| \lesssim 53$ (5678).

MATLAB R12 maintains accuracy for $|b| \lesssim 19$.

3.6 Hyperbolic tangent

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{(e^{2a} + 1)(e^{2a} - 1) + i4e^{2a}\cos(b)\sin(b)}{((e^{2a} + 1)\cos(b))^2 + ((e^{2a} - 1)\sin(b))^2} \quad (13)$$

The hyperbolic tangent maintains accuracy up to a real part $|a| \lesssim 373$. A larger $|a|$ causes either an overflow of e^{2a} to infinity or a truncation of the imaginary part of the result to zero when saving as *double*. For a *single* (*extended*) result, the limit would be $|b| \lesssim 53$ (5678).

MATLAB R12 maintains accuracy for $|a| \lesssim 20$.

4 Inverse transcendental functions

The current implementations of these functions use equation (5) for computing the square root and equation (2) for the logarithm. The number of addends involved would require numerous expressions for obtaining optimal accuracy in each case. Conditional branches to these expressions are unpredictable and would degrade the performance significantly.

The accuracy is improved by computing with the absolute values $|a|$ and $|b|$ and by adjusting the result according to the signs of a and b .

4.1 Inverse cosine

$$\operatorname{acos}(z) = -i \log \left(a + ib + \sqrt{a^2 - b^2 - 1 + 2iab} \right) \quad (14)$$

The real part is adjusted according to the signs of a and b . The imaginary part is always negative.

4.2 Inverse hyperbolic cosine

$$\operatorname{acosh}(z) = \log \left(a + ib + \sqrt{a^2 - b^2 - 1 + 2iab} \right) \quad (15)$$

The imaginary part is adjusted according to the signs of a and b . The real part is always positive.

4.3 Inverse sine

$$\operatorname{asin}(z) = -i \log \left(ia - b + \sqrt{1 - a^2 + b^2 - 2iab} \right) \quad (16)$$

4.4 Inverse hyperbolic sine

$$\operatorname{asinh}(z) = \log \left(a + ib + \sqrt{1 + a^2 - b^2 + 2iab} \right) \quad (17)$$

4.5 Inverse tangent

$$\operatorname{atan}(z) = -\frac{i}{2} \log \frac{1 - a^2 - b^2 + 2ia}{1 + a^2 + b^2 + 2ib} \quad (18)$$

4.6 Inverse hyperbolic tangent

$$\operatorname{atanh}(z) = -\frac{1}{2} \log \frac{1 - a^2 - b^2 - 2ib}{1 + a^2 + b^2 + 2ia} \quad (19)$$