

## 5L-1

## 初等関数の精度低下防止法

井阪 秀高 (神戸日本電気ソフトウェア㈱)

## 1. はじめに

初等関数の計算において、多倍長演算を使わずに、できるだけ精度を保って計算する注意事項をまとめる。さらに、対数ガンマ関数、べき乗関数、 $\sqrt{(x^2+y^2)}$  等については、精度を保つ計算方法を提案する。

## 2. 基本的方法

精度を保つ基本的考え方をまとめると次のようになる。

- ① 近似解  $y$  に微小量  $\Delta y$  を加算して解を更新する。
- ② 途中でビット落ちした値(閾値: 1/16=0.0625)を経由しない。
- ③ 減算で近似区間に還元するときは、入力引数から正確な上位ビット分の減算した後、下位ビット分の減算をする。
- ④  $x \neq 0$  で  $f(x) = 0$  になる関数は  $x$  を中心に最良近似式を作る。
- ⑤ 解の指数部、仮数部別々に計算し、最後に合成する。

また、各関数計算時に使う近似式に対し、16進浮動小数点内部表現計算として共通な方法を以下にまとめる。

## (1) 多項式の場合

$x * (1 + (y * (a_1 + \dots)))$  は、 $x + \{x * y * (a_1 + \dots)\}$  として計算。  
 $x * \{a_0 + (y * (a_1 + \dots))\}$  で  $\{ \}$  内の仮数部がビット落ちするときは、各係数を定数で割り  $\{ \}$  内を 0.5~1.0 になるようにしてから多項式計算をし、最後にその定数を乗じる。

(2)  $x * \{ (\text{多項式 I}) / (\text{多項式 II}) \}$  の場合

$\{ \}$  内の仮数部がビット落ちするときは、 $\{ (\text{多項式 I}) * x \} / (\text{多項式 II})$  と計算するか、近似区間で解が  $x$  に近いとき  $x + x * (\text{多項式 III}) / (\text{多項式 II})$  と式を変形して計算するか、 $|x| < 0.5$  のとき  $(x+x) * \{ (\text{多項式 I}) / \{ 2 * (\text{多項式 II}) \} \}$  と計算する。但し、このとき多項式 I、多項式 II、多項式 III とも  $x$  のとりうる値の範囲内で仮数部ができるだけビット落ちしないように多項式係数に定数を乗じて近似式を作る。

## 3. 三角関数 (sin(x), cos(x), tan(x))

(1)  $y = |x| * 2 / \pi$  とする。

$\sin(x) : n = \{ (y \text{ の整数部}) + 1 \} \wedge 'F \cdot FE'$   
 $\cos(x) : n = (y \text{ の整数部}) \vee '0 \cdot 01'$

ここで、 $(y-n)$  で  $[-1, 1]$  区間に還元したり  $(|x| - n * \pi / 2)$  で  $[-\pi/2, \pi/2]$  区間に還元すると精度が落ちるので、次のように分割減算し  $[-\pi/2, \pi/2]$  区間に還元する。  
 $r = (|x| - n * a) - n * b$ ,  $\cos(x)$  計算時:  $r = -r$   
 $n$ : 整数部,  $r$ : 小数部 \*  $\pi / 2$

$a : \pi / 2$  の仮数部上位半ビット分以下位半ビットを 0 にしたもの

$b : \pi / 2 - a$  ( $a$  の下位半ビット以下のビット列)  
 $\sin(r)$  を計算し、 $n$  の最下位から 2 ビット目が 1 のとき正負逆転する。

$\sin$  は  $x$  が負のときさらに正負逆転する。

(2)  $y = |x| * 4 / \pi$  とする。

このときも (1) と同様に  $[-\pi/4, \pi/4]$  区間に還元する。

$n = \{ (y \text{ の整数部}) + 1 \} \wedge 'F \cdot FE'$

$r = (|x| - n * a) - n * b$

$n$ : 整数部,  $r$ : 小数部 \*  $\pi / 4$

$a : \pi / 4$  の仮数部上位半ビット分以下位半ビットを 0 にしたもの

$b : \pi / 4 - a$  ( $a$  の下位半ビット以下のビット列)  
 $n$  の最下位から 2 ビット目を  $i$

$n$  の最下位から 3 ビット目を  $j$  とする。

$\sin(x)$

$i = 0 : \sin(r), 1 : \cos(r)$  を計算

$j \neq (x \text{ の符号ビット}) : \text{正負逆}$

$\cos(x)$

$i = 0 : \cos(r), 1 : \sin(r)$  を計算

$i \neq j : \text{正負逆}$

$\tan(x)$

$i = 0 : \tan(r), 1 : -1 / \tan(r)$  を計算

$x < 0 : \text{正負逆}$

4.  $x \neq 0$  で  $f(x) = 0$  になる関数(1)  $\log(\Gamma(x))$ 

$x \approx 1, 2$  で  $f(x) \approx 0$  になり桁落ちしやすい。

$x \approx 1$  のとき  $y = x - 1$  として

$\log(\Gamma(x)) \approx -\gamma * y + \sum_{n=2}^{\infty} \{ \zeta(n) * (-y)^n / n \}$

$x \approx 2$  のとき  $y = x - 2$  として

$\log(\Gamma(x)) \approx (1 - \gamma) * y +$

$\sum_{n=2}^{\infty} \{ (\zeta(n) - 1) * (-y)^n / n \}$

を最良近似式化した多項式で求める。

( $\gamma$ : Euler 定数,  $\zeta(n)$ : Riemann のゼータ関数)

(2)  $\text{acosh}(x)$ 

$x \approx 1 (x \geq 1)$  のとき  $f(x) \approx 0$  になり桁落ちしやすい。

$z = \sqrt{(x-1)}, Y = -(x-1)$  として

$\text{acosh}(x) \approx z + z * [ (\sqrt{2}-1) + \sqrt{2} *$

$\sum_{n=1}^{\infty} \{ (2n-1)!! * (Y)^n / \{ (2n)!! * (2n+1) * 2^n \} \}$

] を最良近似式化した多項式で求める。

(3)  $\text{acos}(x)$ 

$x \approx 1 (x \leq 1)$  のとき  $f(x) \approx 0$  になり桁落ちしやすい。

The Algorithms for Keeping Precision of Elementary Functions

Hidetaka Isaka

NEC Software Kobe, Ltd.

$z = \sqrt{1-x}$ ,  $Y = 1-x$ として  
 $\text{acos}(x) \approx z + z * [(\sqrt{2}-1) + \sqrt{2} * \sum_{n=1}^{\infty} [(2n-1)!! * (Y)^n / \{(2n)!! * (2n+1) * 2^n\}]]$   
 ]を最良近似式化した多項式で求める。

#### (4) $\log_2(x)$

$x \approx 1$ のとき  $f(x) \approx 0$  になり桁落ちしやすい。  
 $1/\sqrt{2} < x < \sqrt{2}$  のとき  $y = (x-1)/(x+1)$  として  
 $y$  に対する最良近似式  $f(y)$  を用いる。  
 または、 $x \approx 2^m$ , ( $0.5 \leq m < 1$ ) なる  $m$  を取り出し  
 $m < 1/\sqrt{2}$  のとき:  $b = 1$ ,  $z = 2^m$   
 $m \geq 1/\sqrt{2}$  のとき:  $b = 0$ ,  $z = m$   
 $y = (z-1)/(z+1)$  として  
 $\log_2(x) = (e-b) + f(y)$

### 5. 公式どうり計算すると桁落ちする

#### (1) $x^r$

$\exp(y * \log(x))$  と計算すると  $x$  が大きいとき桁落ちするので、解の指数部と仮数部を別々に計算する。  
 $p$  進内部表現とする。

$x \approx p^e * m$  ( $e$ : 指数部,  $m$ : 仮数部 ( $1/p \leq m < 1$ ), を分離)

if ( $m < 1/\sqrt{p}$ ) then  $m = m * p$ ,  $e = e - 1$

$y \Rightarrow i + d$  ( $i$ : 整数部,  $d$ : 小数部, を分離)

$g = \log_p(m) * y + e * d$

$g \Rightarrow j + b$  ( $j$ : 整数部,  $b$ : 小数部 ( $-1 \leq b < 0$ ), を分離)

$x^r = p^{(e * i + j)} * p^b (e * i + j$ : 指数部,  $p^b$ : 仮数部)

#### (2) $\sqrt{x^2 + y^2}$

$x^2, y^2$  をそのまま計算するとオーバーフローしたり精度が低下することがある。

$|y| \leq |x|$  として  $|x|/n * \sqrt{\{n^2 + n^2 * (y/x)^2\}}$ ,  
 ( $n=1, 0.5, \dots$ ) としたり  $x$  の指数部の値を  $x$ ,  $y$  の指数部から減じて  $\sqrt{x^2 + y^2}$  を計算し、減じた値を解の指数部に加算するのも精度が低下することがある。そこで、近似解  $|x|$  に微小量  $\Delta x$  を加えて解を更新する。

$|y| \leq |x|$  とする。

$r = (y/x)^2$

$t = \sum_{n=0}^{\infty} \frac{r^n}{2^n}$ , ( $0 \leq r \leq 1$ ) を最良近似式化した

$t = r * f(r) / g(r)$  により  $t$  を計算する。

$\sqrt{x^2 + y^2} \approx |x| + |x| * t$

また、精度を上げるためには、

$\sqrt{x^2 + y^2} \approx |x| + |x| * (t * t + r) / (t + t + 2)$  とする。

応用として  $\sqrt{x}$ , ( $0.5 \leq x \leq 1$ ) は、 $r = 2 * (x - 0.5)$

として  $t$  を求め、 $p = 1/\sqrt{2}$ ,  $\sqrt{x} = p + p * t$  とする。

$\text{asinh}(x)$ , ( $x \geq 1$ ) は、 $r = 1/x^2$  として  $t$  を求め、

$\text{asinh}(x) = \log(2x + x * t)$ , ( $x * r = 1/x$ ) とする。

#### (3) $\sqrt{x^2 - y^2}$ , ( $|y| \leq |x|$ )

$\sqrt{x^2 + y^2}$  と同様に  $r = (y/x)^2$  とし

$t = \sum_{n=0}^{\infty} \frac{-r^n}{2^n}$ , ( $0 \leq r \leq a$ ,  $a < 1$ ) を最良近似式化した

$t = r * h(r) / i(r)$  により  $t$  を計算する。

$\sqrt{x^2 - y^2} \approx z = |x| + |x| * t$

また、精度を上げるためには、

$\sqrt{x^2 - y^2} \approx |x| + |x| * (t * t - r) / (t + t + 2)$  とする  
 $\text{acosh}(x)$ , ( $x \geq 1/\sqrt{a}$ ) は、 $\text{asinh}(x)$  と同様に計算できる

### 6. その他の関数

#### (1) $e^x - 1$

$x \approx 0$  で桁落ちするので、最良近似式とデータテーブルを利用する。また、近似区間への還元値 ( $z$ ) を得るときも精度を保つため上位、下位ビット分けて減算する。

$y = x * \log_2(e)$

$|y| < 1$  のとき

$|y|$  の小数点以下  $n$  ビットを取り出し  $i$  とする。

$i = (i + 1) / 2$  (整数演算),  $x < 0$  のとき  $i = -i$

$p = i / 2^{n-1}$ ,  $z = (x - p * a) - p * b$

$e^x - 1 = u(i) + t(i) * (e^x - 1)$

$a$ :  $\log_e(2)$  の仮数部上位半ビット分以下位半ビットを 0 にしたもの

$b$ :  $\log_e(2) - a$  ( $a$  の下位半ビット以下のビット列)

$u(i)$ :  $2^{i/(2^{n-1})} - 1$  のデータテーブル

$t(i)$ :  $2^{i/(2^{n-1})}$  のデータテーブル

$e^x - 1$ :  $z$  に対する最良近似式の値 ( $|z| \leq 1/2^n$ )

$|y| \geq 1$  のとき

$y \Rightarrow n + d$  ( $n$ : 整数部,  $d$ : 小数部 ( $-1 < d \leq 0$ ), を分離)

$|d|$  の小数点下  $n$  ビットを  $i$  とする。

$p = n - i / 2^n$ ,  $z = (x - p * a) - p * b$

$e^x - 1 = 2^n * (v(i) * e^z) - 1$  ( $n$ : 整数部,  $i$ : 仮数部)

$v(i)$ :  $2^{-i/(2^{n-1})}$  のデータテーブル

$e^z$ :  $z$  に対する最良近似式の値 ( $-1/2^n < z \leq 0$ )

#### (2) $\text{atan}(x)$

$x \geq 0$  の範囲で計算し、 $x < 0$  のとき解を負にする。

$x \geq \tan((2n-1)*\pi/4n)$  ( $n$ : 区間数 - 1) のとき

$\text{atan}(x) = \pi/2 - \text{atan}(1/x)$

else

$\tan((2i-1)*\pi/4n) < x \leq \tan((2i+1)*\pi/4n)$

( $i = 0, \dots, n-1$ ), なる  $i$  を見つける。

$\text{atan}(x) = i * \pi/2n + \text{atan}\{(x-s)/(s*x+1)\}$

$s = \tan(i * \pi/2n)$

ここで、16 進内部表現のとき  $\{(2i-1)*\pi/4n < 1\} \wedge \{i * \pi/2n \geq 1\}$  となる区間ができるとその精度が低く、このようにならない  $n$  を決める ( $n=2, 5, 7, 8, 10, 11, 13, \dots$ )。

### 7. あとがき

初等関数を効率よく、精度低下を防ぎながら計算する方法をまとめた。これらの算法は、NECのSXシリーズ等で既に実現し、これらにより精度向上がなされたことを確認している。

最後に、組込み関数実現に関し色々ご協力くださったNEC基本ソフトウェア開発部の片山部長代理、鈴木様に感謝します。

### 8. 参考文献

山内, 宇野, 一松: 電子計算機のための数値計算法 II, 培風館, (1972)

井阪他: ベクトル計算機に於ける三角関数計算法, 情報処理学会第34回全国大会, (1987)