

多倍長計算手法

平成25年度第1四半期

目次

1. はじめに
2. グリーン関数計算
3. 条件数に関して
4. 表現可能な数値範囲の影響とデータ形式
5. DD形式に関して

1. はじめに

平成25年3月8日に行われました
第3回多倍長計算フォーラムで

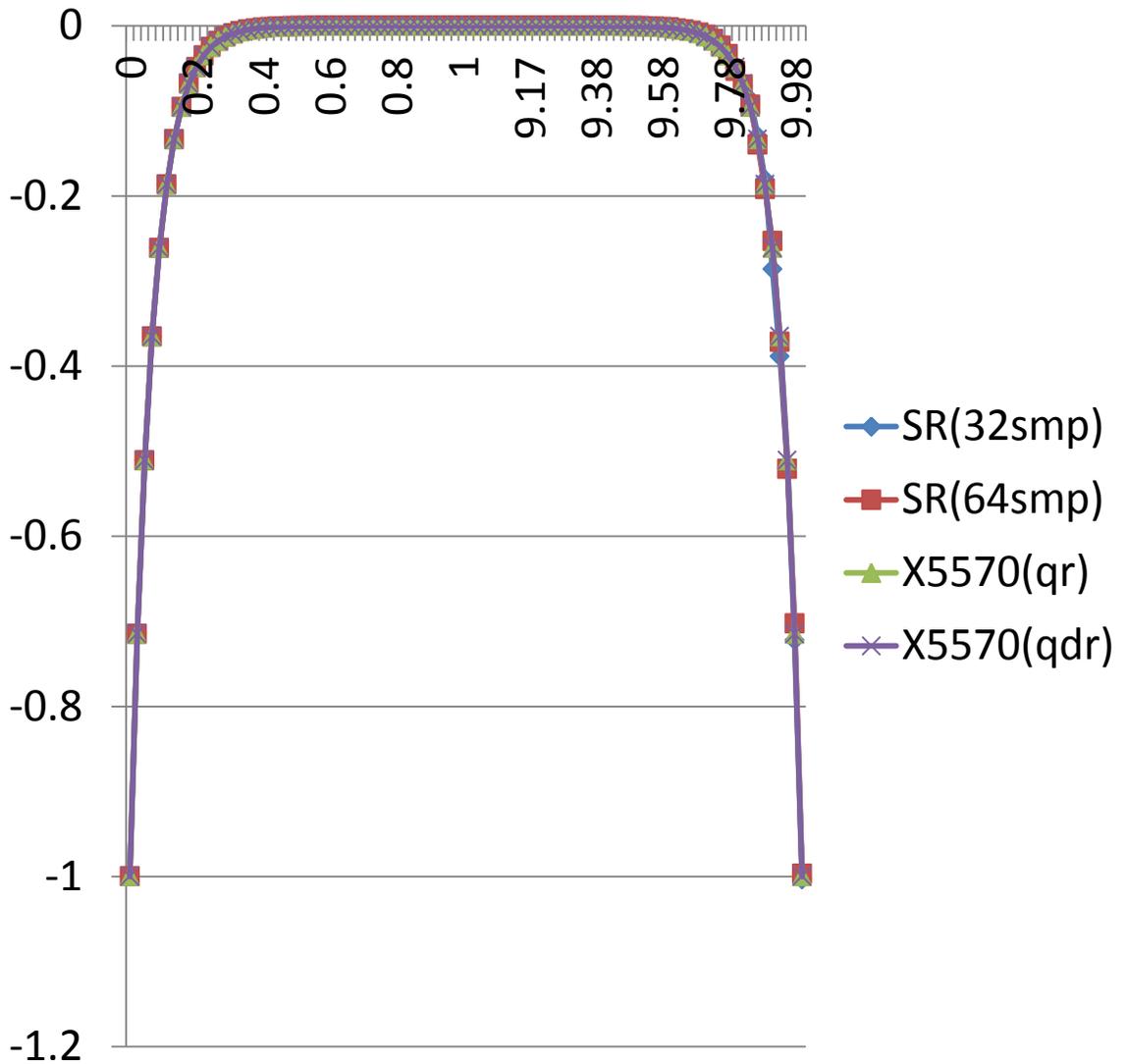
量子モンテカルロ法による物性
スペクトル計算と4倍精度
(KEK物構研 岩野氏)

において、紹介されました4倍精度演算でも
正しい結果が得られない場合がある件の精度
を調査しました。

**結果の検証は静的かつ規則的な
パスに対して $G(\beta) = -1$ でチェックします。
問題の発端は $\beta = 10, u = 5, l = 448$ の
4倍精度演算の結果です。**

この問題を異なるアーキテクチャーの
SR16000/M1とX5570で実行した結果、
若干結果の精度に問題がみられました。
(次ページに結果を掲載)

Beta=10,u=5,l=448 4倍精度精度比較



DT	SR(32smp)	SR(64smp)	X5570(qr)	X5570(qdr)
0	-1.000042	-0.999995	-1.000004	-0.999991
10	-1.00375	-0.996719	-0.999997	-1.000001

DT=0,10 での厳密解は-1.0

2. グリーン関数計算

$$R_{\sigma}(\tau, x) = e^{-\Delta H(\tau)} e^{-\Delta H(\tau-\Delta)} \dots e^{-\Delta H(\Delta)} = B_m B_{m-1} \dots B_1$$

$$\tau = m\Delta, \beta = L\Delta \quad (0 \leq \tau \leq \beta)$$

$$R_{\sigma}(\beta, x) = e^{-\Delta H(\beta)} e^{-\Delta H(\beta-\Delta)} \dots e^{-\Delta H(\Delta)} = B_L B_{L-1} \dots B_1$$

$$\text{グリーン関数} \equiv R_{\sigma}(\tau, x) [I + R_{\sigma}(\beta, x)]^{-1}$$

$$\text{行列表現だと } G_m = (B_m B_{m-1} \dots B_1 B_0) [I + B_L B_{L-1} \dots B_1]^{-1}$$

$m = 0, B_0 = I, B_m B_{m-1} \dots B_1$ の行列積をそのまま計算すると精度が良くありません。

そこで、精度良く計算する方法には以下の2つの方式があります。

(1) QDR方式

$$B_1 = Q_1 D_1 R_1$$

$$B_2 Q_1 = Q_2 D_2 R_2 \quad B_2 B_1 = Q_2 D_2 R_2 D_1 R_1 = Q_2 D_2 D_1 R_2 R_1 = Q_2 D_2 R_2$$

:

:

$$B_m Q_{m-1} = Q_m D_m R_m$$

$$B_m B_{m-1} \dots B_1 = QDR$$

Q, 直交行列, D対角行列, R上三角行列 (対角要素はすべて1)

(2) QR方式

$$B_1 = Q_1 R_1$$

$$B_m Q_{m-1} = Q_m R_m$$

$$B_m B_{m-1} \dots B_1 = Q_m (R_m \dots R_1) = QR$$

Q, 直交行列, R上三角行列

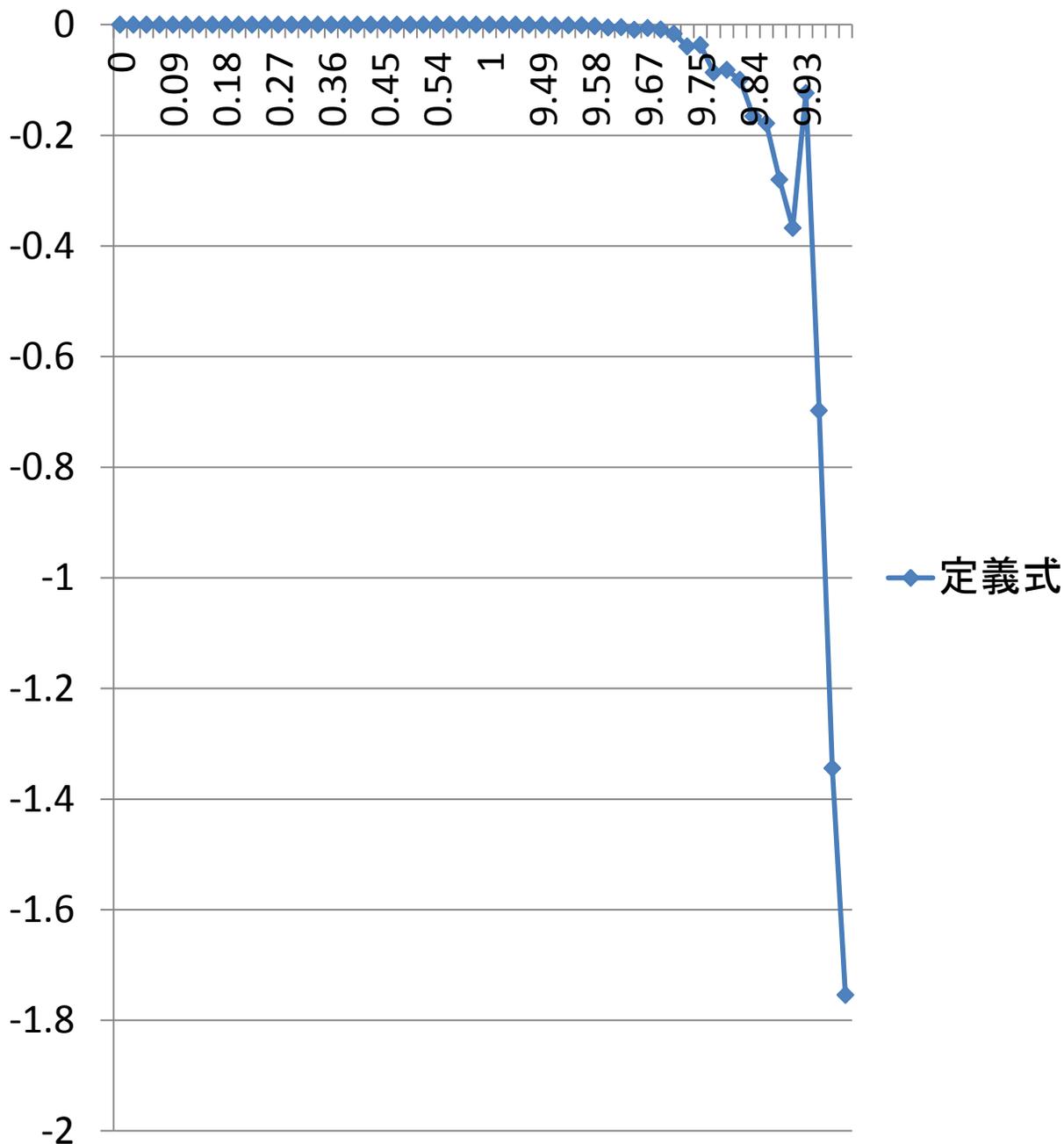
どちらの場合も精度上 $[I + B_L B_{L-1} \dots B_1]^{-1}$ の計算が問題となります。

3つの計算方法による結果を次ページ以降に記しました。

Green 関数

$l=448, \beta=10, u=5$

定義式

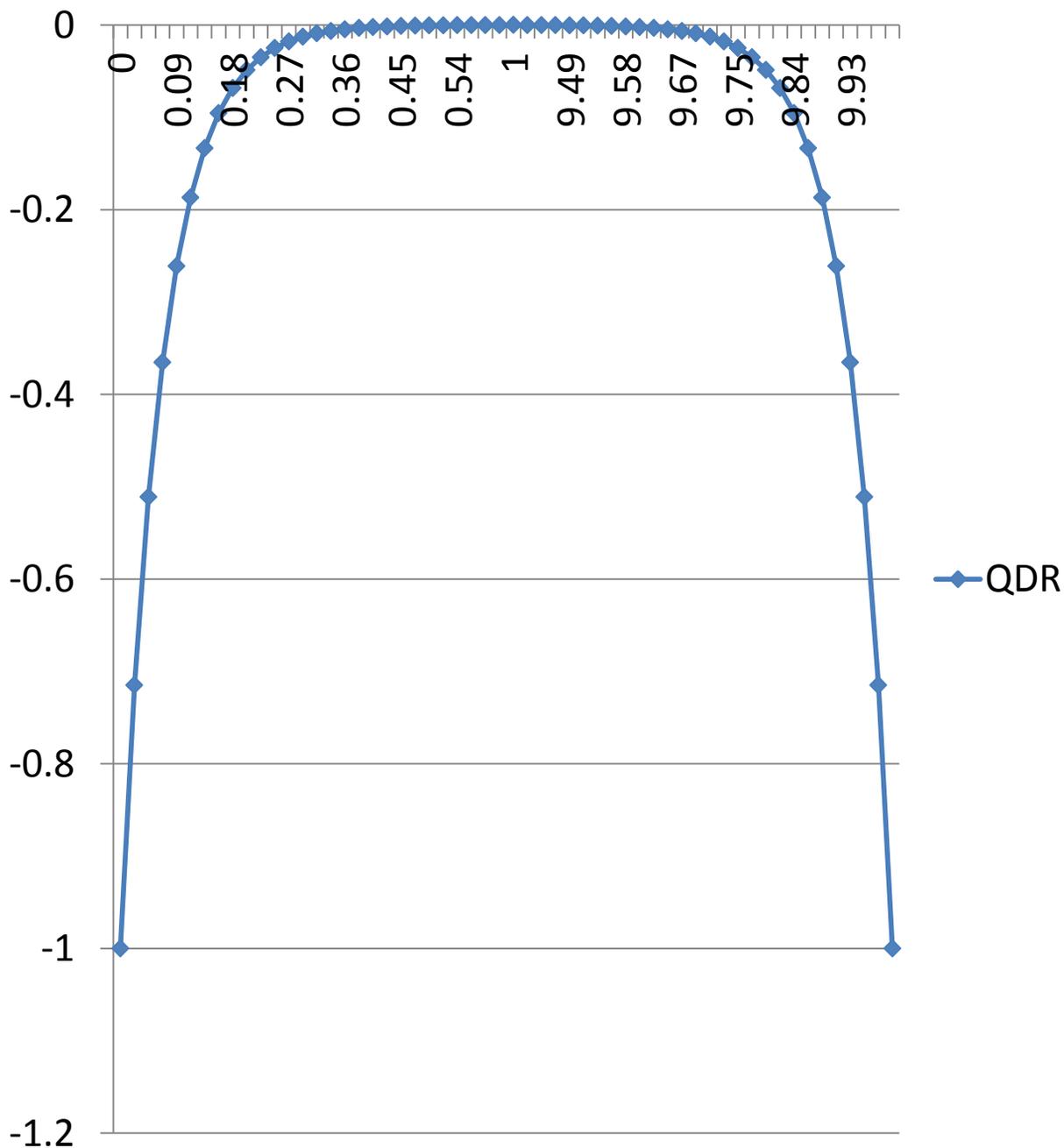


演算精度 4倍精度

Green 関数

$l=448, \beta=10, u=5$

QDR

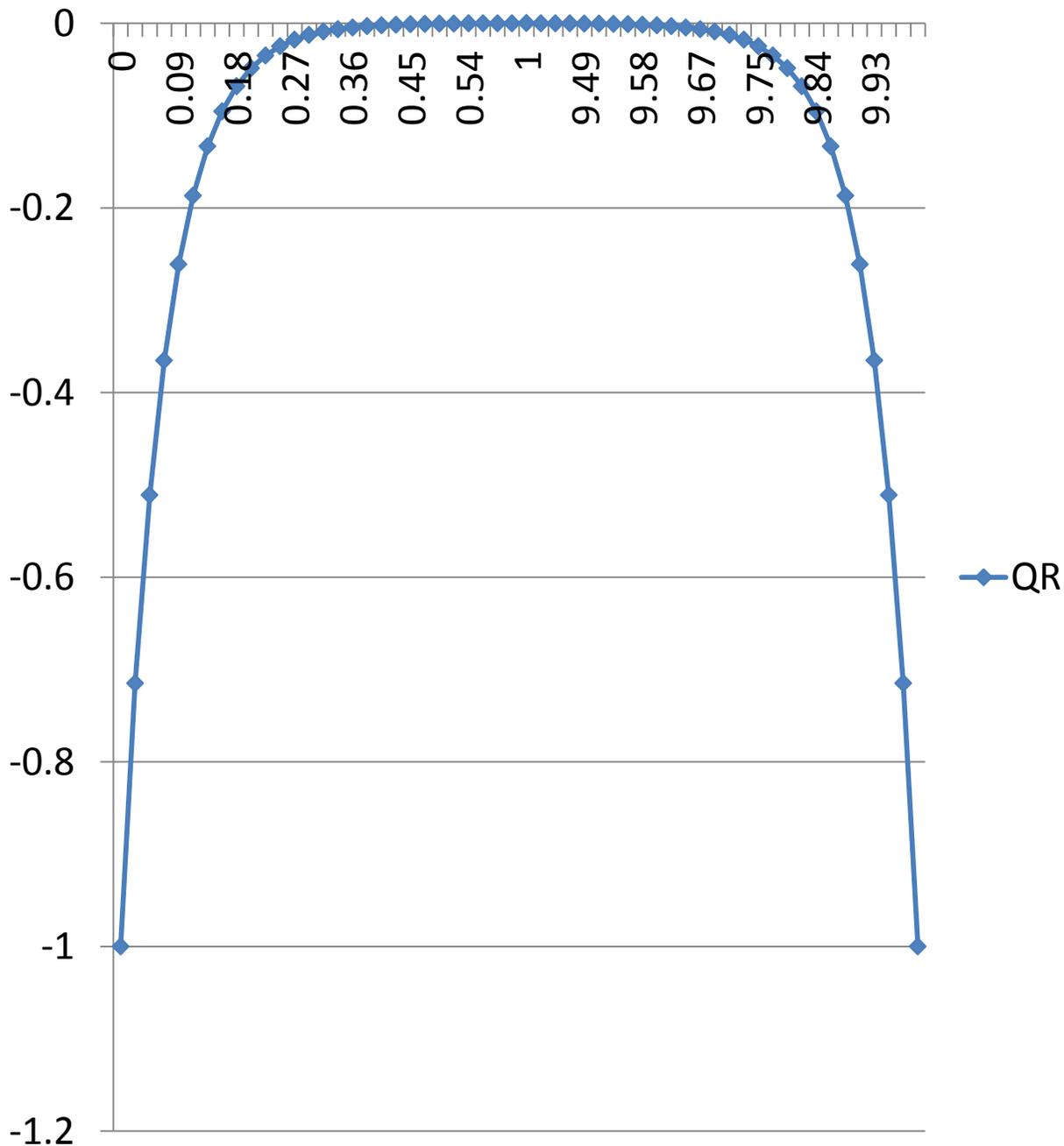


演算精度 4倍精度

Green 関数

$l=448, \beta=10, u=5$

QR



演算精度 4倍精度

3.条件数に関して

今回扱ったケースでは、設定したパラメータから結果の精度が見積もり易くなっています。

パラメータ n, L, β, u のうち $n = 100, L = 448$ と固定しています。

G_m : グリーン関数。

$P = e^{\sqrt{dt \times u \times L}} = e^{\sqrt{\beta \times u \times L}}$, $E(m) = |\text{Tr}(G_m)|$ とすると,

プログラム実行中に現れる数値の最大値 P^2 ,

0以外の絶対値最小値 $1/P^2$ 。

行列やベクトル要素に現れる数値の最大値 P ,

0以外の絶対値最小値 $1/P$ 。

$E(0) = 1, E(L) = 1, m = \frac{L}{2}$ で $E(\frac{L}{2}) = \frac{1}{\sqrt{P}}$ で最小となるので,

条件数は \sqrt{P} となり、演算に必要な最小ビット数は

$\log_2 \sqrt{P} = \frac{1}{2} \log_2 P$ となります。

理論式と実際の実行結果との対応は次ページ以降に示しました。

演算に必要な最小ビット数(L=448)						
β	u=5	u=6	u=7	u=8	u=9	u=10
10	108	118	128	137	145	153
20	153	167	181	193	205	216
演算精度の有効ビット数						
精度	ieee2008 ビット数	dd形式 ビット数				
4倍精度	113	106				
5倍精度	145					
6倍精度	177	159				
7倍精度	209					
8倍精度	241	212				

**DD形式,ieee754 – 2008形式の8倍精度までの
実行結果は演算に必要な最小ビット数と各演算精度
の有効ビット数の関係とよく一致しています。**

dd形式の多倍長計算での実行結果

ここでOKは10桁以上一致,NGは最初の2桁以下しか一致していない
(全桁不一致も含む) 場合を示しています.

$\beta = 10$

6倍精度 $u = 6, 7, 8, 9$ OK, $u = 10$ 7桁一致と微妙

8倍精度 $u = 6, 7, 8, 9, 10$ OK

$\beta = 20$

6倍精度 $u = 5$ 5桁一致と微妙, $u = 6, 7, 8, 9, 10$ NG

8倍精度 $u = 6, 7, 8$ OK, $u = 9$ 6桁一致と微妙, $u = 10$ NG

10倍精度 $u = 9, 10$ OK

6倍精度はieee754 – 2008形式の5倍精度と6倍精度のほぼ中間の有効ビット数

8倍精度はieee754 – 2008形式の7倍精度と有効ビット数は等しい

ieee754 – 2008形式の多倍長計算での実行結果

$\beta = 10$

5倍精度 $u = 6, 7, 8$ OK, $u = 9$ 4桁一致と微妙, $u = 10$ NG

6倍精度 $u = 6, 7, 8, 9, 10$ OK

7倍精度 $u = 6, 7, 8, 9, 10$ OK

8倍精度 $u = 6, 7, 8, 9, 10$ OK

$\beta = 20$

5倍精度 $u = 5, 6, 7, 8, 9, 10$ NG

6倍精度 $u = 5, 6$ OK, $u = 7, 8, 9, 10$ NG

7倍精度 $u = 5, 6, 7, 8$ OK, $u = 9$ 5桁一致と微妙, $u = 10$ NG

8倍精度 $u = 6, 7, 8, 9, 10$ OK

4.表現可能な数値範囲の影響とデータ形式

多倍長計算では,以下の2種類のデータ形式があります。

(1) DD形式

SR16000/M1およびBG/Qシステムなどで使用されている多倍長変数を複数個の倍精度変数の和であらわしたものの。

2個で4倍精度変数,3個で6倍精度変数,
4個で8倍精度変数,5個で10倍精度変数など

(2) ieee形式

ieee754 – 2008の4倍精度データ形式の仮数部部分を拡張したものの。

P倍精度変数 : 符号部1ビット,指数部15ビット,
仮数部 $32 \times P - 16$ ビット

数値表現可能な範囲の制限でDD形式の演算は、オーバーフローや、絶対値の小さい数の精度の低下の影響があります。

u=10,L=448

Q: プログラム実行中に現れる最大数の10進桁数
 P: 行列, ベクトル要素の最大数の10進桁数
 R: 演算に必要な最小ビット数

beta	理論式			16倍精度の実測			SR16000 dd形式
	Q	P	R	Q	P	R	
30	318.4	159.2	264.5	322.8	161.4	264.9	overflow
29	313.1	156.5	260	317.2	158.6	260.4	overflow
28	307.6	153.8	255.5	311.5	155.8	255.8	overflow
27	302.1	151	250.9	305.7	152.8	251.2	NG
26	296.4	148.2	246.2	299.8	149.9	246.4	NG
25	290.7	145.3	241.4	293.8	146.9	241.6	NG
24	284.8	142.4	236.5	287.7	143.8	236.7	NG
23	278.8	139.4	231.6	281.5	140.8	231.6	NG
22	272.7	136.3	226.5	275.1	137.6	226.5	NG
21	266.4	133.2	221.3	268.6	134.3	221.2	5桁一致

ieee754-2008 8倍精度

β	一致10進桁数
21	10 OK
22	9 微妙
23	7 微妙
24	3 微妙

1/Qの計算において, dd形式の0以外の絶対値最小数は
 $2^{-1074} \cong 10^{-323.3} \cong 0.494 \times 10^{-323}$ のため,
 $1074 - 4 \times R$ ビットの精度しか持ち得ません。

aの平方根計算には以下の2つの方法があります。

(1) 逆数平方根 $1/\sqrt{a}$ を求め $\sqrt{a} = a \times (1/\sqrt{a})$ より計算する。

$$x_0 = 1/\sqrt{a}(\text{倍精度}) \text{ から } x_{n+1} = x_n + \frac{x_n(1-ax_n^2)}{2} \text{ で計算。}$$

$$x_n \rightarrow 1/\sqrt{a}.$$

除算がないという利点はあるが、aがおおきな数値の場合、 x_n^2 の計算結果の精度が数値表現の制限の影響をうけ悪くなる。

(2) $x_0 = \sqrt{a}(\text{倍精度})$ から $x_{n+1} = \frac{1}{2}(x_n + \frac{a}{x_n})$ で計算。

$$x_n \rightarrow \sqrt{a}.$$

除算があるが、aがおおきな数値の場合、数値表現の制限の影響はうけない。

平方根計算の方式を(1)から(2)に変更すると、SR16000/M1 DD 形式10倍精度では、 $\beta = 21, 22, 23, 24, 25, 26, 27$ $u=10$ で結果は10進10桁まで一致と精度改善ができました。

$\beta = 28, 29, 30$ $u=10$ はオーバーフローで結果は変わりません。

5.dd形式に関して

a_i : 倍精度変数, ieee754 – 2008 4倍精度変数 仮数部, 指数部が1つ

これは, 2つつなげた場合と3つ以上つなげた場合では異なる処理が必要になる.

$$c_1 = a_1 + a_2$$

$$p_1 = c_1 - a_1$$

$$c_2 = (a_1 - (c_1 - p_1)) + (a_2 - p_1)$$

なら $|c_1| \geq |c_2|$ となる.

$n = 3$

$$t_1, t_2, t_3 \text{ が得られると, } (t_1, t_2, t_3) \rightarrow (t_4, t_5, t_6) \rightarrow (c_1, c_2, c_3)$$

で $|c_1| \geq |c_2| \geq |c_3|$ となる.

$n = 4$

t_1, t_2, t_3, t_4 が得られると,

$$(t_1, t_2, t_3, t_4) \rightarrow (t_5, t_6, t_7, t_8) \rightarrow (t_9, t_{10}, t_{11}, t_{12}) \rightarrow (c_1, c_2, c_3, c_4)$$

で $|c_1| \geq |c_2| \geq |c_3| \geq |c_4|$ となる.

8倍精度の注意事項

DD形式の4倍精度を2つつなげた場合 (DQ形式)

$c_1 = a_1 + a_2, c_2 = a_3 + a_4$ c_i : 4倍精度変数, a_i : 倍精度変数

$|c_1| \geq |c_2|, |a_1| \geq |a_3| \geq |a_4|$ ではあるが,

$|a_2|$ と $|a_3|, |a_4|$ の大小関係は定まらない.

SR16000/M1での8倍精度乗算に関して

倍精度変数を4つつなげたQD形式

4倍精度変数を2つつなげたDQ形式

$$a = a_0 + a_1 + a_2 + a_3, b = b_0 + b_1 + b_2 + b_3$$

$$a = c_0 + c_1, b = d_0 + d_1$$

$$c_0 = a_0 + a_1, c_1 = a_2 + a_3, d_0 = b_0 + b_1, d_1 = b_2 + b_3$$

とすると,

QD形式

$$\begin{aligned} a \times b &= a_0 b_0 \\ &+ a_0 b_1 + a_1 b_0 \\ &+ a_0 b_2 + a_1 b_1 + a_2 b_0 \\ &+ a_0 b_3 + a_1 b_2 + a_2 b_1 + a_3 b_0 \\ &+ a_1 b_3 + a_2 b_2 + a_3 b_1 \end{aligned}$$

DQ形式

$$\begin{aligned} a \times b &= c_0 d_0 + c_0 d_1 + c_1 d_0 \\ &= a_0 b_0 + a_0 b_1 + a_1 b_0 + a_1 b_1 \\ &+ a_0 b_2 + a_0 b_3 + a_1 b_2 + a_1 b_3 \\ &+ a_2 b_0 + a_2 b_1 + a_3 b_0 + a_3 b_1 \end{aligned}$$

DQ形式はQD形式より $a_2 b_2$ ($O(\epsilon^4)$, $\epsilon = 2^{-53}$) だけ精度が悪くなります。

$a_3 = b_3 = 0$ の場合、もつとも精度が悪くなり,

倍精度変数を3つつなげたTD (6倍精度) の乗算結果

$$\begin{aligned} a \times b &= a_0 b_0 \\ &+ a_0 b_1 + a_1 b_0 \\ &+ a_0 b_2 + a_1 b_1 + a_2 b_0 \\ &+ a_1 b_2 + a_2 b_1 \end{aligned}$$

と一致する事になります。