

多倍長計算手法

平成26年度第4四半期

目次

1. はじめに
2. massless 計算
 - 2.1 vtx-1
 - 2.2 vtx-2
 - 2.3 box
3. sinc求積法によるHadamard有限部分計算
 - 3.1 SE,DEの計算方法
 - 3.2 ε -算法
 - 3.3 計算結果精度比較
4. 多倍長減算桁落ちメモ
 - 4.1 3倍精度
 - 4.2 5倍精度
 - 4.3 6倍精度
 - 4.4 7倍精度
5. 仮数部のビット数と乗算のメモ
6. DD,D (long double),DQのインライン化メモ
 - 6.1 FORTRAN
 - 6.2 C

1.はじめに

ファインマンループ積分のmassless計算に関してその発端から数学的根拠などをまとめました。

これに関連して、端点と内部に同時に特異点を持つ積分に対してHadamardの有限部分を sinc 求積法と今までの ε - 算法で実行した時の結果の精度の比較をしています。

また、非対称行列反復計算の調査で作成したメモをまとめて記述しています。

- (1) 整数演算での減算の桁落ち計算
- (2) 仮数部のビット毎の乗算

またDD, D (longdouble), DQソースのインライン化, 並列化のためのメモをまとめています。

2.massless計算

2.1 vtx-1

inf ra vtx の計算において

$$\int_0^1 \int_0^{1-x} \frac{1}{D} dy dx \quad D = -sxy + (x+y)^2 m^2 + (1-x-y)\lambda^2$$

を $s = 500^2$, $m = 0.0005$ で計算する場合, $D \rightarrow D - i\varepsilon$ として計算しますが λ の値が小さくなると, 反復回数が多くなり計算に時間がかかったり, 計算できない ($\lambda = 10^{-200}$ など倍精度, DD形式の4倍精度など) 事が発生するので, 以下の式で計算する事があります。この場合は

$$\int_0^1 \int_0^{1-x} \frac{1}{D^{1-\eta}} dy dx \quad D = -sxy + (x+y)^2 m^2 \quad D \rightarrow D - i\varepsilon, \eta = \frac{1}{\log\left(\frac{1}{\lambda^2}\right)}$$

λ の値が充分小さくないと結果の精度が良くない事があります。

実測値と解析式の相対誤差

λ	実数部	虚数部
10^{-15}	6.23×10^{-3}	1.38×10^{-2}
10^{-150}	6.30×10^{-5}	1.49×10^{-4}
10^{-1500}	6.68×10^{-7}	1.59×10^{-6}
10^{-2500}	2.26×10^{-7}	5.58×10^{-7}
10^{-15000}	6.42×10^{-9}	1.61×10^{-8}

数学的には以下ようになります。

$$I = \int_0^1 \int_0^{1-x} \frac{1}{-sxy + (x+y)^2 m^2 + (1-x-y)\lambda^2} dy dx \quad s < 0$$

$$= \int_0^1 \int_0^1 \frac{p}{-sp^2q(1-q) + p^2 m^2 + (1-p)\lambda^2} dp dq = A + B$$

$$A = \frac{1}{2} \int_0^1 \frac{\log \frac{-sq(1-q) + m^2}{\lambda^2}}{-sq(1-q) + m^2} dq$$

$$B = \int_0^1 \frac{\lambda^2}{2a} \frac{2}{\sqrt{4ac - b^2}} \arctan \frac{\sqrt{4ac - b^2}}{b + 2c} dq$$

$$a = -sq(1-q) + m^2 \geq m^2, b = -\lambda^2, c = \lambda^2, 4ac - b^2 \geq 4m^2 \lambda^2$$

$$|B| \leq \frac{\lambda}{2m^3} \frac{\pi}{2} = \frac{\lambda\pi}{4m^3}$$

$$J = \int_0^1 \int_0^{1-x} \frac{1}{(-sxy + (x+y)^2 m^2)^{1-\eta}} dy dx \quad \eta = \frac{1}{\log(\frac{1}{\lambda^2})}$$

$$= \int_0^1 \int_0^1 \frac{p}{p^{2-2\eta}} \frac{1}{(-sq(1-q) + m^2)^{1-\eta}} dp dq = \frac{1}{2\eta} \int_0^1 \frac{(-sq(1-q) + m^2)^\eta}{-sq(1-q) + m^2} dq$$

$$= \frac{1}{2} \int_0^1 \left[\log\left(\frac{1}{\lambda^2}\right) \frac{1}{-sq(1-q) + m^2} + \frac{\log(-sq(1-q) + m^2)}{-sq(1-q) + m^2} \right] dq$$

$$= \frac{1}{2} \int_0^1 \frac{\log \frac{-sq(1-q) + m^2}{\lambda^2}}{-sq(1-q) + m^2} dq = A$$

$$I \doteq J \quad \text{誤差} = \frac{\lambda\pi}{4m^3}$$

Aの計算から $s < 0$ の場合の解析近似解は以下の様にして導く事ができます。

$$\begin{aligned} A &= \frac{1}{2} \int_0^1 \frac{\log \frac{-sq(1-q) + m^2}{\lambda^2}}{-sq(1-q) + m^2} dq \\ &= \frac{1}{2} \log\left(\frac{m^2}{\lambda^2}\right) \int_0^1 \frac{1}{-sq(1-q) + m^2} dq \\ &\quad + \frac{1}{2} \int_0^1 \frac{\log \frac{-sq(1-q) + m^2}{m^2}}{-sq(1-q) + m^2} dq \\ &= \frac{1}{-s} \log\left(\frac{m^2}{\lambda^2}\right) \log\left(\frac{-s}{m^2}\right) + \frac{1}{2} \frac{1}{-s} \left[\log^2\left(\frac{m^2}{-s}\right) - \frac{\pi^2}{3} \right] \\ &= \frac{1}{-s} \left[\log\left(\frac{m^2}{\lambda^2}\right) \log\left(\frac{-s}{m^2}\right) + \frac{1}{2} \log^2\left(\frac{m^2}{-s}\right) - \frac{\pi^2}{6} \right] \end{aligned}$$

S<0では

$$\inf_{\mathbf{r}, \mathbf{v}, \mathbf{x}}$$

$$s = -500^2, m = 0.0005$$

$$\lambda = 10^{-n}, \eta = \frac{1}{\log\left(\frac{1}{\lambda^2}\right)}$$

n	実測値	解析近似解
20000	0.101794846119901852D+02	0.101794845734671959D+02
30000	0.152693067851922386D+02	0.152693067595264598D+02
40000	0.203591289648179377D+02	0.203591289455857201D+02
50000	0.254489511470131227D+02	0.254489511316449821D+02
60000	0.305387733304930542D+02	0.305387733177042442D+02
70000	0.356285955147071292D+02	0.356285955037635063D+02
80000	0.407184176993800477D+02	0.407184176898227719D+02
90000	0.458082398843588621D+02	0.458082398758820304D+02
100000	0.508980620695518056D+02	0.508980620619412960D+02
110000	0.559878842549004716D+02	0.559878842480005545D+02
120000	0.610777064403659367D+02	0.610777064340598201D+02
130000	0.661675286259212498D+02	0.661675286201190858D+02
140000	0.712573508115471554D+02	0.712573508061783372D+02
150000	0.763471729972295208D+02	0.763471729922376028D+02
160000	0.814369951829577872D+02	0.814369951782968684D+02
170000	0.865268173687238260D+02	0.865268173643561340D+02
180000	0.916166395545213561D+02	0.916166395504153854D+02
190000	0.967064617403454179D+02	0.967064617364746510D+02
200000	0.101796283926192004D+03	0.101796283922533917D+03
210000	0.106886106112057917D+03	0.106886106108593182D+03
220000	0.111975928297940513D+03	0.111975928294652434D+03
230000	0.117065750483837618D+03	0.117065750480711699D+03
240000	0.122155572669747428D+03	0.122155572666770965D+03
250000	0.127245394855668394D+03	0.127245394852830231D+03
260000	0.132335217041599265D+03	0.132335217038889482D+03
270000	0.137425039227538889D+03	0.137425039224948762D+03
280000	0.142514861413486386D+03	0.142514861411008013D+03
290000	0.147604683599440904D+03	0.147604683597067265D+03
300000	0.152694505785401759D+03	0.152694505783126544D+03

と非常に良い結果が得られています。

2.2 vtx-2

さらに, $\int_0^1 \int_0^{1-x} \frac{1}{D^{1-\eta}} dy dx$ $D = xy$ $\eta = \frac{1}{\log(\frac{1}{\lambda^2})}$ の場合を

計算すると,

$$I = \int_0^1 \int_0^{1-x} \frac{1}{D^{1-\eta}} dy dx = \int_0^1 \frac{(1-x)^\eta}{\eta x^{1-\eta}} dx = \frac{1}{\eta} B(\eta, \eta+1) = \frac{\Gamma(\eta)\Gamma(1+\eta)}{\eta\Gamma(1+2\eta)}$$
$$= \frac{1}{\eta^2} \frac{(\Gamma(1+\eta))^2}{\Gamma(1+2\eta)}$$

となり, 二重指数関数型積分で容易に

計算出来ています。結果は積分変数の変換区間(0,1)

での最小値に依存し, 10^{-150} と 10^{-2465} (指数部のビット数が

11か15かに依存)では $\eta = 0.01$ ($\lambda = 10^{-21.7147241}$ に相等)で

その精度は10進1桁と10進24桁(反復法を使用せず)の

差があり, 前者では40回の反復により10進8桁となっ

ています。このケースでは $D \rightarrow D - i\varepsilon$ とする操作は必要

ありませんでした。

2.3 box

少し複雑なケース

$$\eta = \frac{1}{\log\left(\frac{1}{\lambda^2}\right)}$$

$$I = \int_0^1 \int_0^{1-x} \int_0^{1-x-y} \frac{1}{D^{2-\eta}} dz dy dx$$

$$D = xz + y(1-x-y-z)$$

では,

$$I = \frac{1}{1-\eta} \frac{2}{\eta} \left(\frac{2}{\eta} - \frac{\pi^2}{4} 2^{1-\eta} \eta \right) \frac{(\Gamma(1+\eta))^2}{\Gamma(1+2\eta)}$$

$$\frac{(\Gamma(1+\eta))^2}{\Gamma(1+2\eta)} = \frac{1+2\eta}{(1+\eta)^2} \exp\left[\sum_{k=2}^{\infty} (-1)^k (2-2^k) \left(\frac{\zeta(k)-1}{k} \right) \eta^k \right]$$

$$\log \Gamma(1+x) = -\log(1+x) + (1-\gamma)x + \sum_{k=2}^{\infty} (-1)^k \left(\frac{\zeta(k)-1}{k} \right) x^k \text{ より。}$$

$$I = \frac{a_{-2}}{\eta^2} + \frac{a_{-1}}{\eta} + a_0 \text{ とすると, } a_{-2} = 4, a_{-1} = 4, a_0 = -4\left(\frac{5}{12}\pi^2 - 1\right) = -\frac{5}{3}\pi^2 + 4$$

これは以下の事からきています。

$$I = \left[\int_0^1 \frac{(1-q)^{1-\eta} - q^{1-\eta}}{(1-2q)q^{1-\eta}(1-q)^{1-\eta}} dq \right] \times \frac{1}{1-\eta} \left[\int_0^1 \frac{1}{p^{1-\eta}(1-p)^{1-\eta}} dp \right]$$

$$\text{から } \frac{1}{1-\eta} \left[\int_0^1 \frac{1}{p^{1-\eta}(1-p)^{1-\eta}} dp \right] = \frac{1}{1-\eta} \frac{2}{\eta} \frac{(\Gamma(1+\eta))^2}{\Gamma(1+2\eta)}。$$

$$\int_0^1 \frac{(1-q)^{1-\eta} - q^{1-\eta}}{(1-2q)q^{1-\eta}(1-q)^{1-\eta}} dq = \frac{2}{\eta} - \frac{\pi^2}{4} 2^{1-\eta} \eta \text{ となるがより詳細には}$$

次ぎの様になります。

$$\int_0^1 \frac{(1-q)^{1-\eta} - q^{1-\eta}}{(1-2q)q^{1-\eta}(1-q)^{1-\eta}} dq = \int_0^1 \left[\frac{q^\eta}{(1-2q)q} - \frac{(1-q)^\eta}{(1-2q)(1-q)} \right] dq$$

$$= 2 \int_0^1 \frac{q^\eta - (1-q)^\eta}{1-2q} dq + \frac{2}{\eta} = \frac{2}{\eta} + 2^{1-\eta} \int_0^1 \frac{(1-t)^\eta - (1+t)^\eta}{t} dt$$

$$= \frac{2}{\eta} - 2^{1-\eta} 2\eta \left(1 + \sum_{k=1}^{\infty} \frac{(\eta-1)(\eta-2)\dots(\eta-2k)}{(2k+1)!(2k+1)} \right)$$

()の定数項は $-2 \left(1 + \sum_{k=1}^{\infty} \frac{1}{(2k+1)^2} \right) = -2 \left(\sum_{k=1}^{\infty} \frac{1}{(2k-1)^2} \right)$

$$= -2 \frac{3}{4} \zeta(2) = -\frac{\pi^2}{4}$$

()の η の一次の項は、負の数の和の絶対値で計算すると、

$$\sum_{k=1}^{\infty} \frac{\text{分子の項}}{(2k+1)!(2k+1)}$$

$$\text{分子の項} = (2k)! \sum_{r=1}^{2k} \frac{1}{r} < (2k)! [\log(2k+1) + \gamma]$$

$\gamma = 0.5772156649$ (オイラーの定数)

$$\sum_{k=1}^{\infty} \frac{(2k)! [\log(2k+1) + \gamma]}{(2k+1)!(2k+1)} = \sum_{k=1}^{\infty} \frac{\log(2k+1)}{(2k+1)^2} + \gamma \sum_{k=1}^{\infty} \frac{1}{(2k+1)^2}$$

$$\sum_{k=1}^{\infty} \frac{\log(2k+1)}{(2k+1)^2} = \pi^2 \left[\frac{3}{2} \log A - \frac{1}{6} \log(2) - \frac{1}{8} \log(\pi) - \frac{\gamma}{8} \right]$$

$$= 0.4181157384$$

$$A = \lim_{n \rightarrow \infty} \frac{1 \times 2^2 \times 3^3 \times \dots \times n^n}{n^{n^2/2+n/2+1/12}} e^{\frac{n^2}{4}} = 1.282427130$$

$$\log A = 0.2487544703$$

$$\gamma \sum_{k=1}^{\infty} \frac{1}{(2k+1)^2} = \gamma \left[\sum_{k=1}^{\infty} \frac{1}{(2k-1)^2} - 1 \right] = \gamma \left(\frac{\pi^2}{8} - 1 \right) = 0.1348956184$$

これから()内の η の一次の項の係数は -0.5530113568

全体では $-2^{1-\eta} \times 2\eta$ をかけて $2.212\eta^2$ となりこれが一次式まで取った場合の誤差となります。

実測例としては以下のものがあります。

解析近似式を使用して、 $\lambda = 10^{-1000000001} \sim 10^{-1000000045}$

で計算しました。ただしgamma関数はサポートされていない機種もあり、展開式より計算しています。

```
a2= 4.000000000000000065094194454995095
a1= 4.0000000000000001117756660514563709
a0= -12.4493406676118712066445306454111
```

変数変換区間SR16000 $[10^{-100}, 1 - 10^{-100}]$, x5570とx2670では $[10^{-200}, 1 - 2^{-113}]$ の4倍精度では

epsilon**算法の収束具合は**

$\lambda = 10^{-150}$ で

	値	誤差
x5570	0.191143398445070834D + 07	0.140633312313127879D - 05
x2670	0.191143398444577885D + 07	0.210270925388344072D - 07
sr16000	0.191143398444582948Q + 07	0.156658222921909210Q - 05

$\lambda = 10^{-1500}$ で

x5570	0.190895950695064375D + 09	0.919500023320091177D - 01
x2670	0.190895950598091269D + 09	0.477814929853799862D - 01
sr16000	0.190895950686589134Q + 09	0.470812225711499829Q - 01

$\lambda = 10^{-15000}$ で

x5570	0.190871021209837876D + 11	0.771928148538333458D + 04
x2670	0.190871096344010433D + 11	0.183847060803184605D + 03
sr16000	0.190871111864473018Q + 11	0.190871121114962521Q + 11

(注)SR16000では他に変化の小さい要素がある。

となっています。

あとSR16000でのある λ の範囲での a_{-2}, a_{-1}, a_0 は以下の様になっています。

ramda= $10^{-16} \sim 10^{-30}$

a2= 4.00000344724677645844034046757416

a1= 3.99861407619387554668745765583265

a0= -12.2639453402406000000000000020020

ramda= $10^{-31} \sim 10^{-45}$

a2= 4.00000099577869056434484930528836

a1= 3.99939509233538395973479843855165

a0= -12.326990512798999999999999759162

ramda= $10^{-46} \sim 10^{-60}$

a2= 4.00000041264201279800645140064705

a1= 3.99966395945103743579061852729121

a0= -12.358196281822999999999999796046

ramda=10^{-61}~10^{-75}

a2= 4.00000019676625921874246419452751

a1= 3.99979490092391388668793381657374

a0= -12.3781474808740000000000001848136

ramda=10^{-76}~10^{-90}

a2= 4.00000014459869220092907590162562

a1= 3.99983240900972265003642249066011

a0= -12.3848690023040000000000001603800

rambda=10^{-91}~10^{-105}

a2= 4.00000003694782985017339116664484

a1= 3.99996494892422507975597908036244

a0= -12.41545871332175260000000000000001

λ の値が小さくなるほど

a_0 は $-\frac{5\pi^2}{3} + 4$ に近付いています。

3 sin c 求積法による Hadamard 有限部分計算

端点と領域内部で同時に特異点をもつ

場合の計算を, 数理解析研究所講究録

第791巻1992年206-219

Hadamard 有限部分積分に関する DE 公式

緒方秀教, 杉原正顯, 森正武

(東京大学工学部物理工学科)

の資料から以下の問題を測定しました。

$$\text{f.p.} \int_{-1}^1 \frac{F(x)}{(x-\lambda)^2} dx, F(z) = (1-z)^{\frac{1}{4}} (1+z)^{\frac{-1}{4}}$$

$$\text{厳密値} = -\frac{\pi}{2} (1-\lambda)^{\frac{-3}{4}} (1+\lambda)^{\frac{-5}{4}}$$

$$\text{公式 I} = h \sum_{k=-N_1}^{N_2} \frac{F(x_k) / \phi'(x_k)}{(x_k - \lambda)^2} + \left\{ \cot \left[\frac{\pi}{h} \phi(\lambda) \right] F'(\lambda) - \frac{\pi}{h} \frac{\phi'(\lambda)}{\sin^2 [\pi \phi(\lambda) / h]} F(\lambda) \right\}$$

3.1 SE,DEの計算方法

$$\text{DE} : x = \tanh\left(\frac{\pi}{2} \sinh(t)\right) \quad \phi(x) = \sinh^{-1}\left(\frac{1}{\pi} \log\left(\frac{1+x}{1-x}\right)\right)$$

$$s = \frac{1}{\pi} \log\left(\frac{1+x}{1-x}\right), \phi(x) = \log(s + \sqrt{s^2 + 1})$$

$$\phi'(x) = \frac{\frac{1}{\pi} \frac{2}{(1-x^2)}}{\sqrt{\left(\frac{1}{\pi} \log\left(\frac{1+x}{1-x}\right)\right)^2 + 1}}$$

$$\text{SE} : x = \tanh\left(\frac{t}{2}\right) \quad \phi(x) = \log\left(\frac{1+x}{1-x}\right) \quad \phi'(x) = \frac{1}{1-x^2}$$

$$F'(z) = -\frac{1}{2(1-z)^{\frac{3}{4}}(1+z)^{\frac{5}{4}}}$$

$x = 1 - \varepsilon, -1 + \varepsilon$ となる t

$$\text{DE} \quad t = \log\left(\frac{2y}{\pi} + \sqrt{\left(\frac{2y}{\pi}\right)^2 + 1}\right) \quad 2y = \log\left(\frac{1}{\varepsilon}\right)$$

$$\text{SE} \quad t = \ln\left(\frac{2}{\varepsilon}\right)$$

$\varepsilon = 2^{-106}$ の場合, DE : $t = 3.855$, SE : $t = 74.1667$

$h = 0.125$ で測定

ここでsinhの逆関数はサポートされていないので、logを使用した式を使っています。

3.2 ε - 算法

$$\text{f.p} \int_{-1}^1 \frac{F(x)}{(x - \lambda)^2} dx = \lim_{\varepsilon \rightarrow 0} \int_{-1}^1 \frac{F(x)}{(x - \lambda + i\varepsilon)^2} dx$$

で計算しています。

$h = 0.5^8, \varepsilon_0 = 1.0, \varepsilon_n = \left(\frac{1}{1.1}\right)^{n-1}$ で $n = 50$ まで計算

今回は領域を分割しない場合と、領域を

$[-1, \lambda], [\lambda, 1]$ と分割した場合を計算しました。

変数変換区間 $[-1 + \varepsilon, 1 - \varepsilon]$, $\varepsilon = 2^{-106}$ の場合、

領域分割なしでは、 $N = 1976$ ですが、領域分割する

場合、 $\lambda = \pm 0.9$ でオーバーフローを防ぐ処理が必要

になります。

($x = x30(i) * \text{cnt}0 + \text{cnt}1$ で発生する丸め誤差のため)

$N = 1976$ の場合、 $x + 1.0q0 > 0.0, 1.0q0 - x > 0.0$ の場合のみ

計算するとします。また $\varepsilon = 2^{-104.5}$ として $N = 1962$ とする

方法もあります。

3.3 計算結果精度比較

4つの方法で精度を検証した結果は以下の様になっています。

変数変換区間は全て $[-1 + \varepsilon, 1 - \varepsilon]$

$\varepsilon = 2^{-106}$ としています。

λ	DE(分割なし)	DE(分割あり)	sinc (DE)	sinc (SE)
-0.9	0.1400D-04	0.1382Q-06	0.8469D-21	0.1825D-21
-0.8	0.6763D-06	0.5489Q-11	0.6245D-24	0.4563D-22
-0.7	0.5165D-06	0.3596Q-14	0.1757D-24	0.2028D-22
-0.6	0.2143D-06	0.2125Q-15	0.1855D-25	0.1143D-22
-0.5	0.6756D-07	0.1301Q-15	0.9166D-26	0.7301D-23
-0.4	0.3159D-07	0.4492Q-17	0.6154D-26	0.5070D-23
-0.3	0.1333D-07	0.1588Q-16	0.4696D-26	0.3714D-23
-0.2	0.3104D-07	0.2240Q-16	0.5978D-26	0.2852D-23
-0.1	0.1650D-07	0.7873Q-17	0.2845D-26	0.2253D-23
0.1	0.2247D-08	0.4837Q-20	0.1905D-26	0.1508D-23
0.2	0.1024D-07	0.1085Q-17	0.1217D-25	0.1268D-23
0.3	0.1762D-06	0.1065Q-17	0.1360D-26	0.1050D-23
0.4	0.7024D-06	0.2245Q-17	0.1172D-26	0.9314D-24
0.5	0.9061D-07	0.1208Q-16	0.1010D-26	0.8110D-24
0.6	0.7005D-06	0.1356Q-15	0.1099D-26	0.7211D-24
0.7	0.1890D-06	0.5910Q-14	0.5852D-25	0.6299D-24
0.8	0.2144D-06	0.3936Q-11	0.2826D-24	0.5632D-24
0.9	0.5647D-06	0.1723Q-06	0.1583D-21	0.5057D-24

今回の結果ではsinc (DE) 分点数62で最も適しているという結果になりました。

4. 多倍長減算桁落ちメモ

整数演算方式では加減算において内部処理が減算の場合の桁落ちと丸め処理が特に問題となりますので、作成時使用しましたメモを紹介します。

多倍長整数減算 $IZ = IX - IY$ ($IX > IY > 0$)の結果が64ビット整数変数の60ビットのみ使用するとします。

4.1 3倍精度

配列数 2

$$IZ(2) \geq 2^{22} \quad IZ(2) = IZ(2) - 2^{22} \quad \text{桁落ちなし.}$$

丸め位置 $IZ(1)(1:1)$

$$2^{21} \leq IZ(2) < 2^{22} \quad IZ(2) = IZ(2) - 2^{21} \quad \text{丸め位置 } IZ(1)(0:0)$$

丸め後 $IZ(2) = 2^{22}$ 桁落ちなし.

$$IZ(2) \neq 2^{22} \quad \text{桁落ち1}$$

$$2^{20} \leq IZ(2) < 2^{21} \quad IZ(2) = IZ(2) - 2^{20} \quad \text{桁落ち2 丸めなし}$$

$$IZ(2) < 2^{20} \quad IZ(ikk) = IZ(ikk) - 2^{ik}$$

桁落ち $82 - (ikk - 1) \times 60 - ik$

丸めなし。 ($1 \leq ikk \leq 2, 0 \leq ik \leq 59$)

4.2 5倍精度

配列数 3

$$IZ(3) \geq 2^{26} \quad IZ(3) = IZ(3) - 2^{26} \quad \text{桁落ちなし.}$$

丸め位置 $IZ(1)(1:1)$

$$2^{25} \leq IZ(3) < 2^{26} \quad IZ(3) = IZ(3) - 2^{25} \quad \text{丸め位置 } IZ(1)(0:0)$$

丸め後 $IZ(3) = 2^{26}$ 桁落ちなし

$$.IZ(3) \neq 2^{26} \quad \text{桁落ち1}$$

$$2^{24} \leq IZ(3) < 2^{25} \quad IZ(3) = IZ(3) - 2^{24} \quad \text{桁落ち2 丸めなし}$$

$$IZ(3) < 2^{24} \quad IZ(ikk) = IZ(ikk) - 2^{ik}$$

桁落ち $146 - (ikk - 1) \times 60 - ik$

丸めなし。 ($1 \leq ikk \leq 3, 0 \leq ik \leq 59$)

4.3 6倍精度

配列数 3

$$IZ(3) \geq 2^{58} \quad IZ(3) = IZ(3) - 2^{58} \quad \text{桁落ちなし.}$$

丸め位置 IZ(1)(1:1)

$$2^{57} \leq IZ(3) < 2^{58} \quad IZ(3) = IZ(3) - 2^{57} \quad \text{丸め位置 IZ(1)(0:0)}$$

丸め後 $IZ(3) = 2^{58}$ 桁落ちなし.

$$IZ(3) \neq 2^{58} \quad \text{桁落ち1}$$

$$2^{56} \leq IZ(3) < 2^{57} \quad IZ(3) = IZ(3) - 2^{56} \quad \text{桁落ち2 丸めなし}$$

$$IZ(3) < 2^{56} \quad IZ(ikk) = IZ(ikk) - 2^{ik}$$

桁落ち $178 - (ikk - 1) \times 60 - ik$

丸めなし。 ($1 \leq ikk \leq 3, 0 \leq ik \leq 59$)

4.4 7倍精度

配列数 4

$$IZ(4) \geq 2^{30} \quad IZ(4) = IZ(4) - 2^{30} \quad \text{桁落ちなし.}$$

丸め位置 IZ(1)(1:1)

$$2^{29} \leq IZ(4) < 2^{30} \quad IZ(4) = IZ(4) - 2^{29} \quad \text{丸め位置 IZ(1)(0:0)}$$

丸め後 $IZ(4) = 2^{30}$ 桁落ちなし

$$IZ(4) \neq 2^{30} \quad \text{桁落ち1}$$

$$2^{28} \leq IZ(4) < 2^{29} \quad IZ(4) = IZ(4) - 2^{28} \quad \text{桁落ち2 丸めなし}$$

$$IZ(4) < 2^{28} \quad IZ(ikk) = IZ(ikk) - 2^{ik}$$

桁落ち $210 - (ikk - 1) \times 60 - ik$

丸めなし。 ($1 \leq ikk \leq 4, 0 \leq ik \leq 59$)

5. 仮数部のビット数と乗算のメモ 非対称行列の反復解法で使した仮数部の ビット数を調整した乗算のメモを紹介します。

仮数部のビット数に於ける乗算の場合の桁上りと丸め
多倍長整数乗算 $Z = IX \times IY$ ($IX > 0, IY > 0$)の結果
が32ビット整数変数の0ビットのみ使用するとます。
符号部1ビット, 指数部15ビットを想定。

(1)60ビット

配列数 5

$IZ(5) \geq 2$ 桁上りあり 丸め位置 $IZ(3)(0:0)$

$IZ(5) < 2$ 桁上りなし 丸め位置 $IZ(2)(29:29)$

(2)64ビット

配列数 5

$IZ(5) \geq 2^9$ 桁上りあり 丸め位置 $IZ(3)(4:4)$

$IZ(5) < 2^9$ 桁上りなし 丸め位置 $IZ(3)(3:3)$

(3)68ビット

配列数 5

$IZ(5) \geq 2^{17}$ 桁上りあり 丸め位置 $IZ(3)(8:8)$

$IZ(5) < 2^{17}$ 桁上りなし 丸め位置 $IZ(3)(7:7)$

(4)72ビット

配列数 5

$IZ(5) \geq 2^{25}$ 桁上りあり 丸め位置 $IZ(3)(12:12)$

$IZ(5) < 2^{25}$ 桁上りなし 丸め位置 $IZ(3)(11:11)$

(5)76ビット

配列数 6

$IZ(6) \geq 2^3$ 桁上りあり 丸め位置 $IZ(3)(16:16)$

$IZ(6) < 2^3$ 桁上りなし 丸め位置 $IZ(3)(15:15)$

6 DD,D (long double) ,DQのインライン化メモ

E5-2670,E5-2660,Phi5110P等で2つの変数の和で表す演算を使用した時,並列化効果を出すために作成したソースのメモを紹介します。

6.1 FORTRAN

(1) 加算 $c=a+b$

```
p1=a(1)
p2=a(2)
q1=b(1)
q2=b(2)
t1=p1+q1
t2=t1-p1
t3=(p1-(t1-t2))+(q1-t2)
t4=p2+q2
t5=t4-p2
t6=(p2-(t4-t5))+(q2-t5)
t7=t3+t4+t6
t8=t1+t7
t9=t8-t1
t10=(t1-(t8-t9))+(t7-t9)
c(1)=t8
c(2)=t10
```

(2)減算 $c=a-b$

```
p1=a(1)
p2=a(2)
q1=b(1)
q2=b(2)
t1=p1-q1
t2=t1-p1
t3=(p1-(t1-t2))-(q1+t2)
t4=p2-q2
t5=t4-p2
t6=(p2-(t4-t5))-(q2+t5)
t7=t3+t4+t6
t8=t1+t7
t9=t8-t1
t10=(t1-(t8-t9))+(t7-t9)
c(1)=t8
c(2)=t10
```

(3)乘算 $c=a*b$

$$p1=a(1)$$

$$p2=a(2)$$

$$q1=b(1)$$

$$q2=b(2)$$

$$t1=p1*q1$$

$$t2=p1*r$$

$$a1=t2-(t2-p1)$$

$$a2=p1-a1$$

$$t3=q1*r$$

$$b1=t3-(t3-q1)$$

$$b2=q1-b1$$

$$t4=((a1*b1-t1)+a1*b2+a2*b1)+a2*b2$$

$$t5=t4+p1*q2$$

$$t6=t5+p2*q1$$

$$t7=t6+p2*q2$$

$$t8=t1+t7$$

$$t9=t8-t1$$

$$t10=(t1-(t8-t9))+(t7-t9)$$

$$c(1)=t8$$

$$c(2)=t10$$

$$r = 134217729 = 2^{27} + 1(\text{倍精度})$$

$$r = 2^{57} + 1(4\text{倍精度})$$

(4)除算 $c=a/b$

$$p_3=a(1)$$

$$p_4=a(2)$$

$$q_3=b(1)$$

$$q_4=b(2)$$

$$t_s=p_3/q_3$$

$$q_1=t_s*q_3$$

$$q_2=t_s*q_4$$

$$p_1=p_3$$

$$p_2=p_4$$

$$t_1=p_1-q_1$$

$$t_2=t_1-p_1$$

$$t_3=(p_1-(t_1-t_2))-(q_1+t_2)$$

$$t_4=p_2-q_2$$

$$t_5=t_4-p_2$$

$$t_6=(p_2-(t_4-t_5))-(q_2+t_5)$$

$$t_7=t_3+t_4+t_6$$

$$t_8=t_1+t_7$$

$$t_1=t_s$$

$$t_7=t_8/q_3$$

$$t_8=t_1+t_7$$

$$t_9=t_8-t_1$$

$$t_{10}=(t_1-(t_8-t_9))+(t_7-t_9)$$

$$c(1)=t_8$$

$$c(2)=t_{10}$$

6.2 C

(1) 加算 $c1=work11+work12$

```
p1=work11[0] ;
p2=work11[1] ;
q1=work12[0] ;
q2=work12[1] ;
t1=p1+q1 ;
t2=t1-p1 ;
t3=(p1-(t1-t2))+(q1-t2) ;
t4=p2+q2 ;
t5=t4-p2 ;
t6=(p2-(t4-t5))+(q2-t5) ;
t7=t3+t4+t6 ;
t8=t1+t7 ;
t9=t8-t1 ;
t10=(t1-(t8-t9))+(t7-t9) ;
c1[0]=t8 ;
c1[1]=t10 ;
```

(2) 減算 $work1=one-xx$

```
p1=one[0] ;
p2=one[1] ;
q1=xx[0] ;
q2=xx[1] ;
t1=p1-q1 ;
t2=t1-p1 ;
t3=(p1-(t1-t2))-(q1+t2) ;
t4=p2-q2 ;
t5=t4-p2 ;
t6=(p2-(t4-t5))-(q2+t5) ;
t7=t3+t4+t6 ;
t8=t1+t7 ;
t9=t8-t1 ;
t10=(t1-(t8-t9))+(t7-t9) ;
work1[0]=t8 ;
work1[1]=t10 ;
```

(3)乗算 work3=work1*work2

```
p1=work1[0];
p2=work1[1];
q1=work2[0];
q2=work2[1];
t1=p1*q1;
t2=p1*r;
a1=t2-(t2-p1);
a2=p1-a1;
t3=q1*r;
b1=t3-(t3-q1);
b2=q1-b1;
t4=((a1*b1-t1)+a1*b2+a2*b1)+a2*b2;
t5=t4+p1*q2;
t6=t5+p2*q1;
t7=t6+p2*q2;
t8=t1+t7;
t9=t8-t1;
t10=(t1-(t8-t9))+(t7-t9);
work3[0]=t8;
work3[1]=t10;
```

$$r = 134217729 = 2^{27} + 1(\text{倍精度})$$

$$r = 8589934593 = 2^{33} + 1(\text{拡張倍精度})$$

(4)除算 $w_3 = \text{work8} / \text{work2}$

```
p3=work8[0] ;  
p4=work8[1] ;  
q3=work2[0] ;  
q4=work2[1] ;  
ts=p3/q3 ;  
q1=ts*q3 ;  
q2=ts*q4 ;  
p1=p3 ;  
p2=p4 ;  
t1=p1-q1 ;  
t2=t1-p1 ;  
t3=(p1-(t1-t2))-(q1+t2) ;  
t4=p2-q2 ;  
t5=t4-p2 ;  
t6=(p2-(t4-t5))-(q2+t5) ;  
t7=t3+t4+t6 ;  
t8=t1+t7 ;  
t1=ts ;  
t7=t8/q3 ;  
t8=t1+t7 ;  
t9=t8-t1 ;  
t10=(t1-(t8-t9))+(t7-t9) ;  
w3[0]=t8 ;  
w3[1]=t10 ;
```