

多倍長計算手法

平成25年度第4四半期

目次

1. はじめに
2. 整数演算方式による128倍精度演算の作成
 - 2.1 概略
 - 2.2 加減算
 - 2.3 乗算
 - 2.4 作成作業に関して
3. 128倍精度演算の精度
 - 3.1 平方根計算
 - 3.2 ヒルベルト行列
4. 128倍精度演算の性能
5. 量子モンテカルロ法による物性スペクトル計算
 - 5.1 大きな β の場合に関して
 - 5.2 パラメータ領域拡大結果

1.はじめに

量子モンテカルロ法による物性スペクトル計算で正しく計算出来るパラメータ範囲の限界を求めるために、32倍精度変数をつなげる方式で、64,96,128倍精度演算を作成して32倍精度演算でのパラメータ範囲の限界値 $\beta = 250$ $U = 10, N = 100, L = 448$ の条件でx5570で実行した所、64倍精度演算で2日弱、96倍精度演算で1週間、128倍精度演算で2週間かかっています。最終的には64倍精度演算で正しく計算出来るパラメータ範囲は $\beta = 600, U = 10, N = 100$ $L = 448$ まで拡大しています。

ieee754 – 2008形式では表現できる数値範囲から理論的には128倍精度演算まで実行可能ですが、1ケースのテストで2週間もかかるため、128倍精度演算を整数演算方式で行う様にしました。ここで問題となるのはソースステップ数の多さです。

原因としては8倍精度演算等の低精度演算では性能などの問題、32倍精度より低精度演算では繰り返しパターンがほとんどないため *do*文を使用せず、演算をすべて展開するという方式を取っていた事によります。

128倍精度演算となると繰り返しパターンも多くなり、*do*文を使用せず演算を展開してもその性能上の効果もなく、あまり展開が多いと性能上も逆効果となりますので、*do*文を使用してソースステップ数の削減を図りました。

整数演算方式では正整数演算 $a \pm b, a \times b$ が負の値ならない
様にするため、整数型4倍バイト配列は各要素30ビット使用する、
整数型8倍バイト配列は各要素60ビット使用する様にしています。
また8倍精度変数の仮数部のビット数は240ビットと低精度演算では
もっとも扱いやすくなっています。このため、 $8 + 15 \times k$ (k は正整数)
精度演算が同じパターンがでる事になります。また $4n$ 倍精度
のほうが無駄な処理が少なくなるため $8 + 60 \times k$ 倍精度演算が
適していることがわかります。 $k = 2$ とすれば $8 + 60 \times k = 128$
となり、128倍精度演算が適している事がわかります。
最終的には、1つのサブルーチンが350ステップ内で収まり、
加減算、乗算のルーチンでは共通する部分がありため、デバッグ
作業等も含めても、1人で扱える作業量となりました。
ヒルベルト行列の次元数800(条件数 10^{1200})で10進8桁まで一致しました。
これから、128倍精度演算から188倍精度演算、248倍精度演算に
するには定数を変更するだけで良く、実際デバッグ($N = 100$ のヒルベルト
行列の誤差を求める)を含め、188倍精度は2時間、これで慣れた事により
248倍精度演算は1時間で作業が終了しました。
この事から $8 + 60 \times k$ ($k = 0, 1, 2, 3, 4, 5, 6, 7, 8$)倍精度演算を作成し、
加減算、乗算、除算、平方根計算、判定(gt, eq, lt)の合計のステップ数
はコメント、定義文等を含め、 $k = 0$ で1000ステップ、 $k \neq 0$ で1500ステップ
となっています。今回は128倍精度演算を中心に必要におおじて他の
精度演算にかんしても記述しました。

2.整数演算方式による128倍精度演算の作成

2.1概略

除算,平方根計算は以下の様に反復法を用いる事で加減算,乗算の作成に帰着されます。

平方根計算

$$f(x) = x^2 - a, f'(x) = 2x$$

$$x_{n+1} = x_n - \frac{x_n^2 - a}{2x_n} = \frac{1}{2} \left(x_n + \frac{a}{x_n} \right)$$

最終段階 $x_{n+1} = \frac{1}{2} \left(x_n - \frac{a}{x_n} \right) + \frac{a}{x_n}$

初期値 $x_0 = \sqrt{a}$ (4倍精度)

除算

$c = \frac{a}{b} = a \times \frac{1}{b}$ で計算します。

逆数計算は

$$f(x) = \frac{1}{x} - b, f'(x) = -\frac{1}{x^2}$$

$$x_{n+1} = x_n + \frac{\frac{1}{x_n} - b}{-\frac{1}{x_n^2}} = x_n + x_n(1 - bx_n) = x_n(2 - bx_n)$$

で行います。

初期値 $x_0 = \frac{1}{b}$ (4倍精度)

(注) $x_n, 1 - bx_n$ が小さい場合, $x_n(1 - bx_n)$ の計算でアンダーフロー発生の可能性があるので $x_{n+1} = x_n(2 - bx_n)$ としています。

2.2 加減算

加減算は引数を整数型8バイトの60ビットのみ使用する配列に展開します。準備する配列要素の数は $4080 \div 60 + 1 = 69$ となります。引数を正の整数の演算になる様に引数を組み替えます。(この処理は加算及び減算に共通です。)

この処理は低精度演算のものを定数を変更するだけですので前処理サブルーチンに分けます。(ソースステップ数約150) 隠れビットと丸めようのビットのために、先頭配列 $ix(69)$ の値は加算 $2 \leq ix(69) \leq 3$ 減算 $4 \leq ix(69) \leq 7$ となります。

引数の大きい方を配列 ix , 小さい方を配列 iy にいます。

そして指数部の差分だけ配列 iy をシフトして、桁合わせをします。

ここまでは加算及び減算に共通です。ビット位置が1つことなるのでサブルーチンとしての抽出はしていません。

多倍長整数加算または乗算を行い結果を配列 iz にいます。

この段階では整数型8バイト配列 iz は各要素60ビットしか使用しませんここで指数部と丸め処理をします。

加算の場合 $iz(69) \geq 4$ 桁上がりあり, $iz(69) \leq 3$ 桁上がりなし

減算の場合 $iz(69) \geq 4$ 桁下がりなし, $2 \leq iz(69) \leq 3$ 桁下がり1ビット

このケースの場合,ともに丸め処理後,指数部の値と仮数部を配列 iz の

68要素にいます。処理がことなるのは,減算で $iz(69) = 1$ の場合は

桁下がり2ビットで丸め処理は有りません。 $iz(69) = 0$ の場合は

$iz(68), iz(67), \dots$ と調べて最初に1があるビット位置を求める処理が

ある事です。これは $iz(l) \neq 0, iz(l) = 2^m + k (1 \leq k \leq 2^m - 1)$ を実数にすると

m が指数部の値に一致する事を使用します。これで指数部の値と

仮数部を配列 iz の68要素にいます。($iz(69) = 0$ の場合は前につめ、

下位の配列要素には0をいます。)。この後は加減算, 乗算共通で

60ビット68要素配列を64ビット64要素配列に詰め, 128倍精度浮動小数点数を作成します。とくにここもサブルーチン化して抜き出さなくてむ、コメントなども含めても約350ステップ数内に収まっています。

2.3 乗算

2つの引数から積の符号と指数部の値を計算して引数を4080ビットの正整数にします。これを整数型4バイトの30ビットのみ使用する配列に展開します。準備する配列要素の数は $4080 \div 30 + 1 = 137$ となります。二つの引数を配列 ix, iy にいれ多倍長正整数乗算を行い、結果を整数型4倍バイトの30ビット使用する配列 iz にいれます。 $(iz$ の要素数は274)。

$ix(137), iy(137)$ には隠れビット1をいれますので、桁上がりがあるかどうかは $iz(273)$ の値で求まり、これで丸めビットの位置がわかります。 $iz(273) \geq 2$ だと桁上がりがあり、 $iz(272) = 1$ だと桁上がりが有りません。ここで整数型8バイト配列 izz (配列要素数69)を用意して、 $iz \rightarrow izz$ にいれ、加減算と同じ処理をして、128倍精度浮動小数点数を作成します。これらを含めて300ステップ以下で作成できています。

2.4 作成作業に関して

整数演算方式の作成の作業量は連続して使用できるビット数によります。例えば8192ビット使用できるとすれば、ソースステップ数は128倍精度演算は20-30程度で行う事が可能と言えます。また丸め処理は行わない(例えば切り捨て)とすればやはりソースステップ数が大きく削減されます。FPGA等で作成する場合もほとんど同じ様になります。

3.128倍精度演算の精度

量子モンテカルロ法による物性スペクトル計算で演算精度に影響があるのは、逆行列を求める計算とグラムシュミット法の計算です。そこでヒルベルト行列を使用した連立一次方程式の求解と平方根計算の精度を求めました。

また128倍精度演算では、量子モンテカルロ法による物性スペクトル計算では表現可能な数値範囲の制限にかかる可能性があり、問題発生時にその原因が有効ビット数不足なのか、表現可能な数値範囲の制限によるものかの切り分けのため、ヒルベルト行列では188倍精度演算,248倍精度演算,308倍精度演算,368倍精度演算,428倍精度演算,488倍精度演算でも実行しています。

3.2 ヒルベルト行列

ヒルベルト行列のサイズNと条件数には近似的には以下の様になっています。

ヒルベルト行列の条件数				
n	条件数 (ビット数表示)			
400	2022			
450	2276			
500	2530			
550	2785			
600	3039			
650	3293			
700	3547			
750	3801			
800	4056	128倍精度 有効ビット数4081		
850	4310			
900	4564			
950	4818			
1000	5072			
1050	5327			
1100	5581			
1150	5835	188倍精度 有効ビット数 6001		
1200	6089			
1250	6344			
1300	6598			
1350	6852			
1400	7106	248倍精度 有効ビット数 7921		
1450	7361			
1500	7615			

実測結果は以下の様になっています。

ヒルベルト行列の精度		
N	最大誤差	精度
400	1.014E-621	128倍精度
500	1.681E-467	128倍精度
600	1.321E-314	128倍精度
700	1.349E-161	128倍精度
800	7.603E-009	128倍精度
1000	1.013E-280	188倍精度
1100	4.211E-127	188倍精度
1400	1.013E-245	248倍精度

最大誤差のビット数 = 有効ビット数 - 条件数のビット数
となっています。

248倍精度演算では $N = 1400$ が使用できるメモリ容量
の上限となっています。

より高精度演算での精度

N=100でのヒルベルト行列計算			
		条件数0.127E+152	ビット数501
精度	有効ビット数	誤差	
		実測	理論値
308	9841	4.73E-2814	2.80E-2812
368	11761	3.38E-3392	2.52E-3390
428	13681	1.61E-3971	2.66E-3968
488	15601	5.36E-4198	2.80E-4546 (注)
488	15601	1.00E-4547	2.80E-4546
(注)逆数計算で4倍精度の初期値から7回反復した場合の値 308,368,428倍精度は7回反復,488倍精度は8回反復が 必要			

488倍精度以上の精度を作成してもこれ以上精度向上はないと言えます。

4.128倍精度演算の性能

比較のため行列サイズN=400のヒルベルト行列による連立一次方程式の求解で実施しました。

ヒルベルト行列のN=400の実行時間一覧表

方式	精度	x5570	e5430
浮動型	64倍精度	540.603	743.938
浮動型	96倍精度	1551.696	2165.918
浮動型	128倍精度	3467.115	4280.500
整数型	68倍精度	179.059	277.818
整数型	128倍精度	652.896	918.142
整数型	188倍精度	1235.567	1945.560
整数型	248倍精度	2125.342	3337.811

単位: 秒

**32倍精度変数をつなげる方式に比べ、
整数演算方式の性能が非常に良くなっています。**

他の精度での性能

N=100でのヒルベルト行列計算			
実行時間 (秒)			
精度	x5570	e5430	
308	50.298	81.123	
368	70.468	115.572	
428	95.133	156.090	
488	123.216	205.238	
488	123.253	205.404	反復回数7回
N=400 での8倍精度でのヒルベルト行列計算			
Ax=b でAの要素は4倍精度演算で作成			
実行時間 (秒)			
ソース	x5570	e5430	
従来	4.686	5.926	
新規	4.885	6.354	

- 1.反復回数の増加の影響はほとんどない。
- 2.8倍精度演算ではソースステップ数の大幅削減をしても性能低下は5%程度に収まっています。

5.量子モンテカルロ法による物性スペクトル計算

5.1 大きな β の場合に関して

量子モンテカルロ法による物性スペクトル計算

では $B_L \dots B_1 = QDR$ とし $(I + QDR)^{-1}$ を求めるのに、

$D + Q^T R^{-1} = Q'D'R'$ と QDR 分解し、

$Q(D + Q^T R^{-1})R = QQ'D'R'R = Q''D''R'' = I + QDR$ から
逆行列を求めます。

$$P = e^Q \left(Q = \sqrt{\frac{\beta U}{L}} \left(\frac{L+1}{2} \right) \right) \text{とすると,}$$

$\beta U \leq L$ なら結果の絶対値の最小値は $\frac{1}{P}$,

Dの最大数は P^2 となり、グラムシュミット法で
平方数の和の平方根を求めるため、

プログラムに現れる最大数は P^4 となります。

$\beta U \geq L$, では条件数は β/L の大きさ依存する度合いが大きくなる。

これは, $B_L \dots B_1 = QDR$ とし $(I + QDR)^{-1}$ を求める際に発生します。

$D + Q^T R^{-1} = Q'D'R'$ とQDR分解し,

$Q(D + Q^T R^{-1})R = QQ'D'R'R = Q''D''R'' = I + QDR$ から逆行列を求めます。

この $D + Q^T R^{-1}$ のQDR分解の際オーバーフローが発生します。

正確には, $b = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}$ の計算の際

$a_i \times a_i$ の計算で発生します。

オーバーフローを防ぐには,

$a_{\max} = \text{絶対値最大値}(a_i, i = 1, 2, \dots, n)$

として, $b = |a_{\max}| \times \sqrt{c_1^2 + c_2^2 + \dots + c_n^2}$ ($c_i = \frac{a_i}{a_{\max}}$)

とします。ここで

正しい答えの β 最大数における, D'' の最大数 S ,

$S = e^{\sqrt{\beta UL} \times \alpha}$ としてその影響を調査しました。

調査では $n = 20, U = 10, L = 448$ と固定しています。

α 調査結果

条件 $N = 20, U = 10, L = 448$

精度	β の最大数	最大値	α
68倍精度	732	2.518D+1287	1.64
128倍精度	1620	7.778D+2441	2.09
188倍精度	2590	6.198D+3593	2.43
248倍精度	3609	6.580D+4745	2.72
308倍精度	3776	6.591D+4930	2.76

**368倍精度では308倍精度と同じく、
 $\beta = 3776$ では正しく計算出来ていて、
 $\beta = 3777$ では結果が不正となっている事から
308倍精度の結果が限界を示している
と言えます。**

N=20,U=10,L=448と固定すると正しく計算出来るβの最大値は3776でα=2.76であった。ここでN=20,U=10でLを小さくした場合正しく計算できた場合は以下の表のようになっています。

n=20,U=10 での実行結果一覧				
β	L	最大値	$\sqrt{\beta UL}$	α
3776	448	6.571D+4930	1.730D+1786	2.76
4000	320	2.857D+4931	6.014D+1553	3.18
4250	208	6.161D+4919	1.780D+1291	3.81
4500	132	1.077D+4927	2.934D+1058	4.66

この表からβが3776以上になると、βが250増加するとLは0.64倍となる。β = 3776と4000ではβが224の増加なのでLは $250/224 \times 0.64 = 0.714$ 倍となっています。

β₁, U₁, L₁, α₁が既知としてβ₁ < β₂のβ₂を選ぶと、L₂が推定され実行するU₂を決めれば(通常 = U₁)とすると

$$\alpha_2 = \alpha_1 \sqrt{\frac{\beta_1 U_1 L_1}{\beta_2 U_2 L_2}}$$

からα₂が求まり実行可能かどうか分かる。

$$\beta = 4750, L = 84, \alpha = 5.686 \quad \text{最大値} = 10^{4932.6}$$

$$\beta = 5000, L = 54, \alpha = 6.912 \quad \text{最大値} = 10^{4932.5}$$

と計算が困難になる事がわかります。

5.2 パラメータ領域拡大結果

測定結果一覧	x5570			
測定条件 N=100,U=10,L=448				
βと絶対温度Kとの関係				
絶対温度(K) = 10000/β				
β	精度	絶対値最小値		実行時間
		実測値	理論値	(秒)
730	68倍精度	0.824D-384	0.269D-383	57642
1600	128倍精度	0.170D-582	0.214D-582	185097
2500	188倍精度	0.773D-728	0.462D-728	403155
10000/3	248倍精度	0.292D-840	0.978D-841	673631
3770	308倍精度	0.156D-893	0.397D-894	975735

この条件では、32倍精度 $\beta = 250$ (絶対温度40K)から、68倍精度で13.7K、128倍精度で6.26K、188倍精度で4K、248倍精度で3K、308倍精度で2.65Kまでパラメータ領域が拡大されました。これ以上の演算精度を用いても改善はされませんのでこの値が限界と言えます。