

## Accuracy-Recording Quantum-Chemical Calculations (Hokkaido University of Education, Kushiro Campus) Shigeru Obara

### 1 【序】

コンピュータの発達に伴い量子化学計算の研究対象は大型化して来ている。その一方で計算そのものは依然として倍精度実数を使って行われており、精度不足の発生が懸念される。精度不足発生の有無そのものは、別途精度の高い計算を行って比較することにより確認することはできる。しかし、このような別途計算を常に行うことは不便である。それに、きめ細かな解析を行うことも容易でない。通常の量子化学計算を行いながら計算精度も常に記録していくことができるとこれらの不便が解消する。この観点から、どのように精度を記録していけば良いのか等を検討して本研究室で開発した高精度量子化学計算プログラムに精度記録機能を追加した。これについて報告する。

### 2 【有効ビット数の増減要因】

理系の数値計算では有効数字という概念が用いられ、その桁数不足発生の有無が常に注視される。この場合の有効数字は10進数のことであるが、コンピュータでは2進数が使用されるので「2進数の有効数字」という意味で「有効ビット」と表記し、また、「その桁数」という意味で「有効ビット数」と今後表記する。

有効ビット数が減少する主要因は近接数同士の減算が行われることである。より正確には、絶対値が近接する同符号二数の減算、あるいは、異符号二数の加算が行われることである。

$$\begin{array}{r} \phantom{-} \phantom{\phantom{0000000000}} \\ - \phantom{\phantom{0000000000}} \\ \hline \phantom{0000000000} \end{array}$$

これらの演算が行われると、有効ビットの上位が削られて下位だけが残り有効ビット数が減少する(上図)。

この減少要因以外に有効ビット数が減少する要因を、通常、省られることはない。一方、有効ビット数が増加する要因については、通常、全く考慮されることがない。

計算精度を決める以上の要因、つまり、「減算だけによる精度減少」により(18×18)対称行列の対角化計算(16倍精度計算)を行い、計算結果の各要素の有効ビット数を算出した。また、この算出結果の精度を確認するために、高精度計算(32倍精度計算)を別途行って得た結果から精度を算出した。両者を比較すると、極めて残念なことに、差が極端に大きかった。実際には、20ビット強しか精度落ちしていないのに、40ビット強の精度落ちがあると算出され(表1の2行目と最終行)、「減算だけによる精度減少」では精度を低く算出してし

まうことが分った。別の言い方をすると、より正確な精度算出には「精度が増す要因」を探す必要があることが分った。

「確実に精度が増す要因」を四則演算だけから探し出すことはできない。考え得るものは「確率的に精度が増す要因」だけである。この要因として、本研究では以下の三つの演算を取り上げ「精度追加」を行った。

1. 減算精度落ち時追加 減算精度落ちにより  $n$  ビット残ったとする。この場合、精度は  $n$  ビットへ減ったと考えるのではなく、さらに下位の1ビット(0か1かのいずれか)が確率的にある程度正しいので、 $(n + \alpha_1)$  ビットの精度へ減ったと考える。この  $\alpha_1$  が「精度追加」である。
2. 加算繰り上り時追加 例として有効ビット数が3ビットの下記2進数の加算を考える。

$$\begin{array}{r} 101 \\ + 100 \\ \hline 1001 \end{array} \quad (1)$$

和の値は繰り上りにより4ビットになっている。この数の有効ビット数は元々の3ビットよりも多くなったと考えてよいだろう。つまり、加算結果が繰り上る時は有効ビット数が  $\alpha_2$  ビットだけ増すと考える。

3. 丸め繰り上り時追加 例として計算結果が下記の値だったとして、これを小数第1位で丸める場合を考える。

$$111.1 \Rightarrow 1000 \quad (2)$$

丸めた値は4ビットになる。この数の有効ビット数が元々の3ビットよりも多くなると考える。つまり、丸め時に繰り上る時は有効ビット数が  $\alpha_3$  ビットだけ増すと考える。なお、この繰り上りの起きる頻度はかなり低い。しかし、一応、有効ビット数が増す要因の一つとして掲げておく。

### 3 【計算方法】

上記の「精度追加」 $\alpha_1, \alpha_2, \alpha_3$  によって有効ビット数を適切に算出することができるかどうかを確認するために、Rn原子の量子化学計算を行ない、重なり積分のHouseholder法による三重対角化における有効ビット数の算出を行なった。

量子化学計算には本研究室においてJava言語を使用して開発した量子化学計算プログラムを使用し、16倍精度(仮数部496ビット)の計算を解析対象データと

表 1: Rn 原子 16 倍精度計算の重なり積分を Householder 法で三重対角化したときの最低精度要素の有効ビット数の算出値と実際の値

$\alpha$ 値 <sup>†</sup>	主対角要素	副対角要素
減 0.00, 加 0.00, 丸 0.00	450	451
減 0.25, 加 0.00, 丸 0.00	459	461
減 0.50, 加 0.00, 丸 0.00	467	468
減 0.75, 加 0.00, 丸 0.00	472	473
減 1.00, 加 0.00, 丸 0.00	476	477
減 0.00, 加 0.00, 丸 0.00	450	451
減 0.00, 加 0.25, 丸 0.00	456	457
減 0.00, 加 0.50, 丸 0.00	461	462
減 0.00, 加 0.75, 丸 0.00	464	464
減 0.00, 加 1.00, 丸 0.00	466	466
減 0.00, 加 0.00, 丸 0.00	450	451
減 0.00, 加 0.00, 丸 0.25	450	451
減 0.00, 加 0.00, 丸 0.50	450	451
減 0.00, 加 0.00, 丸 0.75	450	451
減 0.00, 加 0.00, 丸 1.00	450	451
	472.54 <sup>‡</sup>	473.58 <sup>‡</sup>

<sup>†</sup> この列の 3 つの数は  $\alpha_1, \alpha_2, \alpha_3$  の値を表わす

<sup>‡</sup> 32 倍精度計算から算出した実際の有効ビット数

した。また、正確な値を得るために 32 倍精度 (仮数部 992 ビット) の計算を行ない、これを真値として用いた。なお、このとき数値計算上の誤差だけを解析に使用するため、解析対象データである 16 倍精度の重なり積分の各要素の仮数部下に 496 ビットの 0 を追加することにより得た行列を 32 倍精度計算で用いる 32 倍精度重なり積分とした。

#### 4 【精度追加と有効ビット数】

それぞれの「精度追加」 $\alpha_1, \alpha_2, \alpha_3$  の種々の値と、それを使用して算出された重なり積分三重対角化行列の最低精度の主対角要素と副対角要素の有効ビット数を表 1 に掲げる。

この表の最左列は  $\alpha_1, \alpha_2, \alpha_3$  の値を表わす。値が 0.00 は追加する有効ビット数が零ビットであることを表わし、1.00 は追加する有効ビット数が 1 ビットであることを表わす。表の中央列と最右列はそれぞれ主対角要素と副対角要素の中で最低精度の有効ビット数の算出値である。また、最終行には 32 倍精度計算から得られた当該要素の実際の有効ビット数を掲げる。

表から  $\alpha_1$  や  $\alpha_2$  の値によって算出された有効ビット数が大きく変わることがわかる。一方、 $\alpha_3$  では当初の予想通り算出有効ビット数に影響が現れない。

$\alpha_1$  と  $\alpha_2$  の種々の値の組み合わせによって三重対角化行列の最低精度要素の有効ビット数の算出値がどう変わるをまとめたのが表 2 である。詳細は当日発表する。

表 2:  $\alpha_1$  と  $\alpha_2$  の種々の値の組み合わせによる Rn 原子 16 倍精度計算の重なり積分三重対角化後の最低精度要素の有効ビット数の算出値と実際の値。実際の値に近い  $[\alpha_1, \alpha_2]$  の値は  $[0.25, 0.75](472, 472)$ 、 $[0.25, 1.00](473, 473)$ 、および、 $[0.75, 0.00](472, 473)$  だった。これは、減算桁落ち時追加精度を小さくして加算繰り上がり時追加精度を大きくする (前者 2 つ)、あるいは、減算桁落ち時追加精度を大きくして加算繰り上がり時追加精度をなしにする、と実際に近い値を算出できると解釈できるが、より確実なものにするには一層の比較・検討が必要だろう。

$\alpha$ 値 <sup>†</sup>	主対角要素	副対角要素
減 0.00, 加 0.00	450	451
減 0.00, 加 0.25	456	457
減 0.00, 加 0.50	461	462
減 0.00, 加 0.75	464	464
減 0.00, 加 1.00	466	466
減 0.25, 加 0.00	459	461
減 0.25, 加 0.25	466	467
減 0.25, 加 0.50	469	470
減 0.25, 加 0.75	472	472
減 0.25, 加 1.00	473	473
減 0.50, 加 0.00	467	468
減 0.50, 加 0.25	471	472
減 0.50, 加 0.50	475	475
減 0.50, 加 0.75	476	477
減 0.50, 加 1.00	478	478
減 0.75, 加 0.00	472	473
減 0.75, 加 0.25	476	476
減 0.75, 加 0.50	478	478
減 0.75, 加 0.75	479	480
減 0.75, 加 1.00	480	481
減 1.00, 加 0.00	476	477
減 1.00, 加 0.25	478	479
減 1.00, 加 0.50	480	481
減 1.00, 加 0.75	481	482
減 1.00, 加 1.00	482	483
	472.54 <sup>‡</sup>	473.58 <sup>‡</sup>

<sup>†</sup> この列の 2 つの数は  $\alpha_1$  と  $\alpha_2$  の値を表わす

<sup>‡</sup> 32 倍精度計算から算出した実際の有効ビット数