



02

数値計算における数学

連立一次方程式の求解法を題材として



岩下武史（北海道大学）

計算機で扱う数学



計算科学等の数値シミュレーションは多くの場合、数学や物理学に基づいている。線形代数や微積分といった多くの応用分野で共通的に用いられる数学もあるが、個々の分野において異なった数学的知識が求められるのが普通である。ここで、これらの個別の数学について述べることは、与えられた紙面から見て現実的ではない。そこで、本稿では数値シミュレーションで「数学」を扱う際に共通的に注意すべき点について述べることにする。

数値シミュレーション（数値計算）では、多くの場合（倍精度）浮動小数点数により、さまざまな量が表現される（図-1）。ここで注意しなければならないのは、数値計算で扱われる“数”は誤差を含んでいるという点である。たとえば、倍数度浮動小数点数は64桁の2進数（64ビット）で1つの数を表すが、 π や $\sqrt{2}$ といった無理数はもとより多くの数が近似的に表され、その値には誤差が含まれる。これらの誤差の影響はしばしば計算の進行に伴って拡大し、最終的にシミュレーションの精度に多大な影響を及ぼすことがある。そこで、本稿では多くのシミュレーションにおいて主たる計算核となっている連立一次方程式の求解部を題材に、計算誤差に関連したいくつかの話題について解説する。なお、本稿

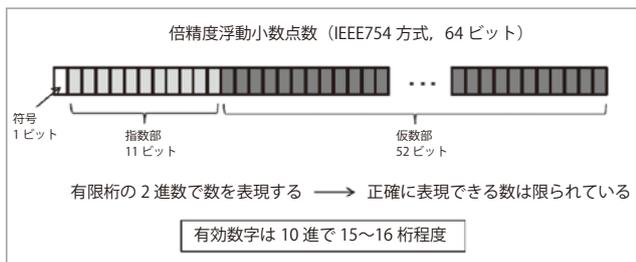


図-1 倍精度浮動小数点数による数の表現

では $n \times n$ の正則な実行列を係数とする実連立一次方程式を扱うものとする。

直接法と反復法



計算機による連立一次方程式の求解法には、直接法と反復法の2種類のカテゴリがある。直接法としてはガウスの消去法やLU分解法等が知られ、行列の変換や代入計算を通じて解ベクトルの要素の各々を直接的に求めていく。一方、初期的に与えられる近似解に対して、ある手順を繰り返し適用することで解を更新し、最終的に十分な精度を持った近似解を得る方法が反復法である。たとえば、ヤコビ法やガウス=ザイデル法、共役勾配法等が反復法としてよく知られている。さて、ここで「反復法では近似解を求めるに過ぎないが、直接法では正しい解を求めることができる」という考えを持つ人がいるが、これは誤った考えである。先に述べたように、計算機上での計算には誤差が含まれるためである。

以下の n 元連立一次方程式について考える。

$$Ax = b \tag{1}$$

ここで、本来の b に代わり、計算機上では誤差の影響で $b + \delta b$ が右辺ベクトルとなったとする。この場合、

$$A(x + \delta x) = b + \delta b \tag{2}$$

と書くことができる。式 (1) (2) より

$$\delta x = A^{-1} \delta b \tag{3}$$

を得る。式 (3) の両辺を $\|b\|$ で割り、さらにノルムの性質 ($\|Ax\| \leq \|A\| \|x\|$) を使って整理すると

$$\frac{\|\delta x\|}{\|x\|} \leq \|A\| \|A^{-1}\| \frac{\|\delta b\|}{\|b\|} \quad (4)$$

を得る。ここで $\kappa(A) = \|A\| \|A^{-1}\|$ は条件数と呼ばれ、式 (4) は右辺ベクトルの相対誤差が小さくても条件数が大きい場合には得られた解に大きな（相対）誤差が含まれる可能性があることを意味している。すなわち、直接法を用いた場合においても、数値や計算の誤差によって通常“正しい解”は求められず、条件数が大きい場合には解ベクトルに大きな誤差が含まれる場合があることになる。条件数が大きい問題は悪条件な問題と呼ばれ、精度の高い解を得ることは解法によらず困難となる¹⁾。

残差による近似解の判定

数値シミュレーションで連立一次方程式を解く場合、解ベクトルがどんな値になるかは事前には分からない。では、得られた解が“十分な精度を持つ”かどうかのように判定すればよいのだろうか。あらかじめ分かっているのは係数行列と右辺ベクトルであるから、通常、以下の残差ベクトル r を用いて判定する²⁾。

$$r = b - A\tilde{x} \quad (5)$$

ここで、 \tilde{x} は計算によって得られた（近似）解である。確かに残差ベクトルが 0 ベクトルとなれば、得られた解は正しい解であると言える。しかし、計算誤差の問題があるため、現実的に残差ベクトルが 0 ベクトルとなることはほとんどない。そのため、反復法では、残差ベクトルのノルムが右辺ベクトルのノルムと比べて十分に小さい場合に（十分な精度を持つ）近似解が得られたとする場合が多い。しかしながら、残差ベクトルが十分に小さかったとしても近似解に大きな誤差が含まれる場合がある。

正しい解 x と近似解 \tilde{x} の間の誤差 e を

$$e = x - \tilde{x} \quad (6)$$

のように表すものとする。このとき、式 (1) (5) (6)

より

$$Ae = r \quad (7)$$

が成り立つ。ここで、誤差ベクトル e が A の固有ベクトルであり、かつその固有値の絶対値が非常に小さいと仮定する。このような場合、 e のノルムがかなり大きな値を持っていたとしても r のノルムは小さい値となる。つまり、このようなケースでは残差ベクトルのノルムを（十分に）小さくしたとしても、依然として大きな誤差が解ベクトルに含まれる可能性が残ることとなる。もう少し詳細に見てみよう。

式 (7) に対して、式 (4) の導出手順と同様に

$$\frac{\|e\|}{\|x\|} \leq \kappa(A) \frac{\|r\|}{\|b\|} \quad (8)$$

を得る。式 (8) は、相対残差（残差ベクトルと右辺ベクトルのノルム比）が小さかったとしても係数行列の条件数が非常に大きい場合、大きな誤差が生じ得ること意味している。つまり、悪条件な問題では相対残差が十分に小さい近似解が得られたとしても安心できないということとなる。

では、残差ベクトルのノルムのほかによい指標がなく、相対残差を使わなければならないとするとどのような基準とすればよいのだろうか？ 残念ながら、この問いに対する一般的な解答はない。それぞれの応用分野において、過去の経験や実験結果との照合等によって適宜決められているのが現状である。ただし、留意すべき事項として、連立一次方程式の求解においてより重要なことは残差ベクトルよりも誤差ベクトルのノルムを小さくすることにあることを挙げるができる。複数の解法を比較する場合には残差ベクトルのノルムの変化だけではなく、あらかじめ解分かっている人工的な問題に対して誤差ベクトルの振舞いをチェックすることも時に重要となる。また、実応用分野では、得られた解が実験結果や物理現象をよく模擬しているかどうか十分に吟味する必要がある。

共役勾配法

連立一次方程式の求解法としてよく利用されている反復法に共役勾配 (CG: Conjugate Gradient) 法がある³⁾。CG法は正値対称な係数行列を持つ連立一次方程式に有効な手法で有限要素解析を始めとした多くの解析で利用されている。また、こうした現状を踏まえて計算機やスーパーコンピュータの性能測定のためのベンチマークとしても活用されている。

CG法の長所は反復ごとに算出される残差ベクトルが互いに直交するという点にある。この性質から、CG法は最大で n 回の反復数で解を導出することができる。つまり、有限の計算量で求解を行うことができるという点からは、直接法と同様の性質を持つということもできる。しかしながらこの性質は誤差のない計算、つまり「紙と鉛筆」による計算の場合であり、数値計算では事情が異なってくる。

CG法の特徴である残差ベクトルの直交性を確保するために、CG法の反復手順には内積演算が含まれる。しかしながら、一般に内積演算の結果は数値誤差の影響を受けやすい。たとえば、第1番目の要素の絶対値が非常に大きく、それ以外の要素の絶対値が小さい、非常に長い N 次元ベクトル f と g の内積を考える。 f および g の i 番目の要素をそれぞれ f_i, g_i と表すと、通常の内積演算の手順では、 f_1 と g_1 の積をまず求め、次に f_2 と g_2 の積をこれに加算し、順次 f_i と g_i の積を計算し、加算する。この場合、 f_1 と g_1 の積、すなわち $f_1 \cdot g_1$ の絶対値が $f_2 \cdot g_2$ の絶対値と比べて非常に大きいと、数値計算上は $f_1 \cdot g_1$ に $f_2 \cdot g_2$ を足しても結果がまったく変わらないということが生じ得る。このような現象が続いた場合、以降の計算は内積演算の結果に反映されないため、内積は $f_1 \cdot g_1$ と計算されることになる。しかしながら、 f と g の要素数が非常に多い場合、いわゆる「ちりも積もれば山となる」で計算結果に反映されなかった $f_2 \cdot g_2 + f_3 \cdot g_3 + \dots$ の (絶対) 値が大きなものとなり、計算結果に大きな誤差が含まれるということが生じ得る。このような内積

演算や他の計算部分による誤差の影響は、しばしば反復が進むにつれて拡大し、CG法における残差ベクトルの直交性を失わせ、近似解が収束しない (解が得られない) という現象を引き起こす要因となる。

前処理

前章でCG法において、内積演算等に起因する計算誤差の影響で求解のプロセスが破綻する場合があることを述べた。このような現象は他の反復法においても見られるもので、実応用上重要な問題である。本問題の対処法として広く用いられているのが前処理と呼ばれる技術である。前処理は反復法の収束性を向上させる手法で、主に求解時間の短縮を目的として利用されている。しかし、反復ごとに増大する計算誤差の影響が深刻化する前に求解プロセスを完了するという点から言えば、反復法における計算誤差の対処法の1つと捉えることもできる。

前処理とは連立一次方程式 (1) を解く代わりに

$$M_1^{-1} A M_2^{-1} (M_2 x) = M_1^{-1} b \quad (9)$$

を解く方法である。 M_1 および M_2 (前処理行列) をうまく選ぶことにより、反復法の収束性を向上させることができる。あいまいな表現であるが、なんらかの意味で前処理後の係数行列 $\tilde{A} = M_1^{-1} A M_2^{-1}$ が単位行列に近い場合、反復法の収束性を改善することができる。つまり $M_1 M_2 \approx A$ となるような M_1, M_2 を選ぶことができれば、 $\tilde{A} \approx I$ (I : 単位行列) となり、前処理としての効果が期待できる。

CG法の場合、前処理後の係数行列 \tilde{A} の対称性を確保するために、 $M_2^{-1} = (M_1^{-1})^T = (M_1^T)^{-1}$ 、すなわち、 $M_2 = M_1^T$ とするのが一般的である。CG法の前処理に関する手順は各反復中で、

$$Mz = r \quad (10)$$

を解くことにより与えられる。ここで、 $M = M_1 M_1^T$ である。したがって、前処理行列には、式 (10) が式 (1) よりも簡単に (より計算量が少なく) 解けることが求められる。また、 M_1 の導出に多大な計算量を必

