

多項式の減次†

平野 菅 保††

浮動小数点演算を用いて代数方程式の解を数値的に求める場合、与えられる代数方程式の係数は有限桁の数値である。そこで、まず近似解が有限桁の数値を係数にもつ代数方程式を満足するための条件を明白にした。有限桁の数値を係数にもつ与えられる代数方程式を満足する既知の近似解を用いて、次数を一つ減じた代数方程式を求める場合、剰余の定理を用いて、多項式の減次を行い、浮動小数点演算では0になるとは限らない、剰余の定数を0としている。このため、代数方程式の係数として与えられる数値の桁数に比較して多くの桁数を用いて数値計算を行わないと、次数を一つ減じた代数方程式から得られる、多項式の減次に用いた既知の近似解よりも絶対値が小さい近似解は、与えられる代数方程式を必ずしも満足しないことを明白にした。そこで、多項式の減次を行う場合、数値計算の途中で得られる数値を用いて、剰余の項の x に関する次数を適切に定めると、多項式の減次に用いた既知の近似解が有限桁の数値を係数にもつ与えられる代数方程式を満足していれば、その代数方程式の係数として与えられる数値の桁数に数桁加えた桁数で数値計算を行っても、次数を一つ減じた代数方程式から、与えられる代数方程式を満足する近似解を求めることができることを説明した。

1. ま え が き

代数方程式の解を数値的に一つずつ順次求め、すべての解を求めることは予想外にむずかしい問題である。これは本来、理論的には数値を無限桁の小数展開として扱わねばならないのに、数値計算をする場合、その数値を有限桁の小数として扱うからである。

とくに電子計算機を使用して数値計算をする場合、用いる数値の桁数を一定にした、いわゆる浮動小数点演算で行われるために、求められた近似解を用いて多項式の次数を減ずる、多項式の減次の数値計算の手順^{1)~3)}を注意深く配慮しないと、簡単な加算の演算段階で、意外に大きな情報損失の現象を生じる。したがって、求められる解がとんでもない間違った数値をとることがある。

2. 有限桁の実係数と実数解

本論文で採用する代数方程式は

$$f_{Tn}(x) = \sum_{i=0}^n a_{Ti} \cdot x^i = 0 \quad a_{Tn} \neq 0 \quad (1)$$

ただし、 $a_{Ti} (i=0, 1, \dots, n)$ は実数であり、解 $x_{Ti} (i=1, 2, \dots, n)$ は実数である。

であり、誤差を含まない真値を表す実数の記号に、すべて a_{Ti}, x_{Ti} のように記号 a, x の右下に True 「真」の意味を表す T をつけている。また、誤差を含まない真値を表す実係数 $a_{Ti} (i=0, 1, \dots, n)$ のみをもつ関数

にも $f_{Tn}(x)$ のように f の右下に T をつけている。

実際に電子計算機で使用している浮動小数点演算で数値計算をする場合、代数方程式(1)の0でない係数 a_{Ti} はそれぞれ有限桁の数値で与えられるので、それぞれ与えられる数値の最下位の桁の次の桁から下位の桁の真の数値は不明である。したがって、不明の桁をすべて0として得られる実数が与えられる代数方程式

$$f_n(x) = \sum_{i=0}^n a_i \cdot x^i = 0 \quad a_n \neq 0 \quad (2)$$

ただし、 $a_i (i=0, 1, \dots, n)$ は実数である。
の0でない係数 a_i である。

代数方程式(2)が代数方程式(1)を代表するためには、代数方程式(2)の0でない係数 a_i が

$$a_{Ti} = a_i + \varepsilon_{Ti} \\ -0.5 \cdot u_i \leq \varepsilon_{Ti} \leq 0.5 \cdot u_i$$

ただし、 u_i は0でない係数 a_i として与えられる数値の最下位の桁の単位である。

を満足することのみ既知である代数方程式(1)の0でない係数 a_{Ti} を代表していなければならない。すなわち、0である係数 a_i をそのまま係数にもち、0でない係数 a_i には

$$-0.5 \cdot u_i \leq \varepsilon_i \leq 0.5 \cdot u_i \quad (3)$$

ただし、 $a_i = 0$ ならば $u_i = 0$ とする。
を満足する任意の誤差 ε_i を加えた実数 $(a_i + \varepsilon_i)$ を係数にもつ代数方程式

$$\sum_{i=0}^n (a_i + \varepsilon_i) \cdot x^i = 0 \quad (4)$$

の解 x'_i よりも、代数方程式(1)の解に近接していると認めることのできる近似解を代数方程式(2)から求めることは、一般にはできない。そこで、本論文で

† Polynomial Deflation by SUGAYASU HIRANO (Department of Mathematical Engineering, College of Industrial Technology, Nihon University).

†† 日本大学生産工学部数理工学科

は、その解 x_j を代数方程式(1)を代表する代数方程式(2)を満足する近似解とする。

2.1 関数値と近似解

代数方程式(4)の解 x_j を代数方程式(2)に代入すると、関数値は

$$f_n(x_j) = \sum_{i=0}^n a_i \cdot x_j^i = - \sum_{i=0}^n \varepsilon_i \cdot x_j^i$$

となり、0になるとは限らない。すなわち、適当な数値計算で代数方程式(2)の近似解 x_j を求め、その近似解 x_j を代数方程式(2)に代入して得られる関数値 $f_n(x_j)$ が条件(3)を満足する適当な誤差 ε_i を用いて

$$f_n(x_j) = \sum_{i=0}^n a_i \cdot x_j^i = - \sum_{i=0}^n \varepsilon_i \cdot x_j^i \quad (5)$$

とできるならば、近似解 x_j は代数方程式(2)を満足する近似解である。

3. 残 差

「代数方程式(2)の近似解が代数方程式(2)を満足している」と判定するためには、(5)の右辺の値、すなわち、残差の誤差を数値計算で求める。

そこで、閉区間(3)に含まれる任意の誤差 ε_i に、代数方程式(2)の近似解 x_j の i 乗、 x_j^i を乗じると $-0.5 \cdot u_i \cdot |x_j^i| \leq \varepsilon_i \cdot x_j^i \leq 0.5 \cdot u_i \cdot |x_j^i|$ が得られる。すなわち、

$$|\varepsilon_i \cdot x_j^i| \leq |\Delta a_i \cdot x_j^i| \quad (6)$$

ただし、 $\Delta a_i = 0.5 \cdot u_i (i=0, 1, \dots, n)$ であり、 $\Delta a_i \neq 0$ である Δa_i を、代数方程式(2)の 0 でない係数 a_i として与えられる数値の最大誤差とよぶ。

である。したがって、代数方程式(2)の近似解 x_j を代数方程式(2)に代入して

$$|f_n(x_j)| \leq \sum_{i=0}^n |\Delta a_i \cdot x_j^i| \quad (7)$$

が成り立てば、この近似解 x_j は条件(5)から代数方程式(2)を満足する。また、不等式(7)の右辺の値を関数値 $f_n(x_j)$ の最大誤差とよび

$$\Delta f_n(x_j)$$

で表す。

3.1 最大影響係数と判定条件

代数方程式(2)を満足する近似解 x_j を代数方程式(2)に代入したとき、求められる関数値 $f_n(x_j)$ の最大誤差 $\Delta f_n(x_j)$ に最も影響を与える誤差項

$$|\Delta a_{\max j} \cdot x_j^{\max j}| = \max_{i=0,1,\dots,n} (|\Delta a_i \cdot x_j^i|) \quad (8)$$

から得られる係数 $a_{\max j}$ を、近似解 x_j の最大影響係数とよぶ。

ここで、近似解 x_j の有効桁数を大きくするために、式(8)で選ばれた x に関して $\max j$ 次の係数 $a_{\max j}$ の有効桁数、すなわち、係数 $a_{\max j}$ として与えられる数値の正しい数値の桁数を、一般には大きくする。

最大影響係数を求める式(8)から

$$|\Delta a_{\max j} \cdot x_j^{\max j}| \leq \sum_{i=0}^n |\Delta a_i \cdot x_j^i| \quad (9)$$

が得られるので、「代数方程式(2)の近似解 x_j が代数方程式(2)を満足している」と、数値計算上容易に判定するため、これ以降、判定条件として条件(7)の代りに

$$|f_n(x_j)| \leq |\Delta a_{\max j} \cdot x_j^{\max j}| \quad (10)$$

を用いる。

4. 剰余の定理

n 次多項式 $f_{Tn}(x)$ (1)は剰余の定理から

$$f_{Tn}(x) = (x - x_r) \cdot f_{T(n-1)}(x) + f_{Tn}(x_r) \quad (11)$$

ただし、 x_r は実数である。

$$f_{T(n-1)}(x) = \sum_{i=0}^{n-1} \bar{b}_{Ti} \cdot x^i \quad (12)$$

$$\bar{b}_{Tn-1} = a_{Tn}$$

$$\bar{b}_{Ti} = \bar{b}_{T(i+1)} \cdot x_r + a_{T(i+1)}$$

$$= \sum_{j=i+1}^n a_{Tj} \cdot x_r^{j-(i+1)}$$

$$(i = n-2, n-3, \dots, 1, 0)$$

$$\bar{b}_{T0} \cdot x_r + a_{T0} = \sum_{j=0}^n a_{Tj} \cdot x_r^j = f_{Tn}(x_r)$$

と表すことができる。

$(n-1)$ 次多項式 $f_{T(n-1)}(x)$ (12)の係数 $\bar{b}_{Ti} (i = n-1, n-2, \dots, 1, 0)$ は n 次多項式 $f_{Tn}(x)$ (1)を高次項より順次 $(x - x_r)$ で除して得られ、剰余としては、 x に関して 0 次の $f_{Tn}(x_r)$ が得られる。

実数 x_r が n 次代数方程式(1)の解であれば、

$$f_{Tn}(x_r) = 0$$

であるから次数を一つ減じた $(n-1)$ 次代数方程式

$$f_{T(n-1)}(x) = \sum_{i=0}^{n-1} \bar{b}_{Ti} \cdot x^i = 0 \quad (13)$$

が剰余の定理の公式(11)から得られる。

4.1 k 次の剰余

n 次多項式 $f_{Tn}(x)$ (1)を高次項および低次項の両方から順次 $(x - x_r)$ で除し、剰余としては、 x に関して k 次の

$$f_{Tn}(x_r) \cdot (x/x_r)^k$$

を得ると、剰余の定理の公式(11)の代りに公式

$$f_{T_n}(x) = (x - x_r) \cdot f_{T_{(k)n-1}}(x) + f_{T_n}(x_r) \cdot (x/x_r)^k \quad (14)$$

$$f_{T_{(k)n-1}}(x) = \sum_{i=0}^{k-1} b_{T_i} \cdot x^i + \sum_{i=k}^{n-1} \bar{b}_{T_i} \cdot x^i \quad (15)$$

$$b_{T_0} = -a_{T_0}/x_r$$

$$b_{T_i} = -(a_{T_i} - b_{T_{i-1}})/x_r$$

$$= -\sum_{j=0}^i a_{T_j} \cdot x_r^{j-(i+1)}$$

$$(i=1, 2, \dots, k-1)$$

$$(\bar{b}_{T_k} \cdot x_r + a_{T_k} - b_{T_{k-1}}) x^k$$

$$= \left(\sum_{j=k+1}^n a_{T_j} \cdot x_r^{j-k} + a_{T_k} + \sum_{j=0}^{k-1} a_{T_j} \cdot x_r^{j-k} \right) \cdot x^k$$

$$= f_{T_n}(x_r) \cdot (x/x_r)^k$$

が得られる。

実数 x_r が n 次代数方程式(1)の解であれば、次数を一つ減じた $(n-1)$ 次代数方程式

$$f_{T_{(k)n-1}}(x) = \sum_{i=0}^{k-1} b_{T_i} \cdot x^i + \sum_{i=k}^{n-1} \bar{b}_{T_i} \cdot x^i = 0 \quad (16)$$

が得られ、 $(n-1)$ 次代数方程式(13)と等しくなる。

したがって、係数 \bar{b}_i (12)と係数 b_i (15)とは

$$\bar{b}_i = b_i \quad (i=0, 1, \dots, k-1)$$

となる。

5. 浮動小数点演算

浮動小数点の数値を用いて数値計算を行うとき、その数値の仮数の桁数を固定する。この固定された桁数を使用桁数とよぶことにする。

数値計算では、段階に応じて使用桁数を変える場合もある。しかし、一段階の数値計算では、その数値計算の全部を通じて使用桁数を固定する。また、使用桁数を変える場合、たとえば、どこかの段階はとくに精密に数値計算をしなければならない場合、その段階のみ大きい使用桁数を用いる。

一方、数値計算では、演算ごとに丸め誤差が結果の数値の最下位の桁に入るので、桁落ちによる桁落ち誤差を除いても、累積丸め誤差に影響されない桁数は使用桁数より数桁小さい。しかし、つねに、この累積丸め誤差を考慮すると複雑になるので、実際に用いられる使用桁数より数桁小さい桁数を、理論的考察をする場合に用いる指定桁数とよび、その指定桁数以内には、桁落ち誤差以外の累積丸め誤差が入らないとする。

そこで、代数方程式(2)の係数 a_i のなかで、0で

ない係数として与えられる数値の与えられる桁数を、それぞれとりあげ、その桁数のなかで最大の桁数を指定桁数とする。

また、桁落ち現象が起こる加減算では、被演算数、演算数、それぞれが含む累積丸め誤差は演算結果に入る。しかし、加減算そのものによる丸め誤差は演算結果には入らない。したがって、演算結果に入っている数値計算による累積丸め誤差の絶対値は、演算結果に入っている、代数方程式(2)の係数 a_i のなかで、0でない係数として与えられる数値が含む誤差(最大誤差 Δa_i)による誤差の絶対値よりも数桁小さい。

そこで、これ以降、理論的考察をする場合、累積丸め誤差を考慮しない。

5.1 指定桁数の増加

代数方程式(2)の係数 a_i ($i=0, 1, \dots, n$) のなかで、0でない係数が有限桁の数値として与えられるとき、それらの数値の与えられる桁の最下位の桁の次の桁から下位の桁へ適当な桁までそれぞれ0を加え、0である係数として与えられる数値とともに、改めて、係数 \bar{a}_i ($i=0, 1, \dots, n$) として与えたとする。

$$f_n(x) = \sum_{i=0}^n \bar{a}_i \cdot x^i = 0 \quad \bar{a}_n \neq 0 \quad (17)$$

代数方程式(17)を満足する近似解 \bar{x}_j を代数方程式(17)に代入すると、条件(3)と同様に

$$-0.5 \cdot \bar{u}_i \leq \bar{\varepsilon}_i \leq 0.5 \cdot \bar{u}_i$$

ただし、 \bar{u}_i は0でない係数 \bar{a}_i として与えられる数値の最下位の桁の単位であり、 u_i (3)との関係は

$$u_i \geq \bar{u}_i \quad (18)$$

である。また、 $\bar{a}_i = 0$ ならば $\bar{u}_i = 0$ とする。

を満足する適当な $\bar{\varepsilon}_i$ を用いて式(5)と同様に

$$f_n(\bar{x}_j) = \sum_{i=0}^n \bar{a}_i \cdot \bar{x}_j^i = -\sum_{i=0}^n \bar{\varepsilon}_i \cdot \bar{x}_j^i \quad (19)$$

が得られる。一方、代数方程式(2)の係数 a_i と代数方程式(17)の係数 \bar{a}_i とはともに実数であり、

$$a_i = \bar{a}_i \quad (i=0, 1, \dots, n)$$

であるから、代数方程式(17)を満足する近似解 \bar{x}_j を代数方程式(2)に代入すると

$$f_n(\bar{x}_j) = \sum_{i=0}^n a_i \cdot \bar{x}_j^i = \sum_{i=0}^n \bar{\varepsilon}_i \cdot \bar{x}_j^i \quad (20)$$

となる。したがって、式(19)と式(20)とから

$$f_n(\bar{x}_j) = \sum_{i=0}^n a_i \cdot \bar{x}_j^i = -\sum_{i=0}^n \bar{\varepsilon}_i \cdot \bar{x}_j^i$$

が成り立つ。また、条件(18)を用いれば、次の定理が成り立つ。

[定理] 代数方程式(17)を満足する近似解は, 代数方程式(2)を満足する.

6. 多項式の減次

剰余の定理により n 次多項式 $f_n(x)$ (2) から

$$f_n(x) = (x-x_1) \cdot f_{(0)n-1}(x) + f_n(x_1) \quad (21)$$

ただし, x_1 は条件(10) ($j=1$) を満足する代数方程式(2)の近似解である.

が得られ, 近似解 x_1 は条件(10)を満足するから

$$|f_n(x_1)| \leq |\Delta a_{\max 1} \cdot x_1^{m+1}| \quad (22)$$

である. つづいて, 次に求める近似解 x_2 が条件(10) ($j=2$) を満足するためには

$$|f_n(x_2)| = |(x_2-x_1) \cdot f_{(0)n-1}(x_2) + f_n(x_1)| \leq |\Delta a_{\max 2} \cdot x_2^{m+2}| \quad (23)$$

でなければならない. したがって, 次数を一つ減じた代数方程式

$$f_{(0)n-1}(x) = 0 \quad (24)$$

から得られる近似解 x_2 が, 必ず条件(23)を満足するためには, 少なくとも

$$|f_n(x_1)| \leq |\Delta a_{\max 2} \cdot x_2^{m+2}| \quad (25)$$

を満足していなければならない.

$|x_2| \geq |x_1|$ の場合, 式(8)から

$$|\Delta a_{\max 1} \cdot x_1^{m+1}| \leq |\Delta a_{\max 1} \cdot x_2^{m+1}| \leq |\Delta a_{\max 2} \cdot x_2^{m+2}| \quad (26)$$

であるから, 条件(22)を用いると, 条件(25)を満足しており, 次数を一つ減じた代数方程式(24)から, 条件(23)を満足する近似解 x_2 を求めることができる. すなわち, 代数方程式(24)から代数方程式(2)を満足する近似解 x_2 を求めることができる.

$|x_2| < |x_1|$ の場合, 式(8)から

$$|\Delta a_{\max 2} \cdot x_2^{m+2}| < |\Delta a_{\max 2} \cdot x_1^{m+2}| \leq |\Delta a_{\max 1} \cdot x_1^{m+1}| \quad (27)$$

であるから, 条件(22)を用いても, 必ずしも, 条件(25)を満足するとは限らない. すなわち, 次数を一つ減じた代数方程式(24)から代数方程式(2)を満足する近似解 x_2 を, 必ずしも求められない.

計算例

$$\begin{aligned} f_{T3}(x) &= (x-\pi \cdot 10^4) \cdot (x-\pi \cdot 10^2) \cdot (x-\pi) \\ &= x^3 + a_{T2} \cdot x^2 + a_{T1} \cdot x + a_{T0} = 0 \end{aligned} \quad (1)$$

$$f_3(x) = x^3 + a_2 \cdot x^2 + a_1 \cdot x + a_0 = 0 \quad (2)$$

$$a_2 = -3.173322700 \times 10^4$$

$$\Delta a_2 = 5.0 \times 10^{-4}$$

$$a_1 = 9.969287400 \times 10^6$$

$$\Delta a_1 = 5.0 \times 10^{-2}$$

$$a_0 = -3.100627700 \times 10^7$$

$$\Delta a_0 = 5.0 \times 10^{-1}$$

係数 a_2, a_1, a_0 は 8 桁の数値で与えられ, 指定桁数を 8 とし, 使用桁数を 10 とする.

代数方程式②に含まれる三つの解のなかで, 絶対値が 2 番目に大きく, 条件(22)を満足する近似解

$$x_1 = 3.141592650 \times 10^2 \quad (3)$$

$$f_3(x_1) = 3.985000000 \times 10 \quad (4)$$

を得たとすると, 次数を一つ減じた 2 次方程式

$$f_{(0)2}(x) = x^2 + \bar{b}_1 \cdot x + \bar{b}_0 = 0 \quad (5)$$

$$\bar{b}_1 = -3.141906774 \times 10^4$$

$$\bar{b}_0 = 9.869617200 \times 10^4$$

が得られ, この 2 次方程式⑤から得られる, 3 次方程式②の二つの近似解

$$x_2 = 3.141592615 \times 10^4 \quad (6)$$

$$x_3 = 3.141596766$$

を 3 次方程式②に代入すると,

$$f_3(x_2) = 1.157795000 \times 10^4 \quad (7)$$

$$f_3(x_3) = 3.984000000 \times 10$$

となる. また, 近似解 x_1 ③, x_2 ⑥, x_3 ⑥の最大影響係数による誤差は, それぞれ式(8)から

$$|\Delta a_{\max 1} \cdot x_1^{m+1}| = |\Delta a_2 \cdot x_1^2| = 4.93480 \times 10$$

$$|\Delta a_{\max 2} \cdot x_2^{m+2}| = |\Delta a_2 \cdot x_2^2| = 4.93480 \times 10^5 \quad (8)$$

$$|\Delta a_{\max 3} \cdot x_3^{m+3}| = |\Delta a_0| = 5.00000 \times 10^{-1}$$

である. したがって, 関数値⑦と最大影響係数による誤差⑧とから, 絶対値が近似解 x_1 ③より大きい近似解 x_2 ⑥は条件(10) ($j=2$)

$$|f_3(x_2)| \leq |\Delta a_{\max 2} \cdot x_2^{m+2}| = |\Delta a_2 \cdot x_2^2|$$

を満足しているが, 絶対値が近似解 x_1 ③より小さい近似解 x_3 ⑥は条件(10) ($j=3$)

$$|f_3(x_3)| \leq |\Delta a_{\max 3} \cdot x_3^{m+3}| = |\Delta a_0|$$

を満足していない. すなわち, 近似解 x_3 ⑥の関数値 $f_3(x_3)$ ⑦と近似解 x_1 ③の関数値 $f_3(x_1)$ ④とは

$$f_3(x_3) \doteq f_3(x_1)$$

であり, 2 次方程式⑤から近似解 x_3 を求めている限り, 近似解 x_3 を求める数値計算に関係ない関数値 $f_3(x_1)$ ④の絶対値は小さくならないので, 不等式(23)からわかるように, 必ずしも

$$|f_3(x_3)| = |(x_3-x_1) \cdot f_{(0)2}(x_3) + f_3(x_1)|$$

$$\leq |\Delta a_0| = |\Delta a_{\max 3} \cdot x_3^{m+3}|$$

の不等式を満足する近似解 x_3 を求められない.

6.1 max 1 次の剰余

n 次多項式 $f_n(x)$ (2) を高次項および低次項の両方

から順次 $(x-x_1)$ で除し、剰余としては、

$$f_n(x_1) \cdot (x/x_1)^{\max 1}$$

ただし、 x_1 は条件(22)を満足する代数方程式(2)の近似解であり、 $\max 1$ は近似解 x_1 の最大影響係数を含む項の x に関する次数である。

を得ると、剰余の定理の公式(21)の代りに公式

$$f_n(x) = (x-x_1) \cdot f_{(\max 1)n-1}(x) + f_n(x_1) \cdot (x/x_1)^{\max 1} \quad (28)$$

が得られる。つづいて、次に求める近似解 x_2 が条件(10) ($j=2$)を満足するためには、

$$|f_n(x_2)| = |(x_2-x_1) \cdot f_{(\max 1)n-1}(x_2) + f_n(x_1) \cdot (x_2/x_1)^{\max 1}| \leq |\Delta a_{\max 2} \cdot x_2^{\max 2}| \quad (29)$$

でなければならない。したがって、次数を一つ減じた代数方程式

$$f_{(\max 1)n-1}(x) = 0 \quad (30)$$

から得られる近似解 x_2 が、必ず条件(29)を満足するためには、すくなくとも

$$|f_n(x_1) \cdot (x_2/x_1)^{\max 1}| \leq |\Delta a_{\max 2} \cdot x_2^{\max 2}| \quad (31)$$

を満足していなければならない。

近似解 x_1 が満足している条件(22)から得られる

$$|f_n(x_1) \cdot (x_2/x_1)^{\max 1}| \leq |\Delta a_{\max 1} \cdot x_1^{\max 1} \cdot (x_2/x_1)^{\max 1}|$$

と等式

$$\Delta a_{\max 1} \cdot x_1^{\max 1} \cdot (x_2/x_1)^{\max 1} = \Delta a_{\max 1} \cdot x_2^{\max 1}$$

および式(8)から

$$|f_n(x_1) \cdot (x_2/x_1)^{\max 1}| \leq |\Delta a_{\max 1} \cdot x_2^{\max 1}| \leq |\Delta a_{\max 2} \cdot x_2^{\max 2}|$$

が得られ、条件(31)を満足しており、次数を一つ減じた代数方程式(30)から、条件(29)を満足する近似解 x_2 を求めることができる。すなわち、代数方程式(30)から代数方程式(2)を満足する近似解 x_2 を求めることができる。

計算例

条件(22)を満足する近似解 x_1 ③が得られたとすると、最大影響係数による誤差⑧から近似解 x_1 ③の最大影響係数 $a_{\max 1}$ を含む項の x に関する次数は

$$\max 1 = 2$$

であるから、次数を一つ減じた2次方程式

$$f_{(2)2}(x) = x^2 + b_1 \cdot x + b_0 = 0 \quad (9)$$

$$b_1 = -3.141906814 \times 10^4$$

$$b_0 = 9.869604514 \times 10^4$$

が得られ、この2次方程式⑨から得られる、3次方程式②の二つの近似解

$$x_2 = 3.141592655 \times 10^4 \quad (10)$$

$$x_3 = 3.141592688$$

を3次方程式②に代入すると

$$f_3(x_2) = 4.023924700 \times 10^5 \quad (11)$$

$$f_3(x_3) = 0(\epsilon) \approx 10^{-3}$$

ただし、関数値 $f_3(x_3)$ は10桁の数値計算では0となるが、 a_0 ②の使用桁数10から考えて、 a_0 ②の11桁目の単位 10^{-3} とした。

となる。関数値⑩と最大影響係数による誤差⑧とから、二つの近似解 x_2 ⑩、 x_3 ⑩はともに条件(10) ($j=2, 3$)

$$|f_3(x_2)| \leq |\Delta a_{\max 2} \cdot x_2^{\max 2}| = |\Delta a_2 \cdot x_2^2|$$

$$|f_3(x_3)| \leq |\Delta a_{\max 3} \cdot x_3^{\max 3}| = |\Delta a_0|$$

を満足している。

6.2 指定桁数の増加

代数方程式(2)の係数 a_i ($i=0, 1, \dots, n$)のうち、0でない係数として与えられる数値の与えられる桁数

$$l_i = [\log_{10} \{|a_i|/(2 \cdot \Delta a_i)\}] + 1 \quad (32)$$

ただし、 $[\]$ は小数点以下を切り捨てるガウス記号であり、 $\Delta a_i \neq 0$ とする。

と

$$l'_i = [\log_{10} \{|(a_i \cdot x_1^i)/(2 \cdot \Delta a_0)\}] + 2 \quad (33)$$

ただし、 x_1 は条件(10) ($j=1$)を満足する代数方程式(2)の近似解であり、 $\Delta a_0 \neq 0$ とする。

とを、 $i=0$ の場合を除いて、比較し

$$l_i < l'_i, \quad i \neq 0 \quad (34)$$

である場合、条件(34)を満足する、それぞれの数値の (l_i+1) 桁目から l'_i 桁目まで0を加え、条件(34)を満足しない数値と0である係数として与えられる数値および係数 a_0 として与えられる数値はそのままにし、改めて、それぞれの数値を係数 \bar{a}_i ($i=0, 1, \dots, n$)として与えたとする。

$$f_n(x) = \sum_{i=0}^n \bar{a}_i \cdot x^i = 0, \quad \bar{a}_n \neq 0 \quad (35)$$

5.1 節の指定桁数の増加における定理により、代数方程式(35)を満足する近似解は必ず代数方程式(2)を満足する。

代数方程式(35)の係数 \bar{a}_i ($i=0, 1, \dots, n$)のなかで、係数 a_0 を除き、0でない係数として与えられる数値の与えられる桁数 \bar{l}_i は、それぞれ

$$\bar{l}_i \geq l'_i, \quad i \neq 0 \quad (36)$$

であり、桁数の最大値

$$\bar{l}_{\max} = \max_{i=0, 1, \dots, n} (\bar{l}_i)$$

ただし、 $\bar{l}_0 = l_0$ である。

が指定桁数となる。

代数方程式(35)の係数 $\bar{a}_i (i=0, 1, \dots, n)$ のなかで、係数 \bar{a}_0 を除き、0でない係数として与えられる数値の最大誤差は、それぞれ

$$\Delta \bar{a}_i \leq |\bar{a}_i/2| \cdot 10^{-l_i+1} \quad i \neq 0 \quad (37)$$

であり、桁数を求める式 l_i (33) および不等式(36)から、 $i=0$ の場合を除き

$$10^{-l_i+1} \leq 10^{-l_i+1} < |(2 \cdot \Delta a_0)/(a_i \cdot x_1^i)|$$

が得られ、これを(37)に代入すると

$$\Delta \bar{a}_i < |\bar{a}_i/2| \cdot |(2 \cdot \Delta a_0)/(a_i \cdot x_1^i)| \quad i \neq 0 \quad (38)$$

となる。一方、代数方程式(2)の係数 a_i と代数方程式(35)の係数 \bar{a}_i とは、ともに実数であり、

$$a_i = \bar{a}_i \quad (i=0, 1, \dots, n) \quad (39)$$

である。また、代数方程式(2)と代数方程式(35)の0次の係数 a_0, \bar{a}_0 として与えられる、それぞれ二つの数値の与えられる桁数 l_0, \bar{l}_0 は等しく

$$\Delta a_0 = \Delta \bar{a}_0 \quad (40)$$

であるから、 a_i (39)、 Δa_0 (40)を不等式(38)に代入し、

$$|\Delta \bar{a}_i \cdot x_1^i| < \Delta \bar{a}_0 \quad (i=1, 2, \dots, n) \quad (41)$$

が得られる。

ここで、条件(10) ($j=1$) を満足する代数方程式(2)の近似解 x_1 と、条件(10) ($x_j = \bar{x}_1$) を満足する代数方程式(35)の近似解 \bar{x}_1 とを、 $i=0$ の場合を除き、ともに桁数を求める式 l_i (33) に代入して得られる二つの数値がそれぞれ等しい場合はそのままよいが、異なった二つの数値がある場合には、改めて、近似解 x_1 の代りに近似解 \bar{x}_1 を桁数を求める式 l_i (33) ($x_1 = \bar{x}_1$) に代入し、それ以降の数値計算をやり直したがつて

$$x_1 \rightarrow \bar{x}_1$$

と変えて、不等式(41)を用いると、条件(10) ($x_j = \bar{x}_1$) を満足する代数方程式(35)の近似解 \bar{x}_1 は

$$|f_n(\bar{x}_1)| \leq |\Delta \bar{a}_0| = \max_{i=0, 1, \dots, n} (|\Delta \bar{a}_i \cdot \bar{x}_1^i|) \quad (42)$$

の関係にある。

2番目に求める代数方程式(35)の近似解 \bar{x}_2 を、次数を一つ減じた代数方程式

$$f_{(0)n-1}(x) = 0 \quad (43)$$

から求めると、

$$\begin{aligned} |\Delta \bar{a}_0| &\leq |\Delta \bar{a}_{\max 2} \cdot \bar{x}_2^{\max 2}| \\ &= \max_{i=0, 1, \dots, n} (|\Delta \bar{a}_i \cdot \bar{x}_2^i|) \end{aligned}$$

であるから、条件(42)から

$$|f_n(\bar{x}_1)| \leq |\Delta \bar{a}_{\max 2} \cdot \bar{x}_2^{\max 2}|$$

を満足する。すなわち、代数方程式 $f_{(0)n-1}(x) = 0$ (43) から

$$\begin{aligned} |f_n(\bar{x}_2)| &= |(\bar{x}_2 - \bar{x}_1) \cdot f_{(0)n-1}(\bar{x}_2) + f_n(\bar{x}_1)| \\ &\leq |\Delta \bar{a}_{\max 2} \cdot \bar{x}_2^{\max 2}| \end{aligned}$$

を満足する近似解 \bar{x}_2 を求めることができる。

このように、多項式の減次の数値計算の手順を変えずに、数値計算の途中で得られる数値を用いて、指定桁数を改めて求め、すなわち、指定桁数を増加させ、次数を一つ減じた代数方程式(43)から代数方程式(2)を満足する近似解 x_2 を求めることができる。

計算例

代数方程式②の係数 a_2, a_1 として与えられる数値、係数 a_0 として与えられる数値の最大誤差 Δa_0 および条件(10) ($j=1$) を満足する代数方程式(2)の近似解 x_1 を、桁数を求める式(33)に代入し、指定桁数を求めると

$$l_{\max} = l'_2 = \left[\log_{10} \left\{ \left| \frac{a_2 \cdot x_1^2}{2 \cdot \Delta a_0} \right| \right\} \right] + 2 = 11$$

ただし、 $|a_2 \cdot x_1^2| = 3.13 \times 10^9$, $\Delta a_0 = 5.0 \times 10^{-1}$ である。

であり、使用桁数を 13 とすると

$$f_3(x) = x^3 + \bar{a}_2 \cdot x^2 + \bar{a}_1 \cdot x + \bar{a}_0 = 0 \quad (12)$$

$$\bar{a}_2 = -3.173322700000 \times 10^4$$

$$\Delta \bar{a}_2 = 5.0 \times 10^{-7}$$

$$\bar{a}_1 = 9.969287400000 \times 10^6$$

$$\Delta \bar{a}_1 = 5.0 \times 10^{-5}$$

$$\bar{a}_0 = -3.100627700000 \times 10^7$$

$$\Delta \bar{a}_0 = 5.0 \times 10^{-1}$$

が得られる。条件(42)を満足する近似解 \bar{x}_1

$$\bar{x}_1 = 3.141592692100 \times 10^2$$

$$f_3(\bar{x}_1) = -4.343900000000 \times 10^{-1}$$

を求めると、次数を一つ減じた2次方程式

$$f_{(0)2}(x) = x^2 + \bar{b}_1 \cdot x + \bar{b}_0 = 0 \quad (13)$$

$$\bar{b}_1 = -3.141906773079 \times 10^4$$

$$\bar{b}_0 = 9.869604243600 \times 10^4$$

が得られる。ここで、この2次方程式⑬の係数 \bar{b}_1 および \bar{b}_0 として与えられる数値の使用桁数をそれぞれ 10 とし得られる2次方程式

$$f_{(0)2}^*(x) = x^2 + \bar{b}_1^* \cdot x + \bar{b}_0^* = 0 \quad (14)$$

$$\bar{b}_1^* = -3.141906773 \times 10^4$$

$$\bar{b}_0^* = 9.869604244 \times 10^4$$

から得られる、3次方程式②の二つの近似解

$$\bar{x}_2^* = 3.141592614 \times 10^4$$

$$\bar{x}_3^* = 3.141592643$$

を代数方程式(2)に代入する.

$$\begin{aligned} f_3(\bar{x}_2^*) &= 1.807590000 \times 10^3 \\ f_3(\bar{x}_3^*) &= -4.400000000 \times 10^{-1} \end{aligned} \quad (15)$$

最大影響係数による誤差⑧ ($x_2 = \bar{x}_2^*$, $x_3 = \bar{x}_3^*$) と関連数値 $f_3(\bar{x}_2^*)$ (15), $f_3(\bar{x}_3^*)$ (15) とから, 二つの近似解 \bar{x}_2^* , \bar{x}_3^* はそれぞれ代数方程式(2)を満足している.

7. ま と め

代数方程式の解を数値的に求めるとき, 代数方程式の係数として与えられる有限桁の数値の最下位の桁の次の桁から下位へ, 非常に多くの桁に0を加え, あたかも, 0を加えた桁まで正しい数値として与えられたように, 数値計算を行う場合がある. しかし, これは多項式の減次の数値計算における簡単な加算の演算段階で起こる桁落ちによる大きな情報損失を防ぐためのものである. 本論文では, この大きな情報損失を防ぐため, 代数方程式の係数として与えられる数値の最下位の桁の次の桁から0を加えるべき桁数を, 数値計算の途中で得られる数値を用いて決定できた.

また, 多項式の減次の数値計算の手順も, 数値計算の途中で得られる数値を用いて決定し, 代数方程式の

係数として与えられる有限桁の数値の最下位の桁の次の桁から下位へ, 多くの桁に0を加え, 数値計算に用いる桁数を特別に増加させなくても, 次数を一つ減じた代数方程式から, 与えられる代数方程式を満足する近似解を求めることができた.

参 考 文 献

- 1) Wilkinson, J. H.: Rounding Errors in Algebraic Processes, p. 55, Her Britannic Majesty's Stationery Office, London (1963).
- 2) 一松 信: 数値解析, p. 159, 税務経理協会, 東京 (1971).
- 3) 伊理正夫: 数値計算, p. 131, 朝倉書店, 東京 (1981).
- 4) 平野菅保: 代数方程式の数値解法, 日本数学会応用数学分科会講演予稿集, p. 94 (1967).
- 5) 平野菅保: 多項式の除算, 日本物理学会応用数学, 力学講演会予稿集, p. 1 (1967).
- 6) 平野菅保: 代数方程式の数値解法, 京都大学数理解析研究所講究録 72, p. 1 (1969).

(昭和58年10月25日受付)

(昭和59年1月17日採録)