

Chapter 1

Mathematical Preliminaries and

Error Analysis

Hung-Yuan Fan (范洪源)

Department of Mathematics,
National Taiwan Normal University, Taiwan

Spring 2016



Section 1.1

Review of Calculus



Def 1.1

A function $f: X \rightarrow \mathbb{R}$ has the limit L at x_0 , denoted by

$$\lim_{x \rightarrow x_0} f(x) = L,$$

if $\forall \varepsilon > 0, \exists \delta > 0$ s.t. $x \in X, 0 < |x - x_0| < \delta \Rightarrow |f(x) - L| < \varepsilon$.



Def 1.2

- ① A function $f: X \rightarrow \mathbb{R}$ is continuous (簡寫: conti.) at $x_0 \in X$ if
$$\lim_{x \rightarrow x_0} f(x) = f(x_0).$$
- ② f is conti. on X if it is conti. at each point of X .
- ③ $C(X) = \{f \mid f \text{ is conti. on } X\}$ denotes the set of all conti. functions defined on X .

Note: if $X = [a, b], (a, b), [a, b)$ or $(a, b]$ with $a < b$, write $C[a, b]$, $C(a, b)$, $C[a, b)$ or $C(a, b]$, respectively.



Def 1.3

A sequence (簡寫: seq.) of real numbers $\{x_n\}_{n=1}^{\infty}$ converges (簡寫: conv.) to the limit x , written

$$\lim_{n \rightarrow \infty} x_n = x, \text{ or } x_n \rightarrow x \text{ as } n \rightarrow \infty,$$

if $\forall \varepsilon > 0, \exists N(\varepsilon) \in \mathbb{N}$ s.t. $n > N(\varepsilon) \Rightarrow |x_n - x| < \varepsilon$.

Thm 1.4 (序列與連續性的關係)

Let f be a real-valued function defined on $\emptyset \neq X \subseteq \mathbb{R}$ and $x_0 \in X$.
The followings are equivalent:

- f is conti. at x_0 .
- \forall seq. $\{x_n\}_{n=1}^{\infty} \subseteq X$ with $\lim_{n \rightarrow \infty} x_n = x_0, \lim_{n \rightarrow \infty} f(x_n) = f(x_0)$.



Def 1.5

- ① A function $f: X \rightarrow \mathbb{R}$ is differentiable (簡寫: diffi.) at $x_0 \in X$ if

$$f'(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0} = \lim_{h \rightarrow 0} \frac{f(x_0 + h) - f(x_0)}{h}.$$

- ② f is cdiff. on X if it is cdiffy6t. at each point of X .
- ③ $C^n(X)$ denotes the set of all functions having n conti. derivatives on X .
- ④ $C^\infty(X)$ denotes the set of functions having derivatives of all orders on X .



Thm 1.6

Let f be a real-valued function defined on X and $x_0 \in X$. Then

$$f \text{ is diff. at } x_0 \implies f \text{ is conti. at } x_0.$$



Thm 1.7 (Rolle's Thm)

$f \in C[a, b]$ and f is diff. on (a, b) . If $f(a) = f(b)$, then $\exists c \in (a, b)$ s.t. $f'(c) = 0$.



Thm 1.10 (Generalized Rolle's Thm)

$f \in C[a, b]$ is n times diff. on (a, b) . If $f(x_i) = 0$ for some $n + 1$ distinct numbers $a \leq x_0 < x_1 < \dots < x_n \leq b$, then
 $\exists c \in (x_0, x_n) \subseteq [a, b]$ s.t. $f^{(n)}(c) = 0$.



Thm 1.8 (MVT)

$f \in C[a, b]$ and f is diff. on (a, b) . Then $\exists c \in (a, b)$ s.t.

$$f'(c) = \frac{f(b) - f(a)}{b - a} \quad \text{or} \quad f(b) - f(a) = f'(c)(b - a).$$



Thm 1.9 (EVT)

If $f \in C[a, b]$, then $\exists c_1, c_2 \in [a, b]$ s.t.

$$f(c_1) \leq f(x) \leq f(c_2) \quad \forall x \in [a, b].$$



Thm 1.11 (IVT)

$f \in C[a, b]$, K is any number between $f(a)$ and $f(b)$
 $\implies \exists c \in (a, b) \text{ s.t. } f(c) = K.$



Thm 1.14 (Taylor's Thm, 泰勒定理)

$f \in C^n[a, b]$, $f^{(n+1)} \exists$ on $[a, b]$ and $x_0 \in [a, b]$.

$\Rightarrow \forall x \in [a, b], \exists \xi(x)$ between x_0 and x s.t. $f(x) = P_n(x) + R_n(x)$, where

$$P_n(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k, \text{(the } n\text{th Taylor poly. for } f\text{)}$$

$$R_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x - x_0)^{n+1}.$$

(**remainder or truncation error** associated with $P_n(x)$)



Remarks

- ① If $\lim_{n \rightarrow \infty} R_n(x) = 0 \quad \forall x \in I$ (I : interval with $x_0 \in I$), then

$$f(x) = \lim_{n \rightarrow \infty} P_n(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k \quad \forall x \in I.$$

We say that the **Taylor series for f about x_0** conv. to f on I .

- ② If $x_0 = 0$, the Taylor series is often called the **Maclaurin series**.



Example 3, p. 11

The second (or third) Taylor poly. for $f(x) = \cos x$ about $x_0 = 0$ is $P_2(x) = P_3(x) = 1 - \frac{1}{2}x^2$, but their truncation errors satisfy

$$|R_2(x)| \leq \frac{|\sin \xi(x)| |x|^3}{6} \leq \frac{|x|^4}{6} = 0.1\bar{6} \cdot |x|^4$$

$$(\because |\sin \xi(x)| \leq |\xi(x)| \leq |x| \quad \forall x \in \mathbb{R})$$

$$|R_3| \leq \frac{|\cos \tilde{\xi}(x)| |x|^4}{24} \leq \frac{|x|^4}{24} = 0.041\bar{6} \cdot |x|^4.$$

(Sharper Bound for $|x| \approx 0!$)



Remark

Two objectives of numerical analysis:

- ① Find an approximation to the solution of a given problem.
- ② Determine a bound for the accuracy of the approximation.
Is this error bound tight and sharp?



Def 1.12 (定積分的定義)

- ① The (Riemann) definite integral of f on $[a, b]$ is defined by

$$\int_a^b f(x) dx = \lim_{\substack{\max \\ 1 \leq i \leq n}} \Delta x_i \rightarrow 0 \sum_{i=1}^n f(z_i) \Delta x_i,$$

where $P = \{a = x_0 < x_1 < \cdots < x_n = b\}$ is any partition of $[a, b]$, $z_i \in [x_{i-1}, x_i]$ and $\Delta x_i = x_i - x_{i-1}$ for $i = 1, 2, \dots, n$.

- ② f is called (Riemann) integrable over $[a, b]$ if the limit exists.

Note: f is conti. on $[a, b] \Rightarrow f$ is integrable over $[a, b]$.



Remark

f is integrable over $[a, b] \Rightarrow$

$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \sum_{i=1}^n f(z_i) x, \quad (x = \frac{b-a}{n})$$

$$\approx \sum_{i=0}^n w_i \cdot f(x_i) x, \quad (w_i: \text{weighting coeff.})$$

with $z_i = x_i$ or x_{i-1} for $i = 1, 2, \dots, n$.



Riemann Sums (黎曼和) with $z_i = x_i \quad \forall i$



Thm 1.13 (定積分的權重均值定理)

$f \in C[a, b]$ and g is an integrable function that **does not** change sign on $[a, b]$. Then $\exists c \in (a, b)$ s.t.

$$\int_a^b f(x)g(x) dx = f(c) \int_a^b g(x) dx.$$

Note: When $g(x) \equiv 1$, we have

$$f(c) = \frac{1}{b-a} \int_a^b f(x) dx \equiv f_{avg},$$

where f_{avg} is the **average value** of f on $[a, b]$.



The Average Value of a Function



Section 1.2

Round-off Errors and Computer Arithmetic

(捨入誤差與電腦算術)



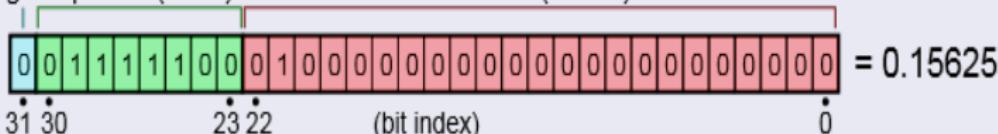
Binary Machine Numbers (二進位機器數字)

IEEE 754-1985 Standard (updated version: IEEE 754-2008)

- ### ① Single Precision Format (32 bits; 單精度)

sign exponent (8 bits)

fraction (23 bits)



- ## ② Double Precision Format (64 bits; 雙精度)

exponent

sign (11 bit)

fraction

(52 bit)



- ### ③ Extended Precision Format (80 bits; 擴充精度)

sign: 1 bit, exponent: 16 bits, fraction: 63 bits



64-bit Floating-Point Representation

- 64-bit representation is used for a real number.
- Each **binary** floating-point number (浮點數) has at least 16 **decimal digits** of precision.
- 1-bit **sign** (符號) s is followed by 11-bit **exponent** (指數) c (**characteristic**, $0 \leq c \leq 2^{11} - 1 = 2047$) and 52-bit binary **fraction** f (**mantissa**: 尾數).



The Normalized Forms (正規化形式或標準化形式)

- Normalized *binary* floating-point form of $x \in \mathbb{R}$ is

$$fl(x) = (-1)^s 2^{c-1023} (1+f)_2 = (-1)^s \left(1 + \sum_{i=1}^k b_i 2^{-i}\right)_{10} 2^{c-1023},$$

where $f = (0.b_1 b_2 \cdots b_k)_2$.

- $\mathfrak{F} = \{fl(y) \mid y \in \mathbb{R}\}$ is a **finite** (and **proper**) subset of \mathbb{R} .
- The difference between two **adjacent** (相鄰的) 64-bit floating-point numbers is $\varepsilon_M = 2^{-52} \approx 2.22 \times 10^{-16}$.

Note: the **machine precision** (or **epsilon**) is

$\varepsilon_M = 2^{-23} \approx 1.19 \times 10^{-7}$ for the single precision format.



Some Examples

- ① Since

$$\begin{aligned}27.56640625_{10} &= \mathbf{11011.10010001}_2 \\&= 1.\mathbf{101110010001}_2 \times 2^4, \text{ (Normalized Form)}\end{aligned}$$

we have $s = 0$, $c = 4 + 1023 = 1027_{10} = \mathbf{10000000011}_2$ and
mantissa $f = 0.\mathbf{101110010001}_2$. Using IEEE 754 format \Rightarrow

0 **10000000011** **101110010001**00 \cdots 0 (補 40 個零!)

- ② Note that

$$0.1_{10} = 0.0\overline{0011}_2 = 1.1\overline{0011}_2 \times 2^{-4}.$$

How to store 0.1_{10} by using IEEE 754 format?



Remarks on IEEE 754 Format

- ① The smallest positive floating-point number (with $s = 0$, $c = 1$ and $f = 0$) is

$$fl_{\min} = 2^{-1022}(1 + 0) \approx 2.2 \times 10^{-308}.$$

- ② The largest one (with $s = 0$, $c = 2046$ and $f = 1 - 2^{-52}$) is

$$fl_{\max} = 2^{1023}(2 - 2^{-52}) \approx 1.8 \times 10^{308}.$$

- ③ $|fl(x)| > fl_{\max} \Rightarrow \text{overflow (上溢位)}$ and $|fl(x)| < fl_{\min} \Rightarrow \text{underflow (下溢位)}$ and **reset** $x = 0$.
- ④ Two zeros **+0** (with $s = 0, c = 0, f = 0$) and **-0** (with $s = 1, c = 0, f = 0$) exist!



- Normalized decimal floating-point form of $y \in \mathbb{R}$ is

$$fl(y) = \pm 0.d_1 d_2 \cdots d_k \times 10^n,$$

where $1 \leq d_1 \leq 9$, $0 \leq d_i \leq 9$ ($i = 2, \dots, k$) and $n \in \mathbb{Z}$. In this case, $fl(y)$: **k-digit decimal machine number**.

- The k -digit $fl(y)$ of a **normalized** real number

$$y = \pm 0.d_1 d_2 \cdots d_k d_{k+1} \cdots \times 10^n$$

can be obtained by **terminating** the mantissa of y at k decimal digits.



Two Methods of Termination

① Chopping: (直接捨去法)

$$fl(y) = \pm 0.\mathbf{d}_1\mathbf{d}_2 \cdots \mathbf{d}_k \times 10^n,$$

i.e. simply chop off the digits $d_{k+1}d_{k+2} \dots$

② Rounding: (四捨五入法)

$$fl(y) = \begin{cases} \pm(0.d_1d_2 \cdots d_k + 10^{-k}) \times 10^n, & d_{k+1} \geq 5 \text{ (Round Up)} \\ \pm 0.\mathbf{d}_1\mathbf{d}_2 \cdots \mathbf{d}_k \times 10^n, & d_{k+1} < 5 \text{ (Round Down)} \end{cases}$$

$\equiv \pm 0.\delta_1\delta_2 \cdots \delta_k \times 10^n$ after chopping.



Example 1, p. 20

Determine the 5-digit (a) chopping and (b) rounding values of

$$\pi = 0.31415926\cdots \times 10^1.$$

Sol:

(a) $f_l(\pi) = 0.31415 \times 10^1$ by chopping.

(b) $f_l(\pi) = (0.31415 + 10^{-5}) \times 10^1 = 0.31416 \times 10^1$ by rounding.



Def 1.15

If p^* is an approximation to p , then

- ① the **absolute error** is $AE(p^*) = |p^* - p|$.
- ② the **relative error** is

$$RE(p^*) = \frac{|p^* - p|}{|p|}, \text{ provided that } p \neq 0.$$

Note: the relative error is independent of the magnitude of p , but the absolute error might vary widely!



Example 2, p. 21

Find the abs. and rel. errors when approximating p by p^* .

- (a) $p = 0.3000 \times 10^1$ and $p^* = 0.3100 \times 10^1$.
- (b) $p = 0.3000 \times 10^{-3}$ and $p^* = 0.3100 \times 10^{-3}$.
- (c) $p = 0.3000 \times 10^4$ and $p^* = 0.3100 \times 10^4$.

Sol:

- (a) $AE(p^*) = 0.1$ and $RE(p^*) = 0.3333 \times 10^{-1}$.
- (b) $AE(p^*) = 0.1 \times 10^{-4}$ and $RE(p^*) = 0.3333 \times 10^{-1}$.
- (c) $AE(p^*) = 0.1 \times 10^3$ and $RE(p^*) = 0.3333 \times 10^{-1}$.

(相對誤差都一樣，但是絕對誤差變化很大!)



Significant Digits (有效位數)

Def 1.16

p^* approximate $p \neq 0$ to t **significant digits** (or **figures**) if
 \exists largest $t \in \mathbb{N} \cup \{0\}$ satisfying

$$RE(p^*) = \frac{|p^* - p|}{|p|} \leq 5 \times 10^{-t}.$$

Note: for any normalized $y = 0.d_1d_2 \cdots \times 10^n \in \mathbb{R}$, its k -digit decimal representation satisfies

$$RE(fl(y)) \leq 10^{-k+1} = 10^{-(k-1)}$$

by using **chopping** (see the textbook), and

$$RE(fl(y)) \leq 0.5 \times 10^{-k+1} = 5 \times 10^{-k}$$

by using **rounding**. (See Ex. 24, p. 31!)



Elementary Floating-Pont Arithmetic

For floating-point representations $fl(x)$ and $fl(y)$ of real numbers x and y , assume that

$$\begin{aligned}x \oplus y &= fl(fl(x) + fl(y)), & x \otimes y &= fl(fl(x) \times fl(y)), \\x \ominus y &= fl(fl(x) - fl(y)), & x \oslash y &= fl(fl(x) \div fl(y)).\end{aligned}$$

Note: in practical computation, we usually have

$$fl(x \mathbf{op} y) = (x \mathbf{op} y)(1 + \delta) \text{ with } |\delta| \leq \varepsilon_M,$$

where $\mathbf{op} = +, -, \times, \div$, and ε_M is the machine precision.



Cancellation of Significant Digits

If $x, y \in \mathbb{R}$ ($x > y$) have the k -digit decimal representations

$$fl(x) = 0.\mathbf{d}_1\mathbf{d}_2 \cdots \mathbf{d}_p \alpha_{p+1} \alpha_{p+2} \cdots \alpha_k \times 10^n,$$

$$fl(y) = 0.\mathbf{d}_1\mathbf{d}_2 \cdots \mathbf{d}_p \beta_{p+1} \beta_{p+2} \cdots \beta_k \times 10^n,$$

then

$$\begin{aligned} fl(x) - fl(y) &= (0.\alpha_{p+1} \alpha_{p+2} \cdots \alpha_k - 0.\beta_{p+1} \beta_{p+2} \cdots \beta_k) \times 10^{n-p} \\ &\equiv 0.\sigma_{p+1} \sigma_{p+2} \cdots \sigma_k \times 10^{n-p}, \end{aligned}$$

i.e. $x \ominus y = fl(fl(x) - fl(y))$ has **at most $k - p$** significant digits,
with the last p digits being either 0 or randomly assigned.



Remark

Suppose that $f_l(z) = z + \delta$ with $|\delta|$ being the absolute error. If $\varepsilon = 10^{-n}$ with $n \in \mathbb{N}$ is a number of **small magnitude**, then

$$\frac{f_l(z)}{f_l(\varepsilon)} \approx (z + \delta) \times 10^n = \frac{z}{\varepsilon} + 10^n \delta.$$

So, the absolute error in computing z/ε is

$$\left| \frac{f_l(z)}{f_l(\varepsilon)} - \frac{z}{\varepsilon} \right| \approx 10^n \cdot |\delta| = |\delta|/\varepsilon.$$



Example 4, pp. 23–24 (1/2)

Given four real numbers

$$x = \frac{5}{7} = 0.\overline{714285}, \quad u = 0.714251$$

$$v = 98765.9, \quad w = 0.111111 \times 10^{-4}.$$

Find 5-digit chopping values of $x \ominus u$, $(x \ominus u) \oslash w$, $(x \ominus u) \otimes v$ and $u \oplus v$.

Sol: The absolute error for $x \ominus u$ is

$$\begin{aligned}|(x - u) - (x \ominus u)| &= |(x - u) - fl(fl(x) - fl(u))| \\&= \left| \left(\frac{5}{7} - 0.714251 \right) - fl(0.71428 \times 10^0 - 0.71425 \times 10^0) \right| \\&= \left| 0.347143 \times 10^{-4} - 0.30000 \times 10^{-4} \right| \\&= 0.47143 \times 10^{-5}.\end{aligned}$$



Example 4, pp. 23–24 (2/2)

The relative error for $x \ominus u$ is given by

$$RE(x \ominus u) = \left| \frac{0.47143 \times 10^{-5}}{0.347143 \times 10^{-4}} \right| = 0.1358 \leq \mathbf{0.136}.$$



Some Tricks

- ① Reformulation of the calculations to avoid the subtraction of two nearly equal numbers.
(改變計算公式以避免相近數字相減)

- ② Rearrangement of the calculations by the nested arithmetic.
(利用巢狀算術技巧以減少四則運算數量)

The lesson: **Think before you compute!**



Illustration of Trick 1

- Distinct real roots of $ax^2 + bx + c = 0$ with $a \neq 0$ and $b^2 - 4ac > 0$ are

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$$

- If $b > 0$ and $4ac \ll b^2$, then
 - $-b + \sqrt{b^2 - 4ac} \approx 0 \Rightarrow$ **Loss of accuracy for computing x_1 !**
 - Rewrite the formula for x_1 by rationalization (有理化)

$$x_1 = \frac{-2c}{b + \sqrt{b^2 - 4ac}}. \text{ (分母不是相近數相減!)}$$

- Use $x_1 x_2 = \frac{c}{a} \Rightarrow x_2 = \frac{c}{ax_1} = \frac{-b - \sqrt{b^2 - 4ac}}{2a}.$



An Example for Trick 1 (1/2)

Example, pp. 25–26

Use 4-digit rounding arithmetic to determine the first root x_1 of $f(x) = x^2 + 62.10x + 1 = 0$.

Sol: Two real roots of $f(x) = 0$ are approximately

$$x_1 = -0.01610723, \quad x_2 = -62.08390.$$

Use 4-digit rounding \Rightarrow

$$fl(\sqrt{b^2 - 4ac}) = fl(\sqrt{(62,10)^2 - (4.000)(1.000)(1.000)}) = 62.06,$$

$$f(x_1) = \frac{-62.10 + 62.06}{2.000} = -0.02000,$$

with the relative error being

$$RE(f(x_1)) = \frac{|-0.01611 + 0.02000|}{|-0.01611|} = 2.4 \times 10^{-1}.$$



An Example for Trick 1 (2/2)

In addition, if we use the reformulation for x_1 , then

$$fl(x_1) = fl\left(\frac{fl(-2c)}{fl(b + \sqrt{b^2 - 4ac})}\right) = fl\left(\frac{-2.000}{62.10 + 62.06}\right) = -0.01610,$$

which has the small relative error 6.2×10^{-4} .

Note: 近似零根 x_1 的精度提升至 3 個有效位數!



An Example of Polynomial Evaluation (1/2)

Example 6, pp. 26–27

Evaluate the 3-digit chopping and rounding values of a poly.

$$f(x) = x^3 - 6.1x^2 + 3.2x + 1.5 \text{ at } x = 4.71.$$

of y are

Sol: The actual value is $y = f(4.71) = -14.263899$. Using 3-digit chopping (rounding arithmetic), we have The 3-digits, after chopping, are

$$f(y) = f\left(((105. - 135.) + 15.1) + 1.5\right) = -13.4. \text{ (Rounding)}$$



An Example of Polynomial Evaluation (2/2)

Hence, the relative errors in computing $f_l(y)$ are

$$RE(f_l(y)) = \left| \frac{-14.263899 + 13.5}{-14.263899} \right| \approx 5.36 \times 10^{-2}, \text{ (Chopping)}$$

$$RE(f_l(y)) = \left| \frac{-14.263899 + 13.4}{-14.263899} \right| \approx 6.06 \times 10^{-2}. \text{ (Rounding)}$$

⇒ Only **one significant digit** for both chopping and rounding

values of $y = f(4.71)!$



Rearrangement of Poly. Evaluation

- Direct Computation: (4 multiplications and 3 additions)

$$f(x) = x \cdot (x \cdot x) - 6.1 \cdot (x \cdot x) + 3.2 \cdot x + 1.5$$

- Nested Computation: (2 multiplications and 3 additions)

$$f(x) = ((x - 6.1) \cdot x + 3.2) \cdot x + 1.5$$

Again, using **3-digit** arithmetic with the **nested form** \Rightarrow

$$RE(f(y)) = \left| \frac{-14.263899 + 14.2}{-14.263899} \right| \approx 4.5 \times 10^{-3}, \text{ (Chopping)}$$

$$RE(f(y)) = \left| \frac{-14.263899 + 14.3}{-14.263899} \right| \approx 2.5 \times 10^{-3}. \text{ (Rounding)}$$



Useful Suggestion

The accuracy of an approximation can be improved if we reduce the number of arithmetic operations.
(減少四則運算的數量可以改進計算解的精度!)

HW of Sec 1.2:

✓ 24,



Section 1.3

Algorithms and Convergence

(演算法與收斂性)



Algorithms and Pseudocodes (虛擬碼)

- An algorithm is a **procedure** that describes a **finite sequence of steps** to be performed in a **specified order**.
- The objective of an algorithm is to implement a procedure for **solving a problem** or **approximating a solution to the problem**.
(演算法目標是求解問題或是得到該問題的數值近似解)
- Pseudocode is an **informal environment-independent** description of the key principles of an algorithm.
- It uses structural conventions of a programming language, but is intended for **human reading** rather than machine reading.



An Example of Pseudocode

To solve the root-finding problem

$$f(x) = ax^2 + bx + c = 0 \quad \text{with} \quad a \neq 0.$$

INPUT coefficients a, b, c .

OUTPUT approximate root x .

Step 1 Compute the discriminant $D = b^2 - 4ac$.

Step 2 Compute approximate root x to $f(x) = 0$ using D .

Step 3 **OUTPUT**(x); **STOP**.



An Illustration of Algorithm

INPUT $N, x_1, x_2, \dots, x_n.$

OUTPUT $SUM = \sum_{i=1}^N x_i.$

Step 1 Set $SUM = 0.$ (累加器初始化)

Step 2 For $i = 1, 2, \dots, N$ do

set $SUM = SUM + x_i.$ (加入下一項)

Step 3 OUTPUT (SUM);

STOP.



Example 1, p. 33

The N th Taylor poly. of $f(x) = \ln x$ about $x_0 = 1$ is

$$P_N(x) = \sum_{i=1}^N \frac{(-1)^{i+1}}{i} (x-1)^i.$$

Construct an algorithm to determine the **minimal** value of N s.t.

$$|\ln(1.5) - P_N(1.5)| < 10^{-5}.$$

Note: From the Alternating Series Thm \Rightarrow

$$|\ln x - P_N(x)| \leq \left| \frac{(-1)^{N+1}}{N+1} (x-1)^{N+1} \right|.$$

So, the **stopping criterion** (停止準則) should be

$$|a_{N+1}| = \left| \frac{(-1)^{N+1}}{N+1} (x-1)^{N+1} \right| < TOL,$$

where TOL denotes the tolerance. (容許誤差)



Algorithm for Example 1

INPUT x 的值，容許誤差 TOL ，最大迭代數 M 。

OUTPUT 多項式次數 N 或錯誤訊息。

Step 1 Set $N = 1$;

$y = x - 1$;

$SUM = 0$;

$POWER = y$;

$TERM = y$;

$SIGN = -1$. (用以改變正負號)

Step 2 While $N \leq M$ do Steps 3–5.

Step 3 Set $SIGN = -SIGN$; (符號交換)

$SUM = SUM + SIGN \cdot TERM$; (累加各項)

$POWER = POWER \cdot y$;

$TERM = POWER/(N + 1)$. (計算下一項)

Step 4 If $|TERM| < TOL$ then (檢驗精度)

OUTPUT (N);

STOP. (計算成功)

Step 5 Set $N = N + 1$. (準備下一次迭代)

Step 6 OUTPUT ('Method Failed'); (計算不成功)

STOP.



Definition

- An algorithm is called **stable** if it satisfies the property that **small** changes in the initial data produce correspondingly **small** changes in the final results.
(初始資料的微小變動 \Rightarrow 計算結果也是微小變化)
- Otherwise, the algorithm is called **unstable**, i.e. **small** changes in the initial data produce **large** changes in the final results.
(初始資料的微小變動 \Rightarrow 計算結果產生大幅變化)



Def 1.17 (誤差的線性與指數成長)

$E_0 > 0$: the magnitude of error at some stage in the calculations,
 E_n : the magnitude of error after n subsequent operations.

- ① The growth of error is called **linear** if $E_n \approx CnE_0$, where the constant $C > 0$ is independent of n .
- ② The growth of error is called **exponential** if $E_n \approx C^n E_0$ for some $C > 1$.



An Unstable Procedure

The sequence $\{p_n\}_{n=0}^{\infty}$ defined by

$$p_n = c_1 \left(\frac{1}{3}\right)^n + c_2 3^n$$

is the general solution to the recursive equation (遞迴方程式)

$$p_n = \frac{10}{3}p_{n-1} - p_{n-2}, \quad n = 2, 3, \dots$$

- $p_0 = 1, p_1 = \frac{1}{3} \Rightarrow c_1 = 1, c_2 = 0$. The solution is

$$p_n = \left(\frac{1}{3}\right)^n.$$

- Use 5-digit rounding $\Rightarrow \hat{p}_0 = 1.0000, \hat{p}_1 = 0.33333$ and hence $\hat{c}_1 = 1.0000, \hat{c}_2 = -0.12500 \times 10^{-5}$. The solution is

$$\hat{p}_n = 1.0000 \left(\frac{1}{3}\right)^n - 0.12500 \times 10^{-5} (3^n).$$



Example of an Unstable Algorithm (2/2)

The absolute error in computing \hat{p}_n is

$$AE(\hat{p}_n) = p_n - \hat{p}_n = 0.12500 \times 10^{-5}(3^n).$$

⇒ An **unstable** procedure with **exponential** growth of errors!



Def 1.18

Suppose that $\{\alpha_n\}_{n=1}^{\infty}$ and $\{\beta_n\}_{n=1}^{\infty}$ are two sequences with $\lim_{n \rightarrow \infty} \alpha_n = \alpha$ and $\lim_{n \rightarrow \infty} \beta_n = 0$. If $\exists K > 0$ and $n_0 \in \mathbb{N}$ s.t.

$$|\alpha_n - \alpha| \leq K|\beta_n| \quad \forall n \geq n_0,$$

then we say that $\{\alpha_n\}_{n=1}^{\infty}$ conv. to α with **rate (or order) of convergence $O(\beta_n)$** , and write

$$\alpha_n = \alpha + O(\beta_n). \quad (\text{as } n \rightarrow \infty)$$

Note: seq. $\{\alpha_n\}_{n=1}^{\infty}$ is often generated by some **iterative method** (迭代法), and it is often compared with $\beta_n = \frac{1}{n^p}$ for $p > 0$.



Example 2, p. 37

For $n \geq 1$, consider two sequences of real numbers

$$\alpha_n = \frac{n+1}{n^2} \quad \text{and} \quad \hat{\alpha}_n = \frac{n+3}{n^3}.$$

Determine their rates of convergence.

Sol: Since

$$|\alpha_n - 0| = \frac{n+1}{n^2} \leq \frac{n+n}{n^2} = 2 \cdot \frac{1}{n} \equiv 2\beta_n,$$

$$|\hat{\alpha}_n - 0| = \frac{n+3}{n^3} \leq \frac{n+3n}{n^3} = 4 \cdot \frac{1}{n^2} \equiv 4\hat{\beta}_n$$

for all $n \geq 1$, it follows that

$$\alpha_n = 0 + \mathbf{O}\left(\frac{1}{n}\right), \quad \hat{\alpha}_n = 0 + \mathbf{O}\left(\frac{1}{n^2}\right).$$



Def 1.19

Suppose that $\lim_{h \rightarrow 0} F(h) = L$ and $\lim_{h \rightarrow 0} G(h) = 0$. If $\exists K > 0$ and $\delta > 0$ s.t.

$$|F(h) - L| \leq K|G(h)| \quad \text{for } 0 < |h| < \delta,$$

then we write

$$F(h) = L + \mathbf{O}(G(h)). \quad (\text{as } h \rightarrow 0)$$

Note: In practice, we often choose $G(h) = h^p$ for $p > 0$, and the largest value of p is expected.



Example 3, p. 38

Show that $\cos h + \frac{1}{2}h^2 = 1 + O(h^4)$.

pf: From Taylor's Thm, $\exists \xi(h)$ between 0 and h s.t.

$$\cos h = 1 - \frac{1}{2}h^2 + \frac{\cos \xi(h)}{24}h^4 \quad \text{for } h \neq 0.$$

Hence, we see that

$$\left|(\cos h + \frac{1}{2}h^2) - 1\right| = \frac{|\cos \xi(h)|}{24}|h^4| \leq \frac{1}{24}|h^4| \quad \text{for } h \neq 0,$$

which gives the desired result by Def.



Thank you for your attention!

