

§1.1 Floating Point Number System

$$f = 2^{-d} \left[\sum_{k=0}^{d-1} b_{d-k} \cdot 2^k \right]$$

$$\underline{d=3} = 2^{-3} \left[\sum_{k=0}^2 b_{3-k} \cdot 2^k \right]$$

$$= 2^{-3} \left[\underbrace{b_1}_{0/1} 2^2 + \underbrace{b_2}_{0/1} 2^1 + \underbrace{b_3}_{0/1} 2^0 \right]$$

$$\pm (1 + \boxed{f}) \cdot 2^e = \pm [1. \boxed{b_1 \ b_2 \ b_3}] \times 2^e$$

$z=0$	1.	0	0	0
$z=1$	1.	0	0	1
$z=2$	1.	0	1	0
\vdots			\vdots	
$z=7$	1.	1	1	1

8 個 7

8個浮点数

8個浮点数

8個浮点数

$e=0$ $e=1$

$e=2$

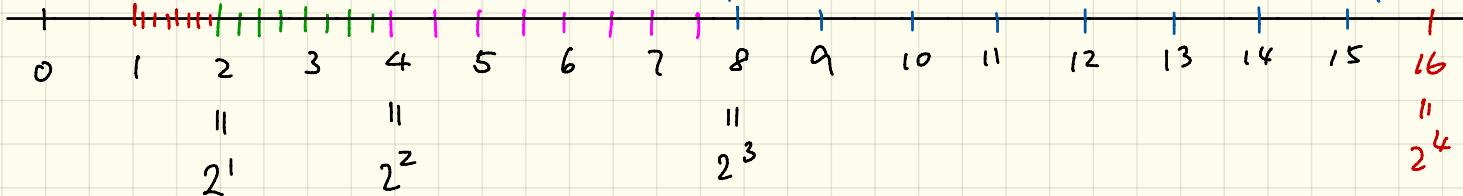
$e=3$

$\times 2^0$

$\times 2^1$

$\times 2^2$

$\times 2^3$



$$\frac{(2^2 - 2^1)}{8}$$

$$\frac{2^3 - 2^2}{8}$$

$$\frac{2^4 - 2^3}{8}$$

$\times 2$

$\times 2$

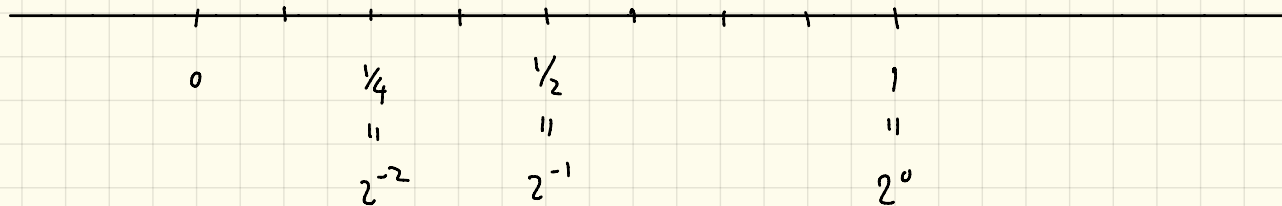
\longleftrightarrow

General case -

$$\frac{2^{e+1} - 2^e}{2^d}$$

$$= \frac{2^e (2 - 1)}{2^d}$$

$$= 2^{e-d}$$



>> ϵps

>> $1 + \epsilon ps$

>> $1 + \epsilon ps / 2$

>> $1 + \epsilon ps / 10$

>> $1 - \epsilon ps$

>> $1 - \epsilon ps / 2$

Answer

$$\gg \text{eps}$$

$$\gg (1 + \text{eps}/2) - 1$$

0

$$\gg 1 + (\text{eps}/2 - 1)$$

$1.1102e-16$

$$\gg 1 + \text{eps}$$

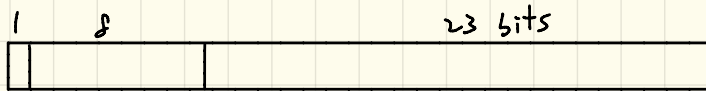
$$\gg 1 + \text{eps}/2$$

$$\gg 1 + \text{eps}/10$$

$$\gg 1 - \text{eps}$$

$$\gg 1 - \text{eps}/2$$

IEEE Single Precision

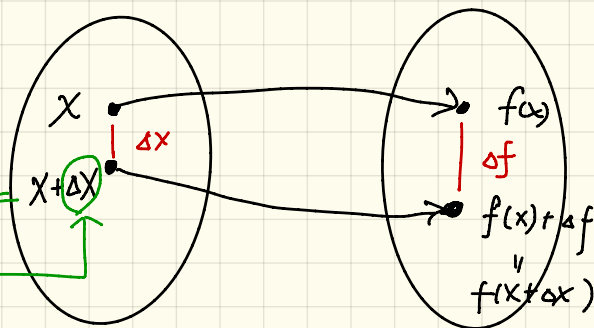


§1.2 Condition Number

In floating number

$$f(x) = x(1+\varepsilon)$$

$$= x + (\varepsilon x)$$



Kappa

$$\kappa = \left| \frac{\Delta f}{\Delta x} \right| = \left| \frac{f(x+\Delta x) - f(x)}{\Delta x} \right| = |f'(x)|$$

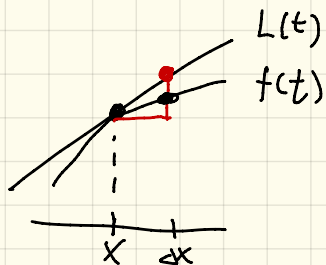
as $\Delta x \rightarrow 0$

(relative) change in output
(relative) change in input

$\kappa f(x)$ 與 (1) data(x) (2) problem (f) 有關

$$\kappa = \frac{\left| \frac{\Delta f}{f} \right|}{\left| \frac{\Delta x}{x} \right|} = \frac{\left| \frac{\Delta x f'(x)}{f(x)} \right|}{\left| \frac{\Delta x}{x} \right|} = \left| \frac{x f'(x)}{f(x)} \right|$$

as $\Delta x \rightarrow 0$



$$\frac{\Delta f}{f(x)} = \frac{f(x+\Delta x) - f(x)}{f(x)} = \frac{\Delta x f'(x)}{f(x)} + O(\Delta x^2)$$

Taylor's
as $\Delta x \rightarrow 0$

幾乎沒有 rounding error
關注源自於問題 f 和
資料 x 的影響

In floating point system, we have

$$\tilde{x} = x(1 + \varepsilon) \quad \text{where } |\varepsilon| \leq \frac{1}{2} \varepsilon_{\max}$$

$$= x + \varepsilon x$$

$$K = \lim_{\varepsilon \rightarrow 0} \frac{\left| \frac{f(x) - f(\tilde{x})}{f(x)} \right|}{\left| \frac{x - \tilde{x}}{x} \right|} = \lim_{\varepsilon \rightarrow 0} \frac{\left| \frac{f(x) - f(x(1 + \varepsilon))}{f(x)} \right|}{\left| \frac{\varepsilon x}{x} \right|} = \lim_{\varepsilon \rightarrow 0} \frac{|f(x) - f(x(1 + \varepsilon))|}{|\varepsilon f(x)|}$$

$$= \lim_{\varepsilon \rightarrow 0} \underbrace{\left| \frac{f(x) - f(\varepsilon x)}{\varepsilon x} \right|}_{f'(x)} \cdot \left| \frac{x}{f(x)} \right| = \left| \frac{x f'(x)}{f(x)} \right|$$

for small ε (not $\rightarrow 0$)

$$K \approx \left| \frac{f(x) - f(x + \varepsilon x)}{\varepsilon f(x)} \right|$$

When the data x is perturbed by a small amount, we expect the relative change to be magnified by a factor of K .

Large condition number suggests when errors cannot be expected to remain comparable in the size of the round off error.

$$\Rightarrow \left| \frac{f(x + \varepsilon x) - f(x)}{f(x)} \right| \approx K \cdot |\varepsilon|$$

Example 1:

$$f(x) = x - c$$

兩個相近數相減, $K \gg 1$

$$\Rightarrow K_f(x) = \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{x}{x-c} \right| \left[\frac{x}{\textcircled{(1)}} \rightarrow x \approx c \right] \uparrow K$$

$$\pi - \frac{355}{113}$$

$$x - c$$

$$= \begin{array}{r} \overset{9}{3.141592653589793} \\ - \overset{9}{3.141592920353983} \\ \hline -0.0000002667641894049666 \end{array} = -2.667 \dots \times 10^{-7}$$

$$K = \left| \frac{3.141592653589793}{-2.667 \dots \times 10^{-7}} \right| \approx 10^7$$

$$\log_{10} K \approx 7$$

有意義

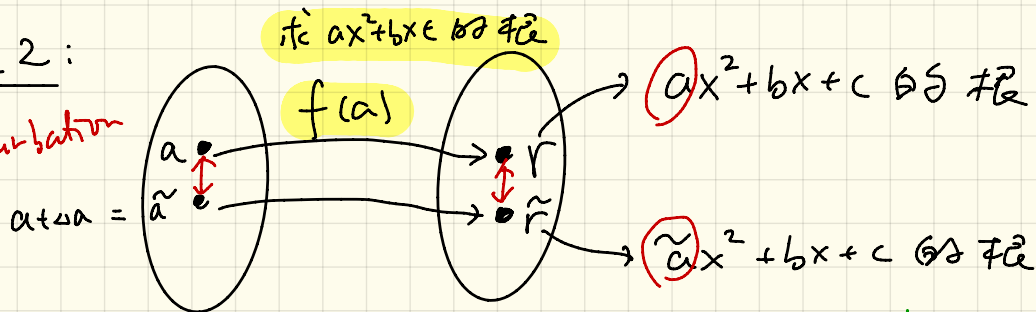
亂數!

Subtractive cancellation: 7位有效位數消失!

$$\text{Try } \pi - \frac{103683}{32989}$$

Example 2:

small perturbation
in a



$$K = \left| \frac{x f'(x)}{f(x)} \right| = \left| \frac{a f'(a)}{f(a)} \right| = \left| a \frac{df(a)}{da} \right| = \left| \frac{a r^2}{p'(r)} \right| = \left| \frac{a r}{p'(r)} \right|$$

Find $r = f(a)$ s.t. $p(r) = ar^2 + br + c = 0$

implicit differentiation

$$\Rightarrow \frac{d}{da} [\quad] = \frac{d}{da}(0)$$

$$\Rightarrow r^2 + a \cdot 2r \cdot \frac{dr}{da} + b \frac{dr}{da} = 0$$

$$\Rightarrow \frac{dr}{da} = \frac{-r^2}{2ar + b} = \frac{-r^2}{p'(r)} = \frac{dp(r)}{dr}$$

$$K = \left| \frac{ar}{p'(r)} \right| = \left| \frac{ar}{2ar+b} \right| = \left| \frac{ar}{\pm \sqrt{b^2-4ac}} \right| = \left| \frac{r}{\boxed{r_1-r_2}} \right|$$

$$r = \frac{-b \pm \sqrt{b^2-4ac}}{2a}$$

$$\Rightarrow 2ar+b = \pm \sqrt{b^2-4ac}$$

$$r_1 - r_2 = \frac{-b + \sqrt{b^2-4ac}}{2a} - \frac{-b - \sqrt{b^2-4ac}}{2a}$$

$$= \frac{\sqrt{b^2-4ac}}{a}$$

$$\Rightarrow \sqrt{b^2-4ac} = a(r_1 - r_2)$$

Note: $p'(r) = 0 \Leftrightarrow r$ is a double root. $K \rightarrow \infty$
 $= 2ar+b = \pm \sqrt{b^2-4ac}$

$p'(r) \approx 0 \Leftrightarrow$ r is nearly a double root

$K \gg 1$ $1.01x^2 - 6x + 9$ 條件數大
誤差大!

§ 1.2 Stability

When error in the result of an algorithm exceeds what conditioning can explain, the algorithm is unstable.

The sensitivity of an algorithm depends on the condition numbers of all of its steps.

Example 1.3.3

$$p(x) = (x - 10^6)(x - 10^{-6}) = 1x^2 - (10^6 + 10^{-6})x + 1$$

$\underbrace{1}_{a}$
 $\underbrace{-(10^6 + 10^{-6})}_{b}$
 $\underbrace{+1}_{c}$

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

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

computed

$$x_1 = 1000000$$

$$x_2 = 1.000007614493370 \times 10^{-6}$$

rel. err.

ϕ

10^{-5}

Calculation	Result	κ
$u_1 = b^2$	$1.00000000002000 \times 10^{12}$	2
$u_2 = u_1 - 4$	$9.99999999980000 \times 10^{11}$	$ u_1 / u_2 \approx 1.00$
$u_3 = \sqrt{u_2}$	999999.9999990000	1/2
$u_4 = u_3 - b$	2000000	$ u_3 / u_4 \approx 0.500$
$u_5 = u_4/2$	1000000	1

$$u_4 = -u_3 - b$$

$$= -999999.9999990000 + 1000000.000001$$

兩相近數相減 $\approx 10^{-11}$
 subtractive cancelling

Big cond. #

$$\kappa = \frac{|u_3|}{|u_4|} \approx \frac{0.5}{10^{-11}} \approx 10^{11}$$

$$\kappa \cdot \epsilon \approx 10^{11} \times 10^{-16} \approx 10^{-5}$$

Function	Condition number
$f(x) = x + c$	$\kappa_f(x) = \frac{ x }{ x+c }$
$f(x) = cx$	$\kappa_f(x) = 1$
$f(x) = x^p$	$\kappa_f(x) = p $
$f(x) = e^x$	$\kappa_f(x) = x $
$f(x) = \sin(x)$	$\kappa_f(x) = \cot(x) $
$f(x) = \cos(x)$	$\kappa_f(x) = x \tan(x) $
$f(x) = \log(x)$	$\kappa_f(x) = 1/ \log(x) $

Solution:

2 相近數相減，會有問題

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a} = \frac{(-b - \sqrt{\quad})(-b + \sqrt{\quad})}{2a(-b + \sqrt{\quad})}$$

$$= \frac{\cancel{b^2} - (\cancel{b^2} - \cancel{4ac})}{\cancel{2a}(-b + \sqrt{\quad})} = \frac{2c}{-b + \sqrt{\quad}} = \frac{\cancel{2c}}{\cancel{2a}x_1}$$

2 相近數相加，
沒問題

$$x_1 = \frac{-b + \sqrt{\quad}}{2a}$$

$$= 1.0000000000000000 \times 10^{-6}$$

由於計算造成的影響

How close to the true answer is your answer?
 $f(x)$ $\tilde{f}(x)$

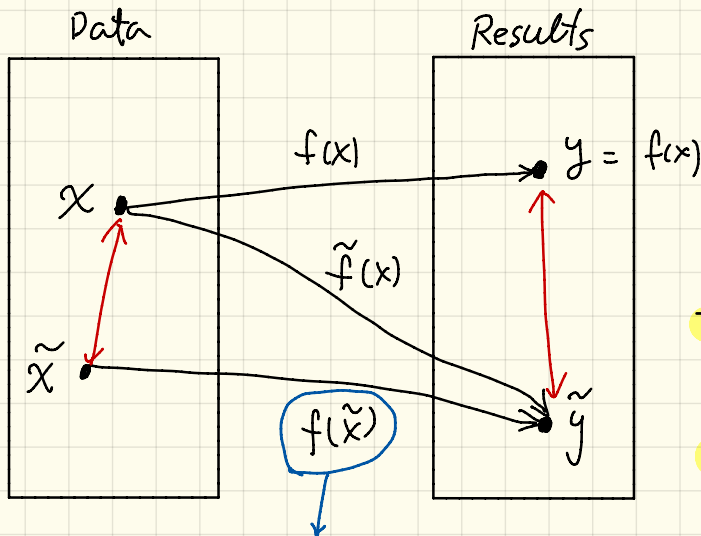
如果 $\exists \tilde{x}$, s.t.

$$f(\tilde{x}) = \tilde{f}(x)$$

Backward error

$$\frac{|\tilde{x} - x|}{|x|}$$

small B. error \Rightarrow



x 沒有 round-off error
 是 "計算" \tilde{f} 造成
 的結果的誤差

$$\frac{|\tilde{f}(x) - f(x)|}{|f(x)|}$$

Forward error

the algorithm ($\tilde{f}(x)$) gives the correct answer
 to nearly the right problem \tilde{x}

How close to the true question is the question you answered?
 x \tilde{x}

$$r = [-2 \ -1 \ 1 \ 1 \ 3 \ 6]' \Rightarrow \text{true roots}$$

$$p = \text{poly}(r) = [1 \ -8 \ 6 \ 44 \ -43 \ -36 \ 36]$$

Backward Error

$$(p_{\text{computed}} - p) / p$$

$$\begin{aligned} &-1 \times 10^{-15} \\ &-4 \times 10^{-15} \\ &-4 \times 10^{-16} \\ &-6 \times 10^{-16} \\ &-2 \times 10^{-15} \\ &-1 \times 10^{-15} \end{aligned}$$

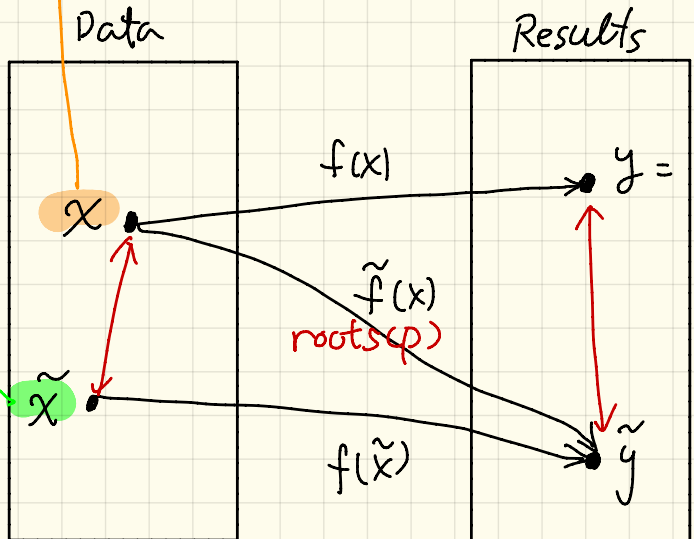
$$\text{Forward error} = \frac{|r - r_{\text{computed}}|}{|r|}$$

$$\begin{aligned} &-6 \times 10^{-16} \\ &-9 \times 10^{-16} \\ &9 \times 10^{-9} \\ &2 \times 10^{-9} \\ &1 \times 10^{-16} \\ &1 \times 10^{-15} \end{aligned}$$

poly(r-computed)

p-computed

$$\begin{aligned} &1. \dots \\ &-8. \dots \\ &6. \dots \\ &44. \dots \\ &-43. \dots \\ &-36. \dots \\ &36. \dots \end{aligned}$$



$$r = \begin{bmatrix} -2 \\ -1 \\ 1 \\ 1 \\ 3 \\ 6 \end{bmatrix}$$

nearly double roots, ill-conditioned

r-computed

$$\begin{aligned} &-2. \dots \\ &-1. \dots \\ &1. \dots \\ &1. \dots \\ &3. \dots \\ &6. \dots \end{aligned}$$

If an algorithm always produces small backward errors, then it is stable.

