

Machine epsilon

Rounding :

$$1.0 = 1.\underbrace{0\dots 0}_{t-1 \text{ 0's}} \times b^0$$

Next number after 1.0 is

$$1.\underbrace{0\dots 0}_{t-2 \text{ 0's}} 1 \times b^0$$

$$1.0 = 1.\overbrace{0\dots 0}^{t-2} \times b^0$$

$$\epsilon_{\text{mach}} = 0.\overbrace{0\dots 0}^{t-2} 005 \times b^0$$

$$1.0 + \epsilon_{\text{mach}} = 1.\overbrace{0\dots 0}^{t-2} 005 \times b^0$$

$$1.0 + \epsilon_{\text{mach}} = 1.\overbrace{0\dots 0}^{t-2} 01 \times b^0$$

base 10 (1's if base 2)

round to the next number after 1.0

$$\text{So } \epsilon_{\text{mach}} = 5 \times b^{1-t} = \frac{1}{2} b^{1-t}$$

— rounding

FORTRAN has built-in (intrinsic) functions

$\epsilonpsilon(1.0) \approx 10^{-7}$ for REAL

$\epsilonpsilon(10.0) \approx 10^{-15}$ for DOUBLE PRECISION

Roundoff Error in Multiplication

①

$$x_1 = \bar{x}_1 + e_1, \quad x_2 = \bar{x}_2 + e_2$$

\bar{x}_1, \bar{x}_2 are the numbers on the computer which represent the actual numbers x_1, x_2 ; e_1, e_2 are the roundoff errors.

Multiply \bar{x}_1 and \bar{x}_2 on the machine,

$$\bar{x}_1 \bar{x}_2 = \overline{\bar{x}_1 \bar{x}_2} + e_{12}^* \quad \left\{ \begin{array}{l} \text{error in } \bar{x}_1 \bar{x}_2 \end{array} \right.$$

$$\left| \text{Relative error in } \bar{x}_1 \bar{x}_2 \right| \leq \epsilon_{\text{mach}} \quad \left(\begin{array}{l} \text{for the numbers} \\ \text{we are} \\ \text{computing with} \end{array} \right)$$

$$\therefore \frac{|e_{12}^*|}{|\bar{x}_1 \bar{x}_2|} \leq \epsilon_{\text{mach}}, \quad \text{i.e. } |e_{12}^*| \leq \epsilon_{\text{mach}} |\bar{x}_1 \bar{x}_2|$$

In practice this is a reasonable estimate not simply an upper bound, i.e.

$$|e_{12}^*| \sim \epsilon_{\text{mach}} |\bar{x}_1| |\bar{x}_2|$$

↑ means "roughly equals".

What is the error in $x_1 x_2$? ②

$$\begin{aligned}
 x_1 x_2 &= (\bar{x}_1 + e_1)(\bar{x}_2 + e_2) \\
 &= \bar{x}_1 \bar{x}_2 + e_1 \bar{x}_2 + e_2 \bar{x}_1 + e_1 e_2 \\
 &= \bar{x}_1 \bar{x}_2 + e_1 \bar{x}_2 + e_2 \bar{x}_1 + \cancel{e_1 e_2} \quad \text{neglect since very small}
 \end{aligned}$$

$$\sim \bar{x}_1 \bar{x}_2 \pm |\bar{x}_1| |\bar{x}_2| \epsilon_{mach} + e_1 \bar{x}_2 + e_2 \bar{x}_1$$

$$\left| \frac{x_1 x_2 - \bar{x}_1 \bar{x}_2}{x_1 x_2} \right| \sim \left| \frac{\bar{x}_1}{x_1} \right| \left| \frac{\bar{x}_2}{x_2} \right| \epsilon_{mach} + \left| \frac{e_1}{x_2} \right| \left| \frac{\bar{x}_2}{x_2} \right| + \left| \frac{e_2}{x_1} \right| \left| \frac{\bar{x}_1}{x_1} \right|$$

$$\sim \epsilon_{mach} + \left| \frac{e_1}{x_1} \right| + \left| \frac{e_2}{x_2} \right|$$

$$\rightarrow \epsilon_{mach} + RE_1 + RE_2$$

\downarrow
 relative error in x_1 relative error in x_2

when entered in computer.

$$\sim 3 \epsilon_{mach}$$

Similarly for division.

Relative Function Test

①

Consider a function $f(x)$ of the form

$$f(x) = \sum_{k=1}^K t_k(x)$$

i.e. f is a sum of K terms.

We seek a zero of f , i.e. x^* such

that $f(x^*) = 0$. This implies

$$\sum_{k=1}^K t_k(x^*) = 0$$

In general we do not expect each $t_k(x^*) = 0$ for $k = 1:K$.

Thus to get 0 there must be cancellation of the $t_k(x^*)$.

In practice if \bar{x} is the approximation we compute to x^* , then

we get $f(\bar{x}) = \sum_{k=1}^K t_k(\bar{x})$. How close to 0 is this?

≈ 0

The relative function test is the ②
following: if $\{x_i\}$ is an approximation
to a zero x^* of $f(x)$, then stop

$$\text{if } \left| \frac{f(x_i)}{t_k(x_i)} \right| < \epsilon$$

where ϵ is a given error tolerance

$$\& \quad |t_k(x^*)| \geq |t_{k-1}(x^*)| \geq \dots \geq |t_1(x^*)|.$$

Another form for the test is

$$\frac{|f(x_i)|}{\sum_{k=1}^k |t_k(x_i)|} = \frac{\left| \sum_{k=1}^k t_k(x_i) \right|}{\sum_{k=1}^k |t_k(x_i)|} < \epsilon.$$

Eg1

$$\begin{aligned}
 t_1(\bar{x}) &= 0. \overbrace{x_1 x_1 x_1 x_1 x_1 x_1 x_1}^7 \text{---} \text{---} \text{---} \\
 t_2(\bar{x}) &= 0. x_2 x_2 \dots x_2 \text{---} \text{---} \text{---} \\
 &\vdots \\
 t_k(\bar{x}) &= 0. x_k x_k \dots x_k \text{---} \text{---} \text{---}
 \end{aligned}$$

↓ unreliable
 after 7th decimal place terms are about 1

$$0 \approx 0.0000000 \text{---} \text{---} \text{---}$$

↓ anything

best outcome possible

∴ the error in the best case will be about 10^{-7} 7 digits before decimal point

Eg2

$$\begin{aligned}
 t_1(\bar{x}) &= \overbrace{x_1 x_1 \dots x_1}^7 \cdot \text{---} \text{---} \text{---} \\
 &\vdots \\
 t_k(\bar{x}) &= x_k x_k \dots x_k \cdot \text{---} \text{---} \text{---}
 \end{aligned}$$

} terms are about 10^7

$$0 \approx 1. \text{---} \text{---} \text{---}$$

Error Analysis of $\delta_h f(x)$

①

$$\delta_h f(x) := \frac{f(x+h) - f(x-h)}{2h}$$

cf $\frac{f(x+h) - f(x)}{h} = \Delta_h f(x)$

We can use $\delta_h f(x)$ or $\Delta_h f(x)$ to approximate $\frac{df(x)}{dx} \equiv f'(x)$.

$\delta_h f(x)$ is more accurate.

b) Truncation error :

$$f(x+h) = f(x) + f'(x)h + \frac{f^{(2)}(x)h^2}{2} + \frac{f^{(3)}(x)h^3}{6} + \frac{f^{(4)}(x)h^4}{24} + \frac{f^{(5)}(\eta_+)h^5}{120}$$

$$f(x-h) = f(x) - f'(x)h + \frac{f^{(2)}(x)h^2}{2} - \frac{f^{(3)}(x)h^3}{6} + \frac{f^{(4)}(x)h^4}{24} - \frac{f^{(5)}(\eta_-)h^5}{120}$$

(2)

$$\delta_h f(x) = f'(x) + f^{(3)}(x) \frac{h^2}{6} + \frac{1}{2} [f^{(5)}(\eta_+) + f^{(5)}(\eta_-)] \frac{h^4}{120}$$

But if $f^{(5)}$ is continuous there exists η such that

$$\frac{1}{2} [f^{(5)}(\eta_+) + f^{(5)}(\eta_-)] = f^{(5)}(\eta)$$

using the intermediate value theorem.

$$\delta_h f(x) = f'(x) + \underbrace{f^{(3)}(x) \frac{h^2}{6} + f^{(5)}(\eta) \frac{h^4}{120}}_{\text{truncation error in } \delta_h f(x)}$$

(c) Roundoff error in formula

Assume $\delta_h f(x)$ is calculated exactly (eg in computer registers), then rounded or chopped. Then

$$\begin{aligned} |\text{roundoff error in } \delta_h f(x)| &\leq \epsilon_{\text{mach}} \times |\delta_h f(x)| \\ &\leq \frac{|f(x+h)| + |f(x-h)|}{2h} \epsilon_{\text{mach}} \end{aligned}$$

$$\approx \frac{|f(x)|}{h} \epsilon_{\text{mach}} \quad (3)$$

This is a bound on the roundoff error but it is a reasonable estimate usually, i.e.

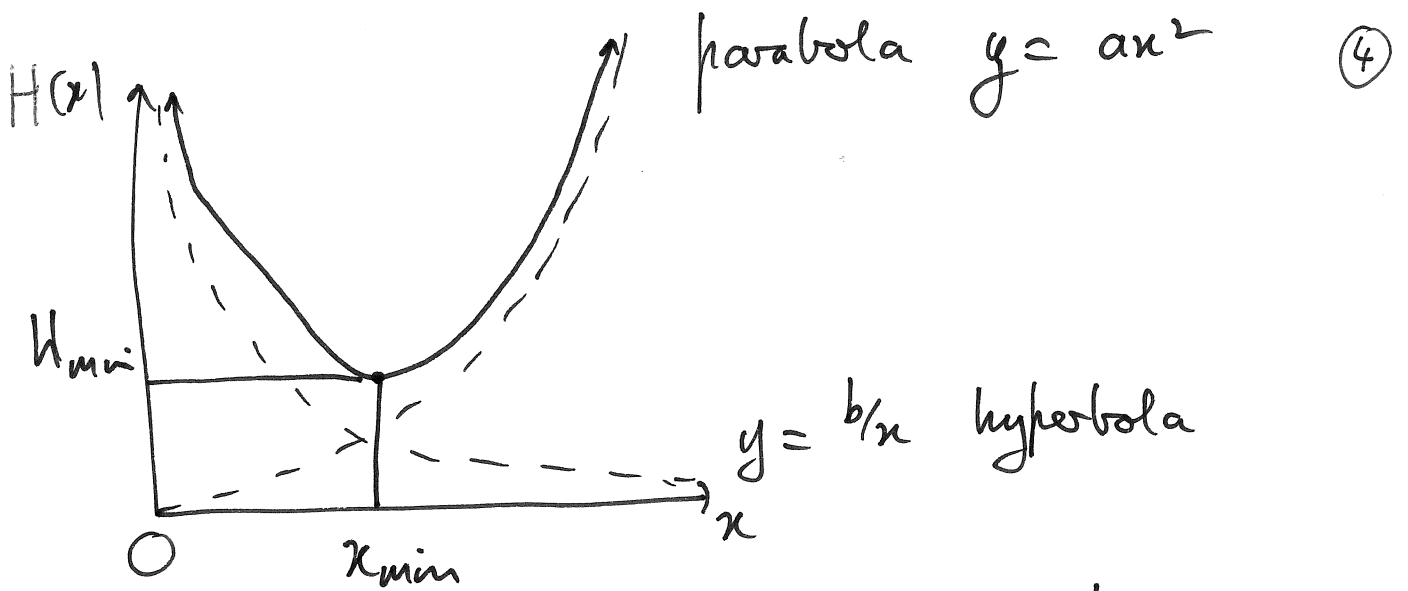
$$\left| \text{roundoff error in } \delta_h f(x) \right| \approx \frac{|f(x)|}{h} \epsilon_{\text{mach}}$$

(d) Total error

$$\left| \text{total error using } \delta_h f(x) \text{ to approximate } f'(x) \right| \leq \left| f^{(3)}(x) h^2 / 6 \right| + \text{truncation}$$

$$\frac{|f(x)|}{h} \epsilon_{\text{mach}} \quad \text{roundoff} \quad + \quad \left| f^{(5)}(x) h^4 / 120 \right| \quad \text{remainder} \quad (*)$$

Consider $H(x) = ax^2 + \frac{b}{x}$, $a, b > 0$



Minimise $H(x)$: $H'(x) = 2ax - \frac{b}{x^2} = 0$

$$\Rightarrow x_{\min} = \sqrt[3]{\frac{b}{2a}}$$

$$H_{\min} = a \left(\frac{b}{2a}\right)^{2/3} + b \left(\frac{b}{2a}\right)^{-1/3}$$

$$= \frac{3}{2} \sqrt[3]{2ab^2}$$

Comparing with the total error bound we see

$$h = x, \quad a = \frac{|f^{(3)}(x)|}{6}, \quad b = |f(x)| \epsilon_{\text{mach}}$$

The minimum bound on the total error occurs when

$$h_{\min} = \sqrt[3]{\frac{|f(x)| \epsilon_{\text{mach}}}{2 |f^{(3)}(x)| / 6}} \quad (= x_{\min})$$

$$\sim \sqrt[3]{\epsilon_{\text{mach}}} h \quad (\text{ignore } \sqrt[3]{3!})$$

where h is a measure of the length scale of $f(x)$:

$$h^3 = \left| \frac{f(x)}{f^{(3)}(x)} \right|.$$

Minimum bound (& estimate) of

the error is

$$\left| \begin{array}{l} \text{total error in} \\ \delta_h f(x) \end{array} \right| \sim \frac{f}{h} (\epsilon_{mach})^{2/3}$$

where f is a typical value of $f(x)$.

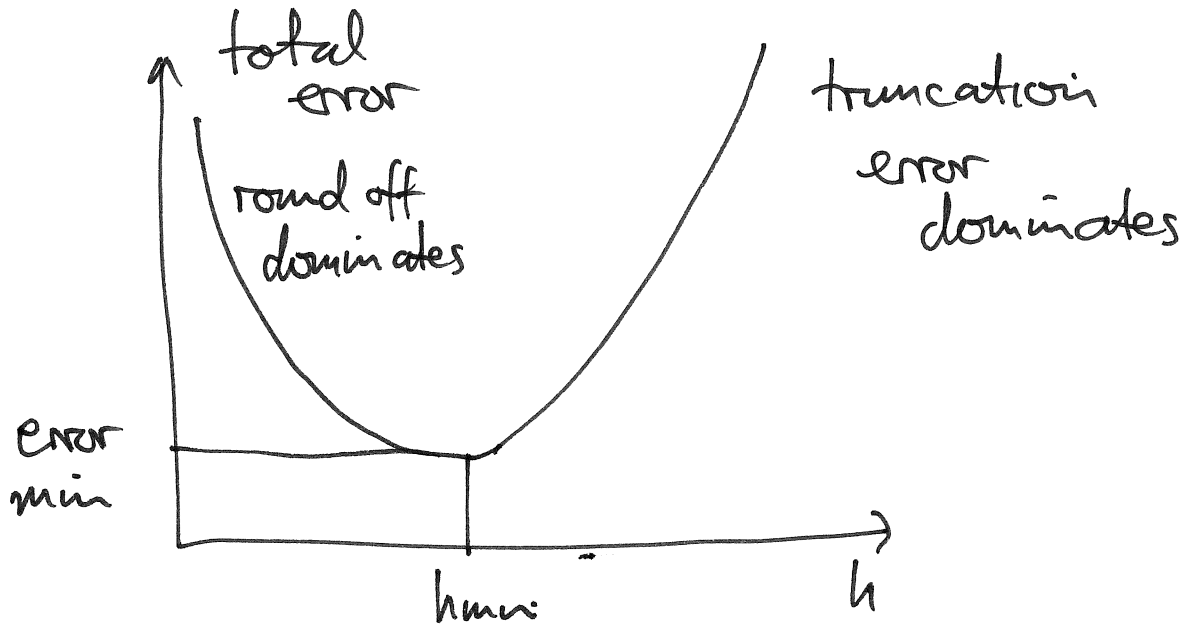
But $f/h \sim f'(x)$ so the

$$\left| \begin{array}{l} \text{minimum relative error in} \\ \delta_h f(x) \end{array} \right| \sim (\epsilon_{mach})^{2/3}.$$

In single precision with $\epsilon_{mach} \approx 10^{-7}$

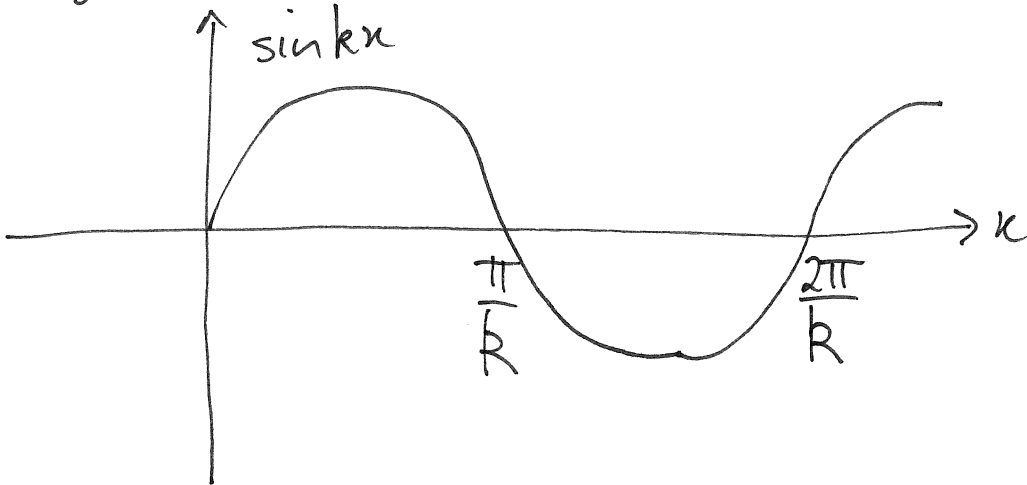
$$(\epsilon_{mach})^{2/3} \approx 10^{-4} \quad \& \quad (\epsilon_{mach})^{1/3} \approx 10^{-2}.$$

General picture

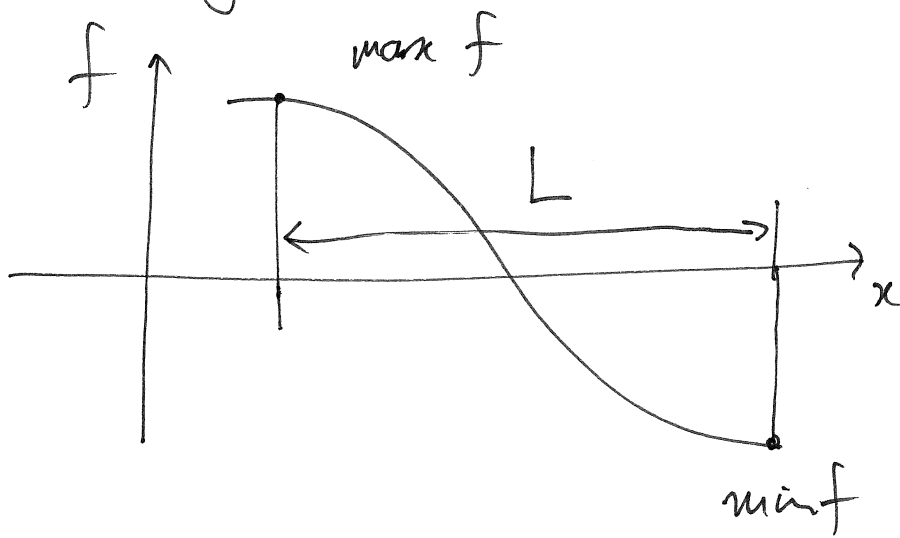


Length scale of a function

Eg $\sin kx$

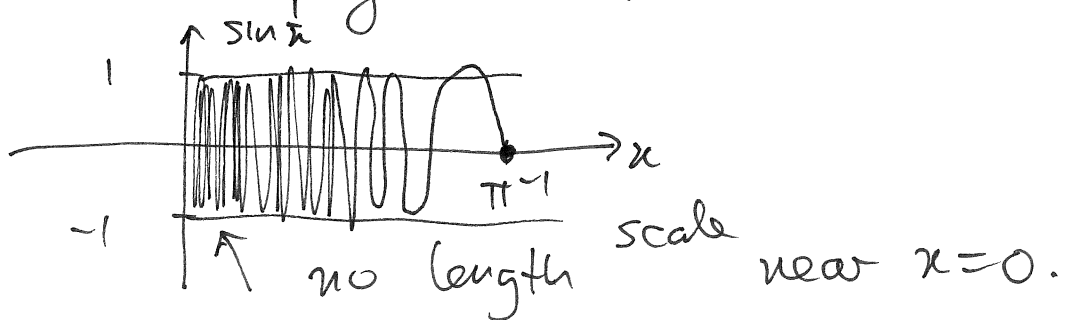


The length scale of $\sin kx$ is $\sim \frac{\pi}{k}$

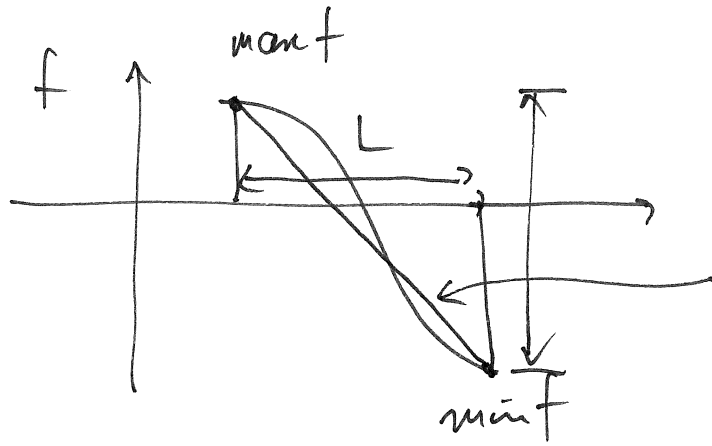


$L =$ length scale of f .

Length scale cannot be defined for every function: eg $\sin \frac{1}{x}$ in $[0, \frac{1}{\pi}]$



We can use L to estimate derivatives ⁽²⁾ of f .



$$\text{slope of line segment} \sim \frac{\text{max } f - \text{min } f}{L} \sim \frac{f}{L}$$

where f is a typical value of the function f .

This gives

$$f^{(1)}(x) = \text{slope of tangent to } f \text{ at } x \sim \frac{f}{L}$$

For higher derivatives we take

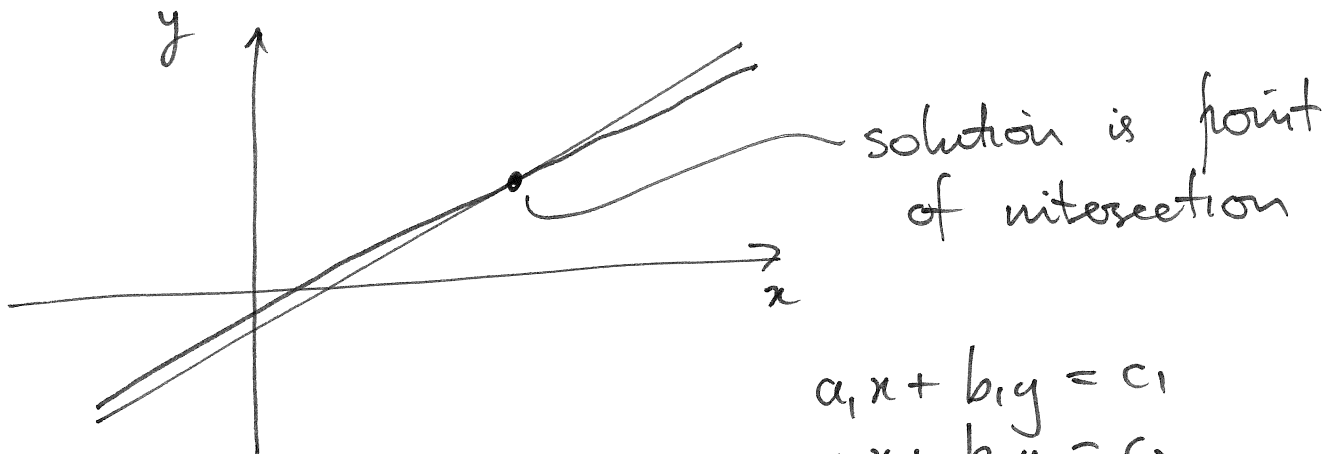
$$f^{(n)}(x) \sim \frac{f}{L^n} \quad \leftarrow \text{order of derivative}$$

Eg $f^{(3)}(x) \sim \frac{f(x)}{L^3}$, assuming $f(x)$ is typical

ie $\frac{f^{(3)}(x)}{f(x)} \sim \frac{1}{L^3}$ — we used this in the analysis of $\Delta_h f(x)$.

III- Conditioned Problems

- ① Two nearly parallel straight lines:



$$a_1x + b_1y = c_1$$

$$a_2x + b_2y = c_2$$

Here data is : $a_1, b_1, c_1, a_2, b_2, c_2$

Here a slight perturbation in the data shifts the point of intersection a long way.

- ② $y'' - y = 0, \quad y(0) = 1, \quad y'(0) = -1$

$$y = c_1 e^{-x} + c_2 e^x$$

$$y(0) = 1 \Rightarrow c_1 + c_2 = 1$$

$$y'(0) = -1 \Rightarrow -c_1 + c_2 = -1$$

$$c_1 = 1, \quad c_2 = 0$$

$$\left. \begin{array}{l} \text{Try } y = e^{px} \\ p^2 - 1 = 0 \\ \Rightarrow p = \pm 1 \\ \Rightarrow e^x, e^{-x} \end{array} \right\}$$

$$\therefore y = e^{-x} \rightarrow 0 \text{ as } x \rightarrow \infty$$

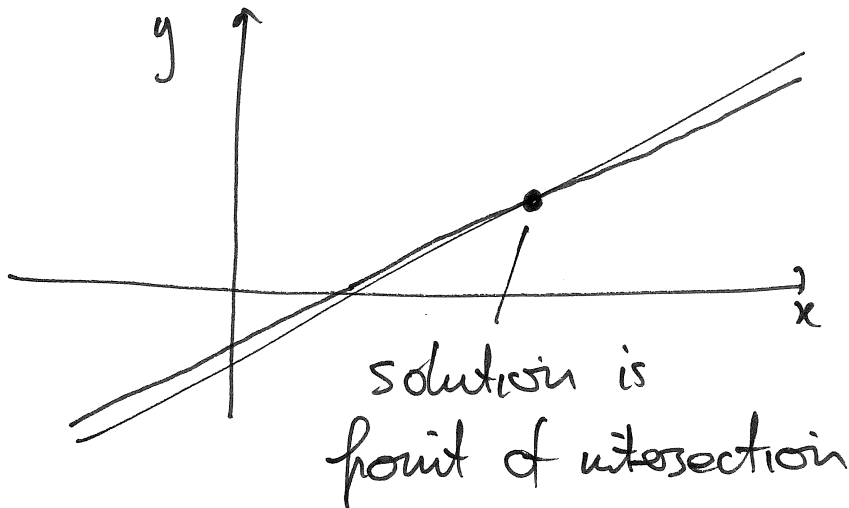
If the data is perturbed so that

$$y(0) = 1 + \delta_1, \quad y'(0) = -1 + \delta_2 \Rightarrow c_1 \approx 1, \quad c_2 \neq 0$$

Then $y \approx e^{-x} + c_2 e^x \rightarrow +\infty$ if $c_2 > 0$ but small if δ_1, δ_2 small

III - Conditioned Problems

- ① Two nearly parallel straight lines:



$$a_1x + b_1y = c_1$$
$$a_2x + b_2y = c_2$$

Here the data is
 $\left. \begin{array}{l} a_1, b_1, c_1 \\ a_2, b_2, c_2 \end{array} \right\}$.

Here a slight perturbation to the data can shift the point of intersection a long way.

② $y'' - y = 0, \quad y(0) = 1, \quad y'(0) = -1.$

$$\frac{d^2y}{dx^2}$$

Problem: determine $y(x)$.

$$y = c_1 e^x + c_2 e^{-x}$$

$$\left. \begin{array}{l} y(0) = c_1 + c_2 = 1 \\ y'(0) = c_1 - c_2 = -1 \end{array} \right\} \Rightarrow c_1 = 0, c_2 = 1$$
$$\Rightarrow y = e^{-x}$$

As $x \rightarrow \infty, y \rightarrow 0$

If the data are perturbed, eg the initial conditions at $x=0$,

$$y(0) = 1 + \delta_1 \quad y'(0) = -1 + \delta_2$$

where δ_1, δ_2 are small,

then $C_1 \neq 0, C_2 \approx 1$

The solution $y = \underset{\neq 0}{C_1} e^{+x} + C_2 e^{-x}$

For x small, $y \approx e^{-x}$.

But as $x \rightarrow \infty$, $y \rightarrow \begin{cases} +\infty & C_1 > 0 \\ -\infty & C_1 < 0 \end{cases}$.

③ Set 1 q 8 - 3 term recurrence relation

$$\phi_{k+1}(x) = \frac{2k}{x} \phi_k(x) - \phi_{k-1}(x)$$

$$\phi_k(x) = C_1 J_k(x) + C_2 Y_k(x)$$

$$\left. \begin{array}{l} \phi_0(x) = J_0(x) \\ \phi_1(x) = J_1(x) \end{array} \right\} \Rightarrow \begin{array}{l} C_1 = 1 \\ C_2 = 0 \end{array} \Rightarrow \phi_k(x) = J_k(x)$$