

# Numerical Solution

Goal: to find  $\underline{A} \underline{x} = \underline{b}$



to be found

If  $\det(\underline{A}) \neq 0$ , you could do

$$\underline{x} = \underline{A}^{-1} \underline{b}$$

You never actually find  $\underline{A}^{-1}$

Two Methods  
=

① Iterative Method.

$\Rightarrow$  Matrix-Vector Projection

$\Rightarrow$  a sequence of vectors

$\Rightarrow$  converge to  $\underline{x}$

GMRES, QMR etc.

$\rightarrow$  Krylov Subspace Method.

## ② Decomposition Method.

Given  $\underline{A}$  from  $\underline{B}$  and  $\underline{C}$

such that  $\underline{A} = \underline{C} \underline{B}$

with finding  $\underline{y}$  so that  $\underline{C} \underline{y} = \underline{b}$

and  $\underline{x}$  such that  $\underline{B} \underline{x} = \underline{y}$  is

easier than  $\underline{A} \underline{x} = \underline{b}$

$$\underline{y} = \underline{C}^{-1} \underline{b}$$

$$\underline{B} \underline{x} = \underline{y} = \underline{C}^{-1} \underline{b}$$

$$\Rightarrow \underline{x} = \underline{B}^{-1} \underline{C}^{-1} \underline{b}$$

$$\underline{x} = \underline{A}^{-1} \underline{b}$$

$\underline{A} = \underline{C} \underline{B}$
---

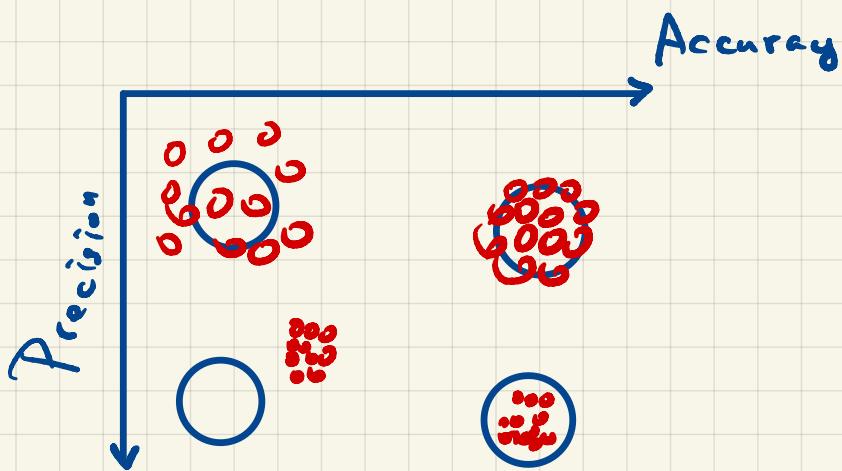
LU  
QR  
SVD

exact  
method

# Accuracy v.s. Precision

Accuracy: How close are the approximations to the truth

Precision : How close are the approximations to each other



Numerical Errors come from

- choice of the model
- numerical approximation
- limits of data representation
- Errors in implementation
- fluctuation in hardware

Truth = Approximation + Error

$$\underline{x}^* = \underline{x} + \underline{\epsilon}$$

$$\Rightarrow \underline{\epsilon} = \underline{x}^* - \underline{x} \leftarrow \text{usually not a good error}$$

Relative error

$$\epsilon_{\text{rel}} = \frac{\underline{x}^* - \underline{x}}{\underline{x}^*}$$

Note: In The iterative methods, you have

a sequence of  $\underline{x}_i$ :  $\underline{x}_0, \underline{x}_1, \dots \rightarrow \underline{x}$

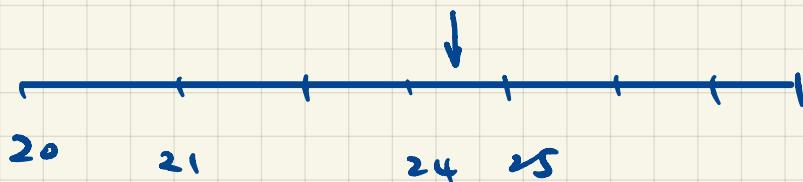
$$\underline{r}_i = \underline{A} \underline{x}_i - \underline{b}$$

↑

residual

not the error

# Significant Figures



$$24 < x < 25$$

$$x = 24.75 ?$$

$$24.525 ?$$

$$24.5 ?$$

Sig. figure tell you how much

useful information you have

Typically the known digit +  
one more

ex).

53800  $\leftarrow$  5 sig. digit.

$5.38 \times 10^4$   $\leftarrow$  3 "

$5.380 \times 10^4$   $\leftarrow$  4 "

$$0.01234 \leftarrow 4 \text{ sig. digits.}$$

$= 1.234 \times 10^{-2}$

Trailing zeros count, leading zero do not

Computer store # in binary:

φ . 1.

%	%	%	%	%	%
---	---	---	---	---	---

Finite # of bits = Finite # of digits.

32 bits = 7 sig. digits

64 bits = 15 sig. digits

## Round - off Errors.

type

integer -1, 1, 0, 1, 2

float / double / long double : decimal  
π.

character , 'a' , 'c' , - .

Consider the division of  $\frac{7}{2} = 3.5$

Let 2, 7 and the answer to be  
integer

$$\Rightarrow \frac{7}{2} = 3$$

$$R.O.E. = 0.5$$

In double precision, you set a minimum error at

$$2^{-52} \approx 2.22 \times 10^{-16}$$

In double precision,

$$| \equiv | + 10^{-20}$$

eps : matlab for machine precision

float point

$$a \neq b$$

$$\text{abs}(a - b) < \text{tol}$$

$\uparrow$   
a # I can  
accept

# Condition Number

The condition number of a matrix  $\underline{A}$  tells us how errors in  $\underline{b}$  influences errors in  $\underline{x}$  when solving  $\underline{A} \underline{x} = \underline{b}$

$$\underline{A}(\underline{x} + \alpha \underline{x}) = (\underline{b} + \alpha \underline{b})$$

↑  
induced  
error in  $\underline{x}$

↑  
error in  $\underline{b}$

Goal:

relate  $\frac{\|\alpha \underline{x}\|}{\|\underline{x}\|} \propto \frac{\|\alpha \underline{b}\|}{\|\underline{b}\|}$

$$\cancel{\underline{A} \underline{x}} + \underline{A} \alpha \underline{x} = \cancel{\underline{b}} + \alpha \underline{b}$$

$$\Rightarrow \underline{A} \alpha \underline{x} = \alpha \underline{b}$$

$$\|\underline{e}b\| = \|\underline{A}\underline{e}x\| \leq \|\underline{A}\| \|\underline{e}x\|$$

Also,  $\underline{e}x = \underline{A}^T \underline{e}b$

$$\Rightarrow \|\underline{e}x\| = \|\underline{A}^T \underline{e}b\| \leq \|\underline{A}^T\| \|\underline{e}b\|$$

$$\|\underline{b}\| = \|\underline{A}x\| \leq \|\underline{A}\| \|x\|$$

$$\frac{\|\underline{e}x\|}{\|x\|} \leq \|\underline{A}^T\| \|\underline{e}b\| \frac{\|\underline{A}\|}{\|\underline{b}\|}$$

$$\frac{\|\underline{e}x\|}{\|x\|} \leq \boxed{\|\underline{A}^T\| \|\underline{A}\|} \frac{\|\underline{e}b\|}{\|\underline{b}\|}$$

Define the condition number or  $K(\underline{A})$

$$K(\underline{A}) = \|\underline{A}^T\| \|\underline{A}\|$$

indicates how error grow

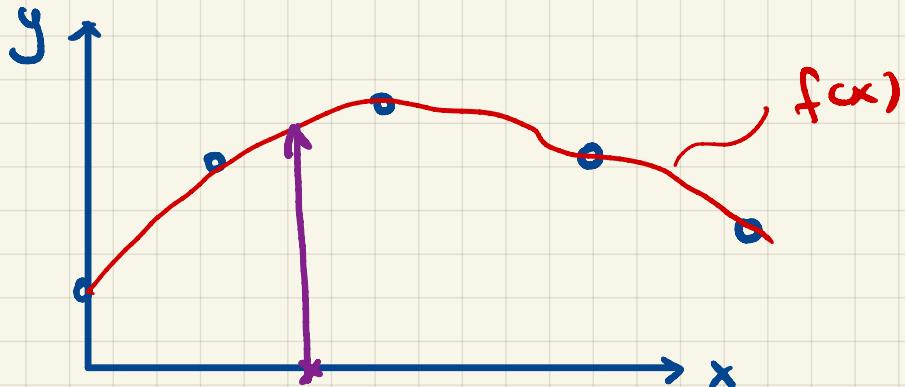
ex).

$$\text{let } \frac{\|\underline{Ax}\|}{\|\underline{b}\|} \sim 10^{-16} \quad \text{with } K(A) \sim 10^6$$

$$\frac{\|\underline{Ax}\|}{\|\underline{x}\|} \leq K(A) \frac{\|\underline{ab}\|}{\|\underline{b}\|} \sim 10^{-10}$$

# Interpolation

Given a data set, how do you determine values between the data.



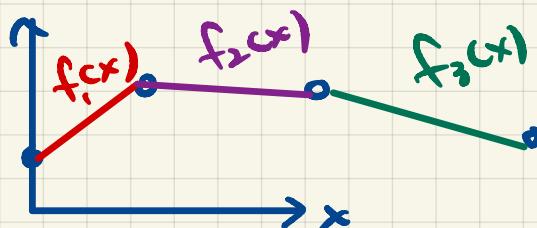
Goal: Find  $f(x)$

① Global interpolation: Use all data

to construct one  $f(x)$

over the entire domain

② Piecewise interpolation : Use individual interpolants between data points



# Polynomial Interpolation - Global

n data pt.

$$(x_1, y_1)$$

:

$$(x_n, y_n)$$

$\Rightarrow$  n-1 polynomial

$$f(x) = a_0 + a_1 x + a_2 x^2 + a_3 x^3 + \dots + a_{n-1} x^{n-1}$$

$a_0 \dots a_{n-1} \Rightarrow$  n coefficients

$$f(x_1) = a_0 + a_1 x_1 + \dots + a_{n-1} x_1^{n-1} = y_1$$

:

$$f(x_n) = a_0 + a_1 x_n + \dots + a_{n-1} x_n^{n-1} = y_n$$

$$\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \dots & x_n^{n-1} \end{bmatrix} \begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_{n-1} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{n-1} \end{bmatrix}$$

Vandermonde Matrix

If  $x_i \neq x_j$  &  $i \neq j$ ,  $\det(V) \neq 0$

$$\Rightarrow \underline{a} = V^{-1} \underline{y}$$

but  $K(V)$  is huge