

Numerical Solutions

Goal is to find \underline{x} such that $A\underline{x} = \underline{b}$

If $\det(A) \neq 0$, A^{-1} exists, $\Rightarrow \underline{x} = A^{-1}\underline{b}$

You never actually find A^{-1} .

2 methods to solve $A\underline{x} = \underline{b}$

- 1) Iteration Methods: Use matrix-vector products & projection to iterate a sequence of vectors \underline{y} that converge to the solution \underline{x}

GMRES, BICGSTAB, QMR, etc.

Most are Krylov Subspace methods.

"Iterative Methods for Sparse Linear Systems" by Yousef Saad.

- 2) Decomposition Methods: Given A find \underline{B} & \underline{C} such that $A = \underline{C}\underline{B}$ then solve $\underline{C}\underline{y} = \underline{b}$ & $\underline{B}\underline{x} = \underline{y}$

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

$$\Rightarrow \underline{C}\underline{B}\underline{x} = \underline{b} \Rightarrow \underline{A}\underline{x} = \underline{b}$$

LU, QR, SVD, etc.

These are "exact" methods

In practice they are not exact due to finite precision

- 1) How computers store information
 - 2) finite precision
 - 3) How errors in A & b influence the errors in x .
-

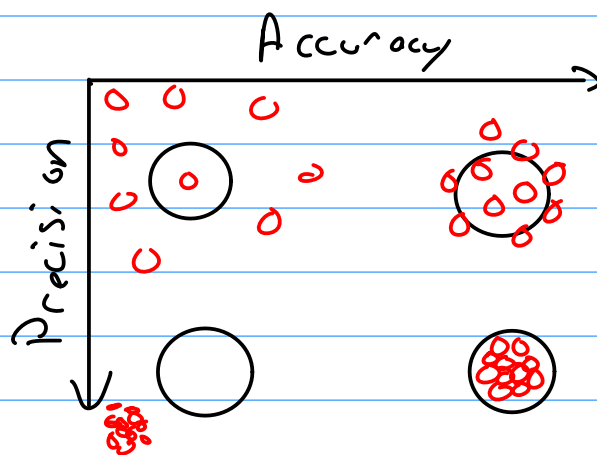
Accuracy vs. Precision

Let the true answer be "truth".

We have approximate solutions

Accuracy: How close are the approximate solutions to truth

Precision: How close are the approximations to each other.



Numerical Errors come from:

- Choice of the model
- Numerical Approximations
- Data representation
- Implementation
- Fluctuations in Hardware

Truth = Approximation + Error

$$x^* = x + \varepsilon$$

$\Rightarrow \varepsilon = x^* - x$ ← not usually a good measure

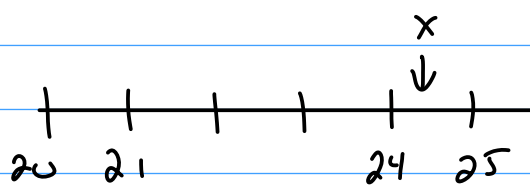
Relative error: $\varepsilon_{rel} = \frac{x^* - x}{x^*}$

Act: In iterative methods you get

$$x_1, x_2, x_3, \dots, x_n \rightarrow \underline{x}$$

$$\underline{r}_i = \underline{A} x_i - \underline{b} = \text{residual (not error!)}$$

Significant Figures



$$24 < x < 25$$

$$x = \cancel{24} \cancel{25}?$$

$$\cancel{24} \cancel{25}?$$

$$24.4!$$

Sig figs tell how much useful information you have: Typically the # of certain digits plus one.

$$\text{ex.) } 53800 \leftarrow 5 \text{ sig figs}$$

$$5.38 \times 10^4 \leftarrow 3 \text{ sig figs}$$

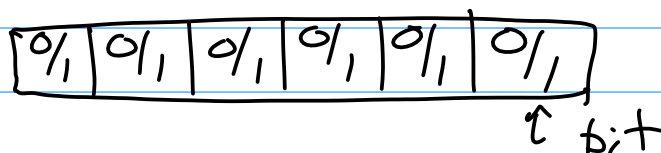
$$5.380 \times 10^4 \leftarrow 4 \text{ sig figs}$$

Trailing zeros count, leading zeros do not.

$$0.01234 \leftarrow 4 \text{ sig figs}$$

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

Computers store numbers in binary: 0 & 1



finite # of bits = finite precision

32-bit : 7 sig figs (single)

64-bit : 15 sig figs (double)

$$\frac{3.25 \times}{1.96 \times} = 1.65 \overline{816} \dots$$

$$\frac{3.259}{1.960} = 1.6627$$

$$\frac{3.250}{1.969} = 1.65$$

Typically, round or chop

1.65432 to 2 sig figs

rounded = 1.66

chop = 1.65

Round-off errors

Data in a computer must have a **t-type**

Integer: ... -1, **0**, 1, 2, ...

float/double/long double: Decimals
1.34, π (up to a #)

Character: 'a', 'b', etc.

Consider $\frac{7}{2} = 3.5$

Let 2 & 7 be integers & the result must also be an integer

$$\Rightarrow \frac{7}{2} = 3 \Rightarrow \text{round-off error of } 0.5$$

In double precision you can represent the difference of numbers larger than

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

In double precision

$$1.0 \times 10^0 = 1.0 \times 10^0 + 10^{-20}$$
$$1. = 1. + 10^{-20}$$

In Matlab: **eps** gives the machine precision

This also why when comparing floating point numbers do not use equals.

$$a \neq b \Rightarrow \text{use } \underbrace{\text{abs}(a-b) < \epsilon)}$$

Round off error example

3 sig fig w/ chopping

$$0.\underline{990} + 0.00\underline{440} + 0.00\underline{490} = 0.9993 \text{ (exact)}$$

$$(0.990 + 0.00440) + 0.00490 \\ 0.994 + 0.00490 = 0.998$$

$$0.990 + (0.00440 + 0.00490) \\ 0.990 + 0.00930 = 0.999$$

Frequent in division & subtraction

You can modify the algorithm to minimize

roots of $x^2 + bx + 1 = 0$ b is large

$$r = (b^2 - 4)^{1/2} \quad x_1 = \frac{b+r}{2} \quad x_2 = \frac{b-r}{2}$$

If b is large $b \approx r$

$$b = 110 \quad r = (110^2 - 4)^{1/2} \approx 109.\overline{9918}$$

$$\text{Truth: } x_1 = 109.99... \\ x_2 = 0.00909166...$$

Solve on 3 sig fig machine + chopping

$$b = 110 \quad r = 109$$

$$x_1 = \frac{b+r}{2.00} = \frac{110+109}{2.00} = \frac{219}{2.00} = 109 \in 3 \text{ digits of accuracy}$$

↑ Actually 109.5

$$x_2 = \frac{b-r}{2} = \frac{110-109}{2.00} = \frac{1.00}{2.00} = 0.500 \in \text{no accurate digits}$$

but, $x_2 = \frac{b-r}{2} = \frac{b-r}{2} \cdot \frac{b+r}{b+r} = \frac{b^2 - r^2}{2(b+r)} = \frac{4}{2(b+r)} = \frac{2}{b+r}$

$$\Rightarrow x_2 = \frac{2.00}{110+109} = \frac{2.00}{219} = 0.00913 \quad \text{vs. } 0.00909$$

↑
one digit of accuracy

Why do we care! Solutions to Ex 2
involve $+$, $-$, $*$, \div

Repeated applications of math operations
introduce errors that accumulate

A **stable algorithm** does not break
down due to these errors

Condition Number

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

Instead of $\underline{A}\underline{x} = \underline{b}$ you solve

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

↑ ↑
error induced error in \underline{b}
by \underline{A}

Goal: Relate $\frac{\|\underline{\Delta x}\|}{\|\underline{x}\|} \propto \frac{\|\underline{\Delta b}\|}{\|\underline{b}\|}$

$$\cancel{\underline{A}\underline{x}} + \underline{A}\underline{\Delta x} = \cancel{\underline{b}} + \underline{\Delta b} \Rightarrow \underline{A}\underline{\Delta x} = \underline{\Delta b}$$

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

$$\text{Also, } \underline{\Delta x} = \underline{A}^{-1}\underline{\Delta b} \Rightarrow \|\underline{\Delta x}\| = \|\underline{A}^{-1}\underline{\Delta b}\| \leq \|\underline{A}^{-1}\| \|\underline{\Delta b}\|$$

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

$$\frac{\|\underline{\Delta x}\|}{\|\underline{x}\|} \leq \|\underline{A}^{-1}\| \|\underline{\Delta b}\| \frac{\|\underline{A}\|}{\|\underline{b}\|}$$

$$\underbrace{\frac{\|\underline{\Delta x}\|}{\|\underline{x}\|}} \leq \underbrace{(\|\underline{A}^{-1}\| \|\underline{A}\|)}_{\uparrow} \underbrace{\frac{\|\underline{\Delta b}\|}{\|\underline{b}\|}}$$

Define the condition number of A as

$$\kappa(A) = \|A\| \|A^{-1}\|$$

$\kappa(A)$ indicates how errors grow

ex.) let $\frac{\|\Delta b\|}{\|b\|} \sim 10^{-16}$ w/ $\kappa(A) \sim 10^{10}$

$$\frac{\|\Delta x\|}{\|x\|} \leq \kappa(A) \frac{\|\Delta b\|}{\|b\|} \sim 10^{-6}$$

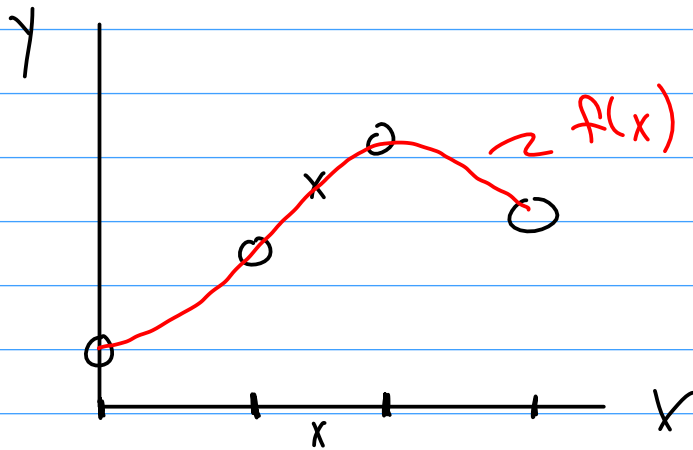
In matlab: `cond`, `rcond`, `condst`

Also why use backslash: $x = A \setminus b$, not

$$x = \text{inv}(A) * b$$

Interpolation

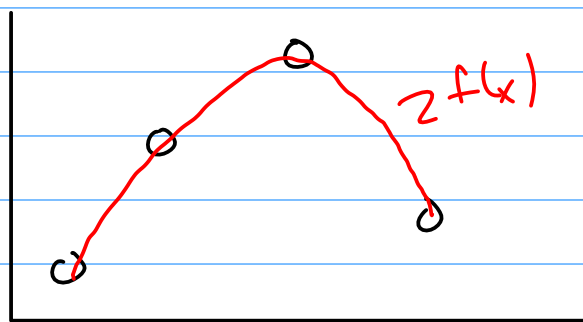
Given a data set, determine values not in the data set.



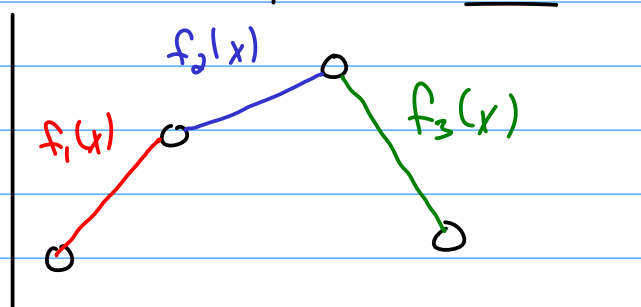
Goal: Get $f(x)$

Two main types:

- 1) Global interpolants: Use all of the data to construct one $f(x)$ over the entire range



- 2) Piecewise Interpolation: One equation between data points



Polynomial Interpolation (Global)

let n -data points $(x_1, y_1) \rightarrow (x_n, y_n)$
be given w/ $x_i \neq x_j$ if $i \neq j$

Construct an $n-1$ polynomial:

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

$a_0 \rightarrow a_{n-1}$ unknowns

Require that $f(x_1) = y_1$, $f(x_2) = y_2$, etc.

$$f(x_1) = a_0 + a_1 x_1 + a_2 x_1^2 + \dots + a_{n-1} x_1^{n-1} = y_1$$

$$f(x_2) = a_0 + a_1 x_2 + a_2 x_2^2 + \dots + a_{n-1} x_2^{n-1} = y_2$$

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

\uparrow known \uparrow unknown

$$\underbrace{\begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & \dots & \dots & x_n^{n-1} \end{bmatrix}}_{\underline{V}} \underbrace{\begin{bmatrix} a_0 \\ a_1 \\ \vdots \\ a_{n-1} \end{bmatrix}}_{\underline{a}} = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\underline{y}}$$

If $x_i \neq x_j$ if $i \neq j$, $\det(\underline{V}) \neq 0$

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

but, $k(\underline{V})$ is huge!

\underline{V} is called the Vandermonde matrix

ex.) let $(300, 1)$, $(400, 1.1)$, $(500, 1.3)$
be given

$$\underbrace{\begin{bmatrix} 1 & 300 & 90000 \\ 1 & 400 & 160000 \\ 1 & 500 & 250000 \end{bmatrix}}_{k(\underline{V}) \sim 6 \times 10^6} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} 1 \\ 1.1 \\ 1.3 \end{bmatrix}$$

It turns out that while the polynomial interpolant is unique, the method to compute it is not.

Next time: Lagrange Interpolation

$$f(x) = \gamma_1 L_1(x) + \gamma_2 L_2(x) + \dots + \gamma_n L_n(x)$$

$L_i(x)$ = Lagrange function / polynomial

