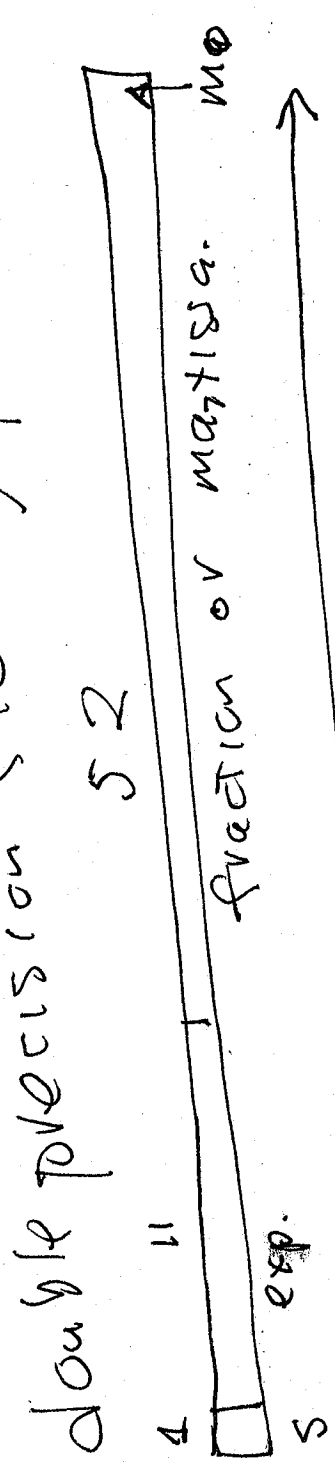


int. format

CADS 21

IEEE 754 = Binary 64
double precision floating point representation



$$2 \times 10^{-23}$$

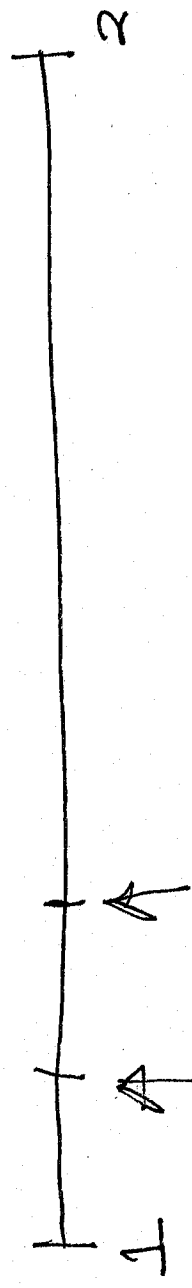
$$X = (-1)^s (1.M_{51} \dots M_0)_2 \times 2^E$$

$$0 \leq E \leq 2047$$

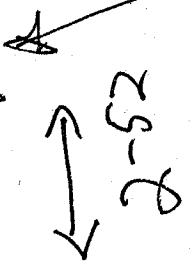
$$X_{MAX} \approx 1.79 \times 10^{308}$$

$$X_{MIN} \approx 2.23 \times 10^{-309}$$

$$0 < X_{MIN}$$

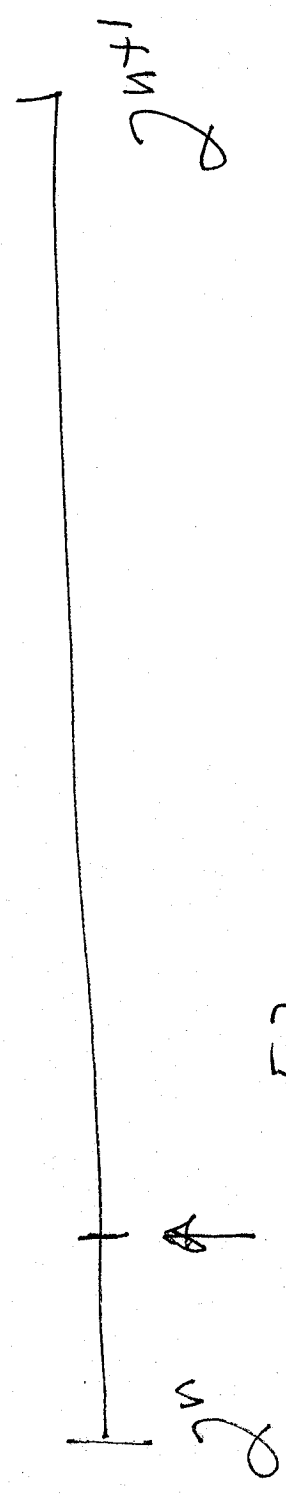


$1 + 2^{-52} = 1 + 2^{-52} \cdot 2^{52} = 2^{52} \cdot 2^{-52} + 2^{52} \cdot 2^{-52} = 2^{52} \cdot (2^{-52} + 2^{-52}) = 2^{52} \cdot 2^{-51} = 2^{52-51} = 2^1 = 2$



next floating point numbers

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



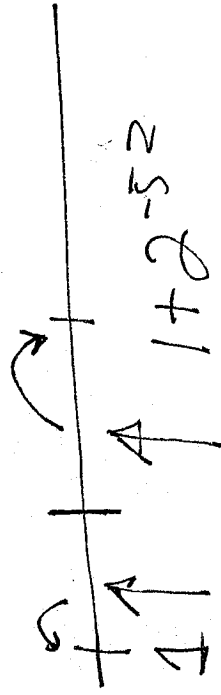
$1 + 2^{n-52}$

$2^{52} \rightarrow 1 + 2$

53
2

$F = \text{floating point numbers} - \text{finite set}$

Given $x \in \mathbb{R}$, $fl(x)$ is the nearest floating point number via rounding



Basic Axiom - for all $x \in \mathbb{R}$, $fl(x)$ satisfies

$$\frac{|x - fl(x)|}{|x|} \leq \frac{1}{2} 2^{-52} \approx 1.1 \times 10^{-16}$$

machine epsilon

$$= \epsilon_{\text{mac}}$$

Warning: IBM

Sammit

4

200 PLOPS

= 200 peta floating point operations / second

$$200 \times 10^{15}$$

$$\boxed{2 \times 10^{17}}$$

each flop can yield 1×10^{-16} error

$$\underline{1 \text{ second}} \quad (2 \times 10^{17}) (1 \times 10^{-16})$$

$$= \underline{\underline{20 \text{ error}}}$$

Stability

5

~~Inputs~~
Inputs

$f \rightarrow y$

Soln.

$$(A, b) \mapsto x \text{ with } Ax = b$$

$F \rightarrow$ into floats.
Specific algorithm

$$F(x) \xrightarrow{F/f} F(y)$$

Stability = fairly accurate translation

Least Squares

$M \times N$

$m > n$

$$\text{Solve } A\vec{x} = \vec{b}$$

$$A_{11}x_1 + A_{12}x_2 + \dots + A_{1n}x_n = b_1$$

|

|

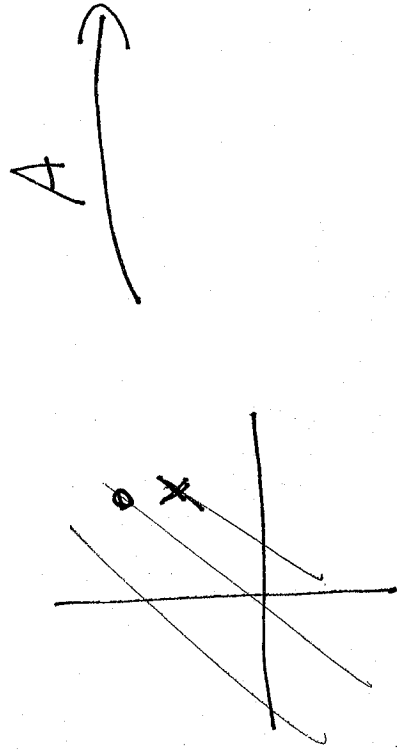
|

$$A_{m1}x_1 + A_{m2}x_2 + \dots + A_{mn}x_n = b_m$$

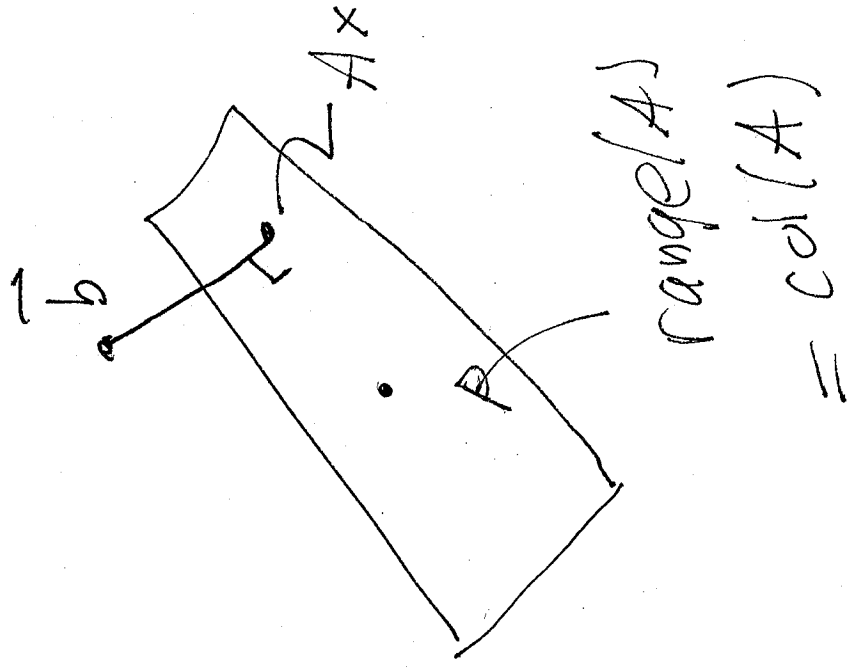
more equations than unknowns.

As a map $A: \mathbb{R}^n \rightarrow \mathbb{R}^m$ $m > n$.

$n=2$ $m=3$.



$A \rightarrow$



- ① Find x which minimizes $\|Ax - b\|_2$
- ② Ax will be the perpendicular projection of b onto $\text{col}(A)$ so $Ax = P\vec{b}$

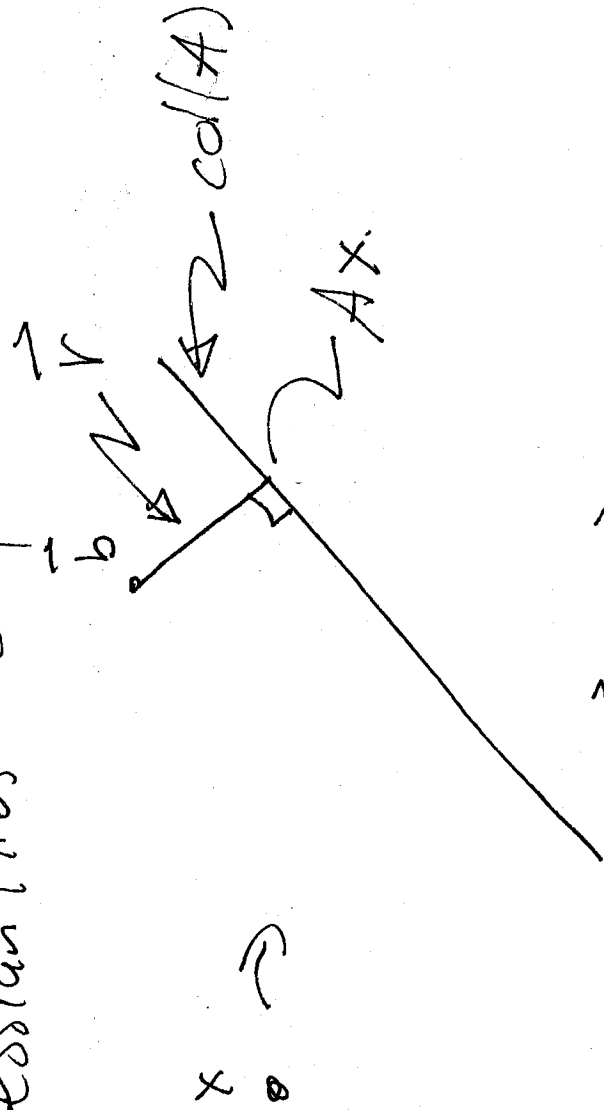
For (1) Let $\Phi(x) = \|Ax - b\|_2^2 \Rightarrow$

Minimize Φ via calculus i.e. Find x_0

with $\nabla \Phi(x_0) = 0$. turns out that

at x_0 the Hessian $\nabla^2 \Phi(x_0)$ is pos def so min.

(2)



$$\text{residual} = \vec{r} = \vec{b} - Ax$$

For minimum, $\vec{r} \perp \text{col}(A)$

so if $A = \begin{bmatrix} a_1 & a_2 & \dots & a_n \end{bmatrix}$

$$\vec{r} \perp \text{col}(A)$$

$$\vec{r} \perp a_j \text{ all } j$$

$$a_j^T \vec{r} = 0 \text{ all } j$$

$$\begin{bmatrix} a_1^T \vec{r} \\ \vdots \\ a_n^T \vec{r} \end{bmatrix} = \vec{0}$$

$$=$$

$$\begin{bmatrix} a_1^T \\ a_2^T \\ \vdots \\ a_n^T \end{bmatrix} \vec{r}$$

$$A^T \vec{r} =$$

So Minimum of \hat{r} is Characterized 10

by $A^T \hat{r} = \vec{0}$

$$A^T (\vec{b} - A\hat{x}) = 0$$

normal equations

$$A^T \vec{b} = A^T A x$$

$$A^T A \hat{x} = A^T \vec{b}$$

first try, compute $(A^T A)^{-1}$

$$\Rightarrow \hat{x} = (A^T A)^{-1} A^T \vec{b} \quad \underline{\underline{\text{unique soln}}}$$

Theorem: $m \geq n$, $A^T A$ is invertible "

if and only if A has full rank i.e. $\text{rank} = n$
i.e. n ind col.

We prove the contra positive.

$A^T A$ is singular $\Rightarrow A$ is rank deficient

Assume $A^T A$ sing \Rightarrow there is $\vec{v} \neq 0$ with:

$$(A^T A) \vec{v} = \vec{0} \Rightarrow V^T A^T A \vec{v} = \vec{0}$$

$$\text{so } (A \vec{v})^T (A \vec{v}) = 0 \text{ or } \|A \vec{v}\|_2^2 = 0$$

$$\Rightarrow A \vec{v} = 0$$

so A is rank def since $\vec{v} \neq 0$.

Converse

$\text{rank}(A) < n \Rightarrow$ There is

12

a $\vec{v} \neq 0$ with $A\vec{v} = 0 \Rightarrow$

$A^T A \vec{v} = 0 \Rightarrow A^T A$ is singular

Since $\vec{v} \neq 0$.

$m \geq n$

and $\text{rank}(A) = n$

$I \neq$

$\Rightarrow A^+ = (A^T A)^{-1} A^T$ pseudoinverse

\Rightarrow soln to the normal eq. is

$\vec{x} = A^+ \vec{b}$.

Lemma If $M=N$ and A is invertible

then

$$A^+ = A^{-1}$$

13

$$\begin{aligned} \text{Pf} & \quad (A^T A^{-1} A^T) A \\ & = A^{-1} (A^T A^T) A \\ & = A^{-1} A = I. \end{aligned}$$