

$$V \xrightarrow{T} W$$

$$V_{B_1} \xrightarrow[B_1]{[T]} W_{B_1'}$$

$$V_{B_2} \xrightarrow[B_2]{[T]} W_{B_2'}$$

$$P = \begin{bmatrix} I \\ B_2 & B_1 \end{bmatrix} \quad [I]_{B_2'} = Q.$$

$$B_2 \begin{bmatrix} T \\ B_2' \end{bmatrix}_{B_2'} = Q_{B_1'} \begin{bmatrix} T \\ B_1 \end{bmatrix} P^{-1}$$

$$V = \mathbb{R}^3 \xrightarrow{T} \mathbb{R}^2$$

$$T \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} x+y \\ y-z \end{pmatrix}$$

$$\boxed{B_1 = \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right\} \quad B_1' = \left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}}$$

$$B_2 \begin{bmatrix} T \\ B_2' \end{bmatrix} = \left[T \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}_{B_2} \mid T \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}_{B_2} \mid T \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}_{B_2} \right]$$

$$= \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & -1 \end{bmatrix}$$

$$B_2 = \left\{ \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} \right\} \quad B_2' = \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix} \right\}$$

$$P^{-1} = \begin{bmatrix} I \\ B_1 & B_2 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad Q = \begin{bmatrix} I \\ B_2' & B_1' \end{bmatrix} = \begin{bmatrix} Y_2 & Y_L \\ Y_L & -Y_2 \end{bmatrix}$$

$$\boxed{\begin{bmatrix} T \\ B_2 & B_2' \end{bmatrix} = Q \begin{bmatrix} T \\ B_1' & B_L \end{bmatrix} P^{-1}}$$

* If $V=W$, then
 $Q=P$

$$\boxed{\begin{bmatrix} T \\ B_2 & B_2' \end{bmatrix} = Q \begin{bmatrix} T \\ B_1' & B_1 \end{bmatrix} Q^{-1}}$$

* $A \in M_{n \times m}(\mathbb{R})$

$A: \mathbb{R}^m \rightarrow \mathbb{R}^n$
 (a_{ij})

$$\begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix} \mapsto \begin{pmatrix} \sum_{j=1}^m a_{1j}x_j \\ \vdots \\ \sum_{j=1}^m a_{nj}x_j \end{pmatrix}$$

$\left\{ \begin{array}{l} \text{linear} \\ \text{transformation} \end{array} \right\} \longleftrightarrow \left\{ \begin{array}{l} \text{all} \\ \text{matrices} \end{array} \right\}$

* Eigen values and Eigen vectors

- Eigen value, vectors corresponding to distinct eigen values are linearly independent.
- * Let $E_\lambda = \text{subspace spanned by eigen vectors of } \lambda$
- Algebraic multiplicity of eigen values = it's multiplicity as a root of characteristic polynomial.

geometric multiplicity of $\lambda = \dim E_\lambda$

In general geom. multiplicity (α) \leq alg. multiplicity of λ .

* Diagonalizability: A is diagonalizable if \exists invertible matrix P s.t. $A = P^{-1} \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{pmatrix} P$.

eigen values of A $\lambda_1, \lambda_2, \dots, \lambda_k$

alg. mult of $\lambda_i = 1 + i$ (alg. mult of $\lambda_j \geq 1$ for some j)

\downarrow
 A is diagonalizable.

alg. mult = geom. mult
+ λ_j

$|$
 A is diagonalizable

alg. mult \neq
geom. mult
for some

λ_j

A is not
diagonalizable

When A is not diagonalizable.

"almost" diagonal JCF
diag + N.

* Similar matrices?

* Inner products

(1) $V \times V \rightarrow \mathbb{R}$

$v, w \mapsto \langle v, w \rangle$

① $\langle v, w \rangle \geq 0$. $\langle v, v \rangle = 0 \Leftrightarrow v=0$.

②

* Norms

$\| \cdot \| : V \rightarrow \mathbb{R}_{\geq 0}$

$$\|v\| = \sqrt{\langle v, v \rangle}$$

* Standard inner product on V

① Over \mathbb{R}^n : dot product

$$\begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} = \sum x_i y_i$$

② Over \mathbb{C}^n :

$$\sum x_i \bar{y}_i$$

* Different norms on vectors

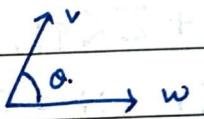
$$\textcircled{1} \quad \|v\| = \langle v, v \rangle^{1/2} = (\sum v_i^2)^{1/2}$$

\textcircled{2} Holder norm

$$\textcircled{3} \quad v = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix} \quad \|v\| = (\sum |v_i|^p)^{1/p} \quad p \in \mathbb{Z}_{>0}$$

* Angles between vectors

$$\cos \theta = \frac{\langle v, w \rangle}{\|v\| \|w\|}$$



If $\theta = 90^\circ$, then we say v is orthogonal to w .

$$v \perp w \Leftrightarrow \langle v, w \rangle = 0$$

* Orthogonal set

* Orthonormal set.

linearly independent?

Given any set $S = \{v_1, \dots, v_n\}$, we can generate a set of orthonormal vectors $\{q_1, \dots, q_n\}$ such that

$$\text{span } \{v_1, \dots, v_k\} = \text{span } \{q_1, \dots, q_k\}$$

General fact:

If $\{q_1, \dots, q_n\}$ is a set of orthonormal vectors in V ($\dim V = m$), then for every $v \in V$ consider

$$r = v - \langle q_1, v \rangle q_1 - \langle q_2, v \rangle q_2 - \dots - \langle q_n, v \rangle q_n.$$

— / —

Then $\langle r, q_1 \rangle = 0$.

$$\langle r, q_2 \rangle = 0$$

$$\vdots$$
$$\langle r, q_n \rangle = 0.$$

r is orthogonal to all these vectors.

$$r = v - \sum \langle q_i, v \rangle q_i$$

Thus, for any $v \in V$, we may express v as

$$v = r + \sum \langle q_i, v \rangle q_i$$

$$\text{where } \langle r, q_i \rangle = 0 \quad \forall q_i$$

orthogonal decomposition of v

If the set $\{q_1, \dots, q_n\}$ is a basis for V then $r=0$ &

$$v = \sum_{i=1}^n \langle q_i, v \rangle q_i$$

(orthogonal decomposition of v).

- Linear algebra Sahai Bist.
- Schaum's outline - (linear algebra)
(Lipschitz).

11

* Triangulable matrices

→ Every matrix over \mathbb{C} is triangulable.

* Self adjoint operators

Symmetric (over \mathbb{R}) $A = A^t$

Hermitian (over \mathbb{C}) $A = A^*$

adjoint of a linear transformation.

$$\langle Tx, y \rangle = \langle I, Ty \rangle$$

* Symmetric / Hermitian matrix (self adjoint)

① Eigen values of a self adjoint matrix are real, non-negative

* Eigen vectors of a symmetric matrix are linearly independent and orthogonal.

→ Spectral theorem for self-adjoint matrices

* If A is self-adjoint then \exists invertible matrix P s.t.
 $P^{-1} A P$ is diagonal.

Moreover P is orthogonal (unitary).

* Vector norms

Suppose V is a vector space over \mathbb{C}

Defn: A norm is a function $\|\cdot\| : V \rightarrow \mathbb{R}$ s.t.

$$\textcircled{1} \quad \|\mathbf{v}\| \geq 0, \quad \|\mathbf{v}\| = 0 \text{ iff } \mathbf{v} = 0.$$

$$\textcircled{2} \quad \|\alpha \mathbf{v}\| = |\alpha| \|\mathbf{v}\| \quad \forall \alpha \in \mathbb{C} \text{ & } \mathbf{v} \in V.$$

$$\textcircled{3} \quad \|\mathbf{v} + \mathbf{w}\| \leq \|\mathbf{v}\| + \|\mathbf{w}\| \quad \forall \mathbf{v}, \mathbf{w} \in V.$$

Example:

p-norm on V

$$\text{let } \mathbf{v} = \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}, \quad \|\mathbf{v}\|_p = \left(\sum |v_i|^p \right)^{\frac{1}{p}}$$

$$\text{if } p=1$$

$$\|\mathbf{v}\|_1 = \sum |v_i|$$

$$\text{if } p=\infty$$

$$\|\mathbf{v}\|_\infty := \max_i \{|v_i|\}$$

$$\text{if } p=2$$

$$\|\mathbf{v}\|_2 = \left(\sum |v_i|^2 \right)^{\frac{1}{2}}$$

(Euclidean norm).

$$p_1 \geq 1 \text{ & } p_2 \geq 1$$

$$m \|\mathbf{v}\|_{p_2} \leq \|\mathbf{v}\|_{p_1} \leq n \|\mathbf{v}\|_{p_2}$$

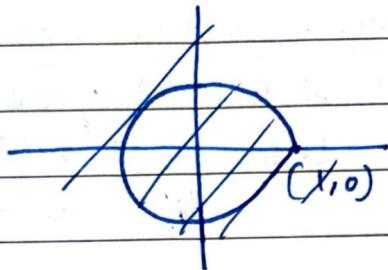
*

Closed unit balls in p-norms

$$B = \{ \mathbf{v} \in V \mid \|\mathbf{v}\| \leq 1 \}$$

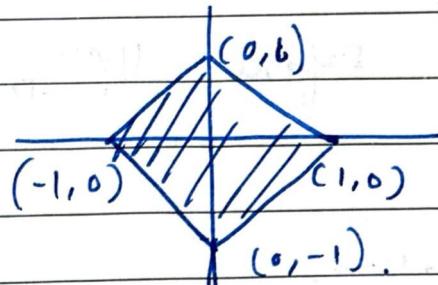
- $p=1$

$$v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix} \quad \|v\| = |v_1| + |v_2|$$



$$\|v\| = |v_1| + |v_2|$$

$$z = x + y$$



- $p=2$.

$$v = \begin{pmatrix} v_1 \\ v_2 \end{pmatrix}$$

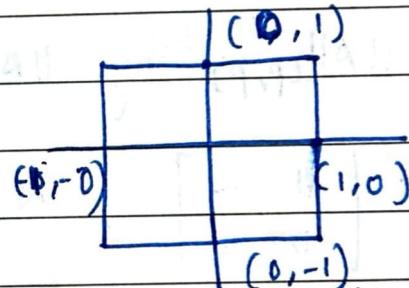
$$\|v\| = \sqrt{|v_1|^2 + |v_2|^2}$$



- $p=\infty$

$$\|v\|_\infty = \max \{ |v_1|, |v_2| \}$$

$$1 \leq p \leq \infty$$



Matrix Norm

Let $A \in M^{n \times m}(\mathbb{C})$

Then a matrix norm is a function $\|\cdot\| : M^{n \times m}(\mathbb{C}) \rightarrow \mathbb{R}$

Satisfying (1) $\|A\| \geq 0$, $\|A\| = 0 \Leftrightarrow A = 0$.

(2) $\|\alpha A\| = |\alpha| \|A\| \quad \forall \alpha \in \mathbb{C}$

(3) $\|A+B\| \leq \|A\| + \|B\| \quad \forall A, B \in M^{n \times m}$.

Additionally if AB is defined

$\|AB\| \leq \|A\| \|B\|$ whenever AB is defined.

Example:

Induced matrix norm

Let $A \in M^{n \times m}(\mathbb{C})$

then $A : {}^p \mathbb{C}^m \rightarrow {}^q \mathbb{C}^n$

$$\|x\|_p \quad x \longrightarrow Ax \quad \|Ax\|_q$$

$$\text{Define } \|A\|_{(p,q)} = \sup_{\substack{x \in \mathbb{C}^m \\ x \neq 0}} \frac{\|Ax\|_q}{\|x\|_p}$$

Since,

$$\rightarrow \|Ax\| = |\alpha| \cdot \|x\| \quad \text{we can define } \|A\|_{(p,q)} = \sup_{\|x\|_p=1} \|Ax\|_q$$

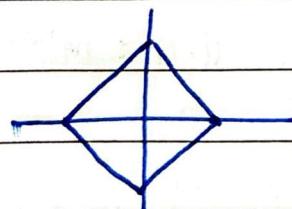
In practice, usually $p=q$, and we will denote

$$\|A\|_{(p,p)} \text{ by } \|A\|_p.$$

eg: $A = \begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix}, A : \mathbb{R}^2 \rightarrow \mathbb{R}^2$

$$\begin{bmatrix} 1 & 2 \\ 0 & 2 \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x+2y \\ 2y \end{pmatrix}$$

$$p=1$$



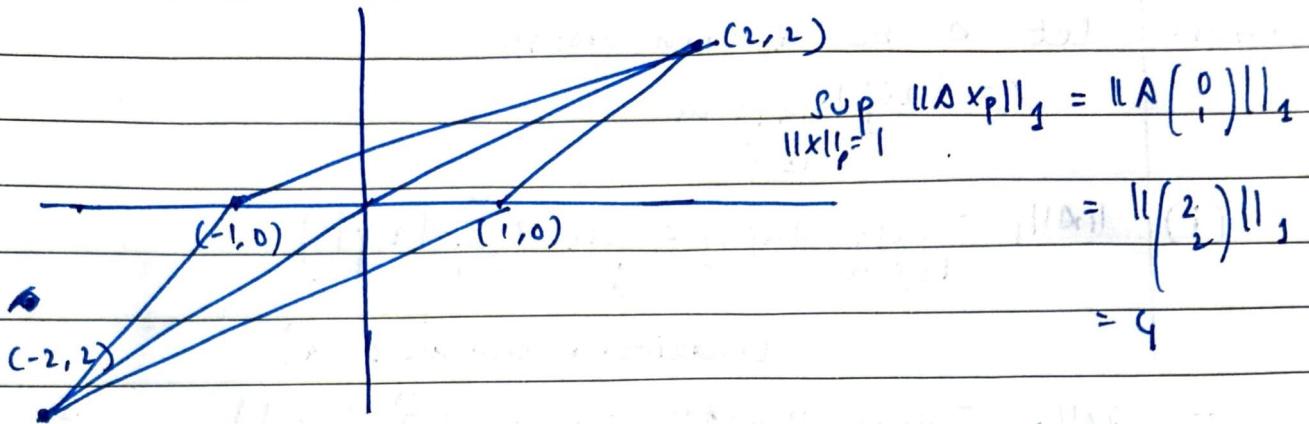
$$\mathbb{R}^2 \xrightarrow{A} \mathbb{R}^2$$

$$(-2, -2)$$

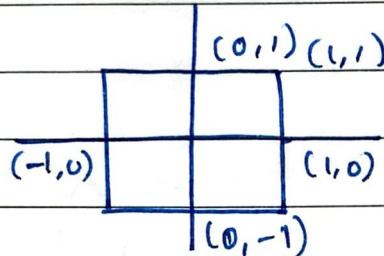
$$(2, 2)$$

$$(1, 0)$$

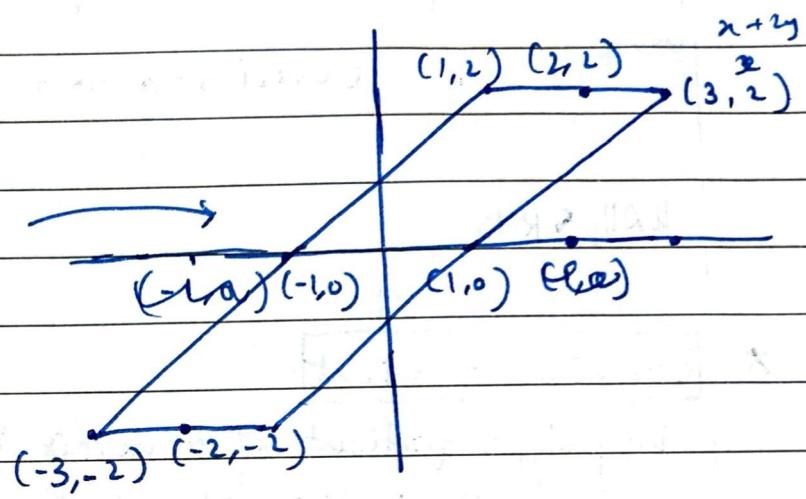
$$(-1, 0)$$



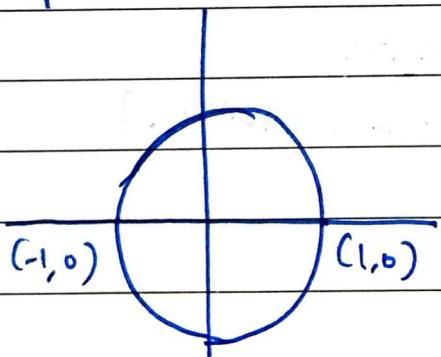
$$p=\infty.$$



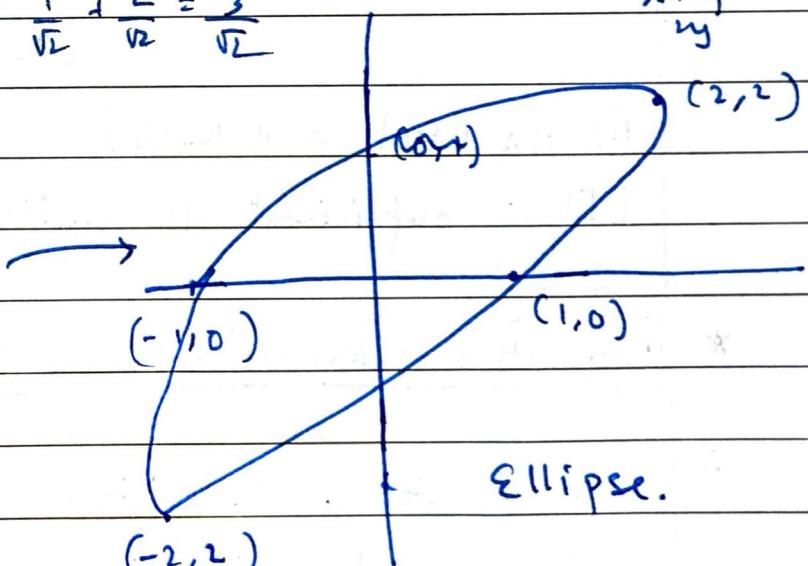
$$\sup_{\|x\|_\infty=1} \|Ax\|_2 = 3.$$



$$\frac{1}{\sqrt{2}} + \frac{1}{\sqrt{2}} = \frac{3}{\sqrt{2}}$$



$$\|A\|_2$$



$$\sup_{\|x\|_2=1} \|Ax\|_p = 2\sqrt{2} \approx 2.82?$$

Theorem: let A be a $m \times n$ matrix

$$A = (a_{ij}) \quad \begin{matrix} 1 \leq i \leq m \\ 1 \leq j \leq n \end{matrix}$$

$$(i) \|A\|_1 = \max_{1 \leq j \leq n} \|a_{j*}\|_1 = \max_j \left(\sum_{i=1}^n |a_{ij}| \right)$$

(maximum column sum)

$$(ii) \|A\|_\infty = \max_i \|a_i^*\|_1 = \max_i \left(\sum_{j=1}^n |a_{ij}| \right) \quad a_i = i^{th} \text{ row of } A.$$

(maximum row sum).

$$\|A\|_1 \leq \text{RHS}$$

* Rayleigh quotient

Rayleigh quotient of a matrix A is defined as

$$R_A : V \setminus \{0\} \rightarrow \mathbb{C}$$

$$v \mapsto \frac{v^* A v}{v^* v} = \frac{\langle Av, v \rangle}{\langle v, v \rangle} \quad v \neq 0$$

$$\text{① } R_A(\alpha v) = \alpha R_A(v) \quad \forall \alpha \neq 0.$$

(It is sufficient to consider v s.t. $\|v\|=1$.)

* Properties of Rayleigh quotient

i) If A is a Hermitian matrix, with e-values $\lambda_1 \leq \lambda_2 \dots \leq \lambda_n$ and associated eigen vectors p_1, \dots, p_n satisfying

$$p_i^* p_j = \delta_{ij} \quad (\delta_{ij} \text{ (Kronecker delta)})$$

$$\delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

vectors
are
orthonormal

then $R_A(\lambda_K)$

$$1) R_A(P_K) = \lambda_K$$

$$2) \lambda_K = \max_{v \in V_K} R_A(v)$$

$$V_K = \text{span}\{p_1, \dots, p_K\}$$

$$3) \lambda_K = \min_{v \in V_{K-1}} R_A(v), V_{K-1}^{\perp} = \text{span}\{p_K, \dots, p_n\}$$

$$v = v_{K-1} \oplus v_{K-1}^{\perp}$$

direct sum.

Theorem: A is Hermitian with e-values $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$

4 e-vectors $\{p_1, \dots, p_n\}$ s.t. $\langle p_i, p_j \rangle \delta_{ij}$

then

$$(i) R_A(P_K) = \lambda_K$$

$$(ii) \lambda_K = \max_{v \in V_K} R_A(v) \quad V_K = \text{span}\{p_1, \dots, p_K\}$$

Proof:

let $V = [P_1 | P_2 | \dots | P_n]$ then $V^* A V = \text{diag}(\lambda_i) = D$

for an $v \in V, v \neq 0$, let $U^{-1} v$ so that $v = Uw$

$$\begin{aligned} R_A(v) &= \frac{v^* A v}{v^* v} = \frac{(Uw)^* A (Uw)}{(Uw)^* (Uw)} = \frac{w^* (U^* A U) w}{w^* w} \\ &= w^* D w = R_D(w). \end{aligned}$$

$$\text{for } v \in V_K, v = \begin{pmatrix} \vdots \\ \alpha_1 \\ \vdots \\ \alpha_K \\ \vdots \\ 0 \end{pmatrix}$$

$$v = \sum_{i=1}^K \alpha_i p_i, w = U^{-1} v = \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_K \\ \vdots \\ 0 \end{pmatrix}$$

$$R_A(v) = R_D \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_K \\ \vdots \\ 0 \end{pmatrix} = (\bar{\alpha}_1, \dots, \bar{\alpha}_K, 0, \dots, 0) \begin{pmatrix} \lambda_1 & & & & \alpha_1 \\ & \ddots & & & \vdots \\ & & \lambda_n & & \alpha_n \\ & & & \ddots & \vdots \\ & & & & 0 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_K \\ \vdots \\ 0 \end{pmatrix}$$

$$(x_1 - \dots - x_k \ 0 \ \dots \ 0) \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_k \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \frac{\sum_{i=1}^k \lambda_i |x_i|^2}{\sum_{i=1}^n |x_i|^2}$$

(1) $R_A(P_K) = \lambda_K$ since $P_K = \sum_{i=1}^K x_i p_i$
 where $x_K = 1$
 $\quad \quad \quad \alpha_i = 0 \ \forall i \neq K.$

for $v \in V_K$

$$R_A(v) = \frac{\sum_{i=1}^K \lambda_i |x_i|^2}{\sum_{i=1}^n |x_i|^2}$$

$$(2) \max_{v \in V_K} R_A(v) = \max_{v \in V_K} \left\{ \frac{\sum_{i=1}^K \lambda_i |x_i|^2}{\sum_{i=1}^n |x_i|^2} \right\}$$

Recall: $\lambda_0 \leq \lambda_1 \leq \dots \leq \lambda_n$

each:

$$\rightarrow \lambda |x_i|^2 \leq \lambda_K |x_i|^2 \ \forall i \leq K$$

$$\rightarrow \sum \lambda |x_i|^2 \leq \lambda_K (\sum |x_i|^2).$$

$$\forall v \in V_K \quad R_A(v) \leq \lambda_K \Rightarrow \lambda_K = \max_{v \in V_K} \{R_A(v)\}$$

* Calculating $\|A\|_2$

Let A be a $n \times n$ matrix

$$\|A\|_2 = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\|Ax\|_2}{\|x\|_2} = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \sqrt{\langle Ax, Ax \rangle}$$

Squaring

$$\|A\|_2^2 = \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{\langle Ax, Ax \rangle}{\langle x, x \rangle} = \frac{(Ax)^* (Ax)}{x^* x}$$

$$= \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} \frac{x^* A^* Ax}{x^* x}$$

$$= \sup_{\substack{x \in \mathbb{C}^n \\ x \neq 0}} R_{A^* A}(x).$$

$A^* A$ is hermitian, always, so

Hence $R_{A^* A}(v) \leq$ Largest e-value of $A^* A$

(attained at p_1)

$$\therefore \text{Consequently } \sup R_{A^* A}(x) = p(A^* A) = \|A\|_2^2$$

Note that: $A^* A$ is always positive semi-definite

$$\boxed{\|A\|_2 = \sqrt{R_{A^* A}(A^* A)}} \rightarrow \text{largest singular value of } A.$$

1. Solve linear eqs.
 2. Compute eigenvectors & eigenvalues

Aim of
Numerical linear
algebra.

* Singular Value Decomposition

$$A = U \Sigma V^*$$

mn. \downarrow diagonal with singular
unitary value on the diagonal.

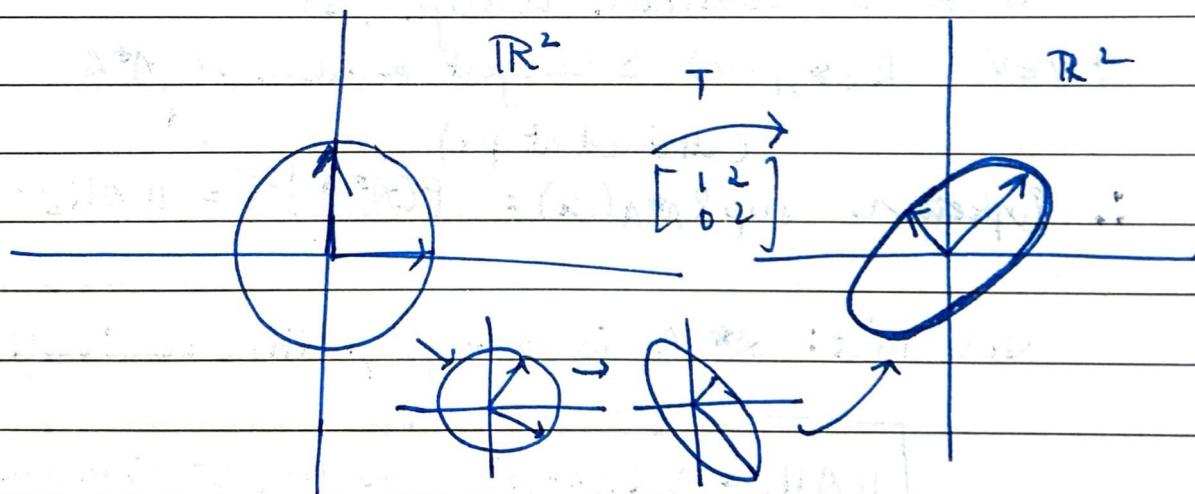
$$V^* V = I$$

$$V^{-1} = V^*$$

The cols of V are orthogonal for \mathbb{R}
Unitary if working over \mathbb{C}

→ Orthogonal (over \mathbb{R}) $\Rightarrow O^t O = O O^t = I$

→ Unitary (over \mathbb{C}) $\Rightarrow V^* V = V V^* = I$.



Any linear transformation from $\mathbb{R}^m \rightarrow \mathbb{R}^n$

$$\mathbb{R}^m \xrightarrow{T} \mathbb{R}^n$$

unit ball \mapsto hyperellipse
(n-dimensional ellipsoid).

$$A(x) = U \sum V^*(x)$$

↴ rotation
scaling / stretching
 ↴ rotating

$$A_{mn} = U \Sigma V^*$$

↓
 diagonal
 singular
 values on
 diagonal)
 unitary ←

$$A : \mathbb{R}^n \rightarrow \mathbb{R}^m.$$

$$A_{\max} = V_{\max} \sum_{n \geq 0} V^*(n)$$

$r = \text{rank}(A) = \#$ non-zero singular values of Σ . ($\sigma_1 - \sigma_r$)
 $m > n$ tall matrix $\Sigma \rightarrow$ singular values $\neq 0$.

$$\begin{array}{c}
 \text{Full} \\
 \text{Estimate} \\
 \text{SVD}
 \end{array}
 \left[\begin{array}{c} A \\ \hline m \\ n \end{array} \right] = \left[\begin{array}{c} \Sigma \\ \hline r \text{ cols} \\ m \times n \end{array} \right] \left[\begin{array}{c} U \\ \hline m \\ n \end{array} \right] \left[\begin{array}{c} V^T \\ \hline n \\ m \end{array} \right]$$

$$m \leq n$$

$$\begin{bmatrix} & & & \\ & & & \\ & & & \end{bmatrix}_{m \times n} = \begin{bmatrix} & & & \\ & & & \\ & & & \end{bmatrix}_{n \times n} \begin{bmatrix} & & & \\ & & & \\ & & & \end{bmatrix}_{m \times n}$$

$m \times n$ matrix A ($m \geq n$)

$$A = U \sum_{\text{diagonal}} V^*$$

unitary

$$\begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix}$$

U_i : left singular vectors of A

σ_i : singular values of A

V_i : right singular value vectors of A .

$$\begin{bmatrix} U_1 & | & U_2 & | & \cdots & | & U_m \end{bmatrix}_{m \times m} \begin{bmatrix} \sigma_1 \\ \sigma_2 \\ \vdots \\ \sigma_n \end{bmatrix}_{n \times n} \begin{bmatrix} V_1^* \\ V_2^* \\ \vdots \\ V_m^* \end{bmatrix}_{m \times n}$$

$$A = U \Sigma V^*$$

$$AV = U\Sigma$$

$$A \begin{bmatrix} v_1 & v_2 & \dots & v_n \end{bmatrix} = \begin{bmatrix} u_1 & u_2 & \dots & u_n \end{bmatrix} \begin{bmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_n & \end{bmatrix}$$

$$\begin{bmatrix} Av_1 & Av_2 & \dots & Av_n \end{bmatrix} = \begin{bmatrix} \sigma_1 u_1 & \dots & \sigma_n u_n \end{bmatrix}$$

$Av_i = \sigma_i u_i \quad \forall 1 \leq i \leq n$

Let $A = U \Sigma V^*$ be the SVD of A .

The eigen values of $A^* A$ are σ_i^2 . The right singular vectors are the corresponding orthonormal vectors of $A^* A$

i.e. $A^* A v_i = \sigma_i^2 v_i \quad \forall 1 \leq i \leq n$

$$A = U \Sigma V^*$$

$$(A^* A) = (U \Sigma V^*)^* (U \Sigma V^*) = V \Sigma^* U^* U \Sigma V^* = V \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix} V^*$$

eigen decomposition of $A^* A$

Hermitian.

The eigenvalues of $A A^*$ are σ_i^2 & ($m-n$) zeroes. The left singular vectors are the orthonormal eigenvalues of $A A^*$.

$$\begin{bmatrix} A \\ \vdots \\ \end{bmatrix}_{n \times m} \quad \begin{bmatrix} A^* A \\ \vdots \\ \end{bmatrix}_{n \times n} \quad \begin{bmatrix} A A^* \\ \vdots \\ \end{bmatrix}_{m \times m}$$

$$(A A^*) u_i = \sigma_i^2 u_i \quad \forall 1 \leq i \leq n.$$

Proof:

$$\begin{aligned} AA^* &= (U\Sigma V^*)(V\Sigma V^*) \\ &= U\Sigma V^* (V\Sigma^* V^*) \\ &= U\Sigma^2 V^* \end{aligned}$$

← eigen decomposition of

$$AA^* = \begin{matrix} \text{m-n orthonormal} \\ \text{vectors} \end{matrix} \left[\begin{matrix} U_1 & \cdots & U_n \end{matrix} \right] \begin{matrix} \Sigma^2 \\ \vdots \\ 0 \end{matrix} \left[\begin{matrix} V_1^* \\ \cdots \\ V_m^* \end{matrix} \right]$$

(m-n) zeroes.

*

Solving a system of linear equations

- Gaussian Elimination.
- Method
- algorithm
- conditioning
- stability
- work (no. of computation) $\xrightarrow{\quad}$ floating pt. arithmetic.
- QR decomposition $\xrightarrow{\quad}$ projection matrices
 $\xrightarrow{\quad}$ Gram-Schmidt orthonormalization.
- SVD to solve system of linear eqns.
- least squares methods.
- Iterative methods.

* Gaussian Elimination

$$AX = B$$

$$A_{m \times n} X_{n \times 1} = B_{m \times 1}$$

linear system of equations

$$\begin{bmatrix} a_{11} & \dots & a_{1n} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{bmatrix}_{m \times n} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}_{m \times 1}$$

linear system of eqns

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

$$a_{m1}x_1 + \dots + a_{mn}x_n = b_m$$

m eqns in n variables.

Gaussian elimination

Let A be a $m \times m$ matrix.

Row reduction

$$3x + y + 5z = 12$$

$$x + y - z = 9$$

$$2x - 4y + z = 5$$

$$\left[\begin{array}{ccc|c} 3 & 1 & 5 & 12 \\ 1 & 1 & -1 & 9 \\ 2 & -4 & 1 & 5 \end{array} \right]$$

1. Augmented matrix

$$\left[\begin{array}{ccc|c} 3 & 1 & 5 & 12 \\ 1 & 1 & -1 & 9 \\ 2 & -4 & 1 & 5 \end{array} \right]$$

Row Reduced form

① Upper triangular matrix
(at least).

Row Reduced echelon form

① Upper triangular matrix
+ first non-zero element of row
should be 1!

$$R_3 \leftrightarrow R_2$$

$$\left[\begin{array}{ccc|c} 3 & 1 & 5 & 12 \\ 2 & 1 & -1 & 9 \\ 1 & 1 & 5 & -1 \end{array} \right]$$

$R_1 \leftrightarrow R_2$

$$\left[\begin{array}{ccc|c} 1 & 1 & -1 & 9 \\ 3 & 1 & 5 & 12 \\ 2 & -4 & 1 & 5 \end{array} \right]$$

E_{12} *

* **Permutation matrix** ($m \times m$)

$P_{ij} =$

$\underline{\underline{P_{12}}}$

$3 \times 3 \quad I_{3 \times 3} \rightarrow P_{12}$

$$\left[\begin{array}{ccc} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right] \rightarrow \left[\begin{array}{ccc} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{array} \right]$$

P_{12} : can be defined as exchange of rows 1 & 2 or exchange of columns 1 & 2.

P_{ij} = permutation matrix obtained by exchanging the i^{th} & j^{th} columns (1 rows) of $I_{m \times m}$.

Rows of 1st matrix \times Col's of right matrix

e.g.

$$\left[\begin{array}{ccc} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{array} \right] \left[\begin{array}{ccc} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right] = \left[\begin{array}{ccc} 4 & 5 & 6 \\ 1 & 2 & 3 \\ 7 & 8 & 9 \end{array} \right]$$

essentially exchanging
the 1st & 2nd rows.

$P_{ij} A =$

↓ Right left multiplication with P_{ij}
jth row is going to be exchanged with ith rows of A.

$$\left[\begin{array}{ccc} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{array} \right] \left[\begin{array}{ccc} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{array} \right] = \left[\begin{array}{ccc} 2 & 1 & 3 \\ 5 & 4 & 6 \\ 8 & 7 & 9 \end{array} \right]$$

Right multiplication with P_{ij}

↓ jth col is exchanged with ith col of A.

ΔP_{ij}

$$\left[\begin{array}{ccc|c} 3 & 1 & 5 & 12 \\ 1 & 1 & -1 & 9 \\ 2 & -4 & 1 & 5 \end{array} \right] \xrightarrow{\substack{R_1 \leftrightarrow R_2 \\ P_{12}[A:b]}} \left[\begin{array}{ccc|c} 1 & 1 & -1 & 12 \\ 3 & 1 & 5 & 9 \\ 2 & -4 & 1 & 5 \end{array} \right]$$

$$R_2 = R_2 - 3R_1 \quad R_3 = R_3 - 2R_1$$

$$\left[\begin{array}{ccc|c} 1 & 1 & -1 & 9 \\ 0 & -2 & -8 & -15 \\ 0 & -6 & 3 & -13 \end{array} \right] \xrightarrow{-\frac{R_2}{2}, -\frac{R_3}{3}}$$

$$\left[\begin{array}{ccc} 1 & 0 & 0 \\ -3 & 1 & 0 \\ -2 & 0 & 1 \end{array} \right] \left[\begin{array}{ccc} 1 & 1 & -1 \\ 3 & 1 & 5 \\ 2 & -4 & 1 \end{array} \right] \rightarrow \left[\begin{array}{ccc} 1 & 1 & -1 \\ 0 & -2 & 8 \\ 0 & -3 & 3 \end{array} \right]$$

$$R_2 - 3R_1 \wedge R_3 - 2R_1$$

$$A \xrightarrow{P_{12}} (P_{12} A) \xrightarrow{L_1} (L_1 P_{12} A) \xrightarrow{L_2} (L_2 L_1 P_{12} A)$$

$$L_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -3 & 1 \end{bmatrix}$$

upper
triangular
matrix

$$L_2 L_1 P_{12} A = U$$

a_{11} is non-zero. $\Rightarrow a_{11} \neq 0$.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

\Rightarrow L_i 's are always lower triangular

L_1

$$\begin{bmatrix} 1 & 0 & 0 \\ -\frac{a_{21}}{a_{11}} & 1 & 0 \\ -\frac{a_{31}}{a_{11}} & 1 & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22}' & a_{23}' \\ 0 & a_{32}' & a_{33}' \end{bmatrix}$$

L_2 is applied to $L_1 A$.

assuming
that
 a_{22}' is
non-zero.

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{a_{32}'}{a_{22}'} & 1 \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22}' & a_{23}' \\ 0 & a_{32}' & a_{33}' \end{bmatrix}$$

$$= \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ 0 & a_{22}' & a_{23}' \\ 0 & 0 & a_{33}'' \end{bmatrix}$$

* Exercise

4x4 matrix

See that L_1, L_2, L_3 obtained commute. with each other.

$$L = L_1 L_2 (A)$$

$$L_2 L_1 (A) = L$$

*

Gaussian Elimination without pivoting

Assumption the pivots encountered at every step of reduction are non-zero.

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ \vdots & & & \\ a_{n1} & a_{n2} & \dots & a_{nm} \\ & & & a_{nn} \end{bmatrix}$$

we can find lower Δ^r matrices L_1, L_2, \dots, L_{n-1} such that

$$L_{n-1} \left(\dots L_2 (L_1 A) \right) = U \text{ (upper } \Delta^r \text{)}$$

$$L_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ -\frac{a_{21}}{a_{11}} & 1 & & & & \\ \vdots & \vdots & & & & \\ -\frac{a_{n1}}{a_{11}} & 0 & \dots & 0 & 1 & \end{bmatrix}_{n \times n}$$

$$L_2 = \begin{bmatrix} 1 & 0 & & & & \\ 0 & 1 & & & & \\ \vdots & -\frac{a_{32}}{a_{22}} & 1 & & & \\ 0 & -\frac{a_{n2}}{a_{22}} & & 1 & & \end{bmatrix}_{n \times n}$$

$$L_K = \begin{bmatrix} 1 & & & & \\ 0 & 1 & & & \\ \vdots & & 1 & & \\ 0 & & & l_{K+1, K} & \\ 0 & 0 & l_{n, K} & \cdots & 1 \end{bmatrix}$$

$$l_{k+1, k} = -\frac{a_{k+1, k}^{(k-1)}}{a_{k, k}^{(k-1)}}$$

$$l_{k+1, k} = -\frac{a_{k+1, k}^{(k-1)}}{a_{k, k}^{(k-1)}}$$

$$l_{jk} = -\frac{a_{jk}^{(k-1)}}{a_{k, k}^{(k-1)}} \quad (j > k+1)$$

Let

$$L = L_n \cdot L_{n-1} \cdots L_2 \cdot L_1 \stackrel{\text{check}}{=} \begin{bmatrix} 1 & & & & \\ -\frac{a_{11}}{a_{11}} & 1 & & & \\ & \ddots & \ddots & & \\ & & -\frac{a_{22}}{a_{22}} & 1 & \\ & & & \ddots & \\ -\frac{a_{n1}}{a_{11}} & -\frac{a_{n2}}{a_{22}} & \cdots & -\frac{a_{nn}}{a_{nn}} & 1 \end{bmatrix}$$

$$LA = U$$

$$A = L^{-1}U$$

now,

let L be L^{-1} & L^{-1} be L .

$$L^{-1}A = U$$

$$A = LU$$

LU factorization of A .

* Solving a system of linear eqns.

$$Ax = b.$$

- factorize $A = LU$

$$Ax = b$$

$$(LU)x = b.$$

$$\text{i.e. } L(Ux) = b.$$

$$\downarrow \\ y$$

- Solve $Ly = b$. for y . (by substitution)

- Solve $Ux = y$. (by back substitution).

* Gaussian Elimination with pivoting

If at the k^{th} step of Gaussian Elimination, the element

$a_{k,k}^{(k-1)} = 0$, then you need "pivoting"

(pivot at the
 k^{th} step).

$$A^{k-1} = \left[\begin{array}{cccc} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22}^{(1)} & \cdots & a_{2n} \\ 0 & 0 & & \\ \hline k & \cdots & 0 & \\ & & * & \\ & & * & \\ & 0 & 0 & * \end{array} \right] \quad \{$$

Strategy - I (Partial)

look at the elements below the k^{th} row & take the absolute value exchange with the one that has largest abs. value. (Row exchange).

* Gaussian Elimination is always done with pivoting because it's an unstable algorithm in itself. (theoretically). (even with pivot).

* Partial Pivoting

If during G.E., the pivot at the k^{th} step is zero, then choose $a_{k,k}^{(k-1)}$ to be the largest among $|a_{jk}^{(k-1)}|$, $k \leq j \leq n$ (for stability)

$$L_n P_n (L_2 P_2 (L_1 (P_1 A))) = \Delta$$

Combining all the steps

$$(L_n' L_{n-1}' \dots L_1' P_n P_{n-1} \dots P_1) \Delta = \Delta$$

$$L^{-1} P A = \Delta$$

$$\boxed{P A = L \Delta}$$

* Complete pivoting

Pivots are chosen among elements $\{a_{ij}\}_{\substack{1 \leq k \leq n \\ j \leq k \leq n}}$

By looking at their absolute value.

$$L_n P_n \dots L_2 P_2 L_1 P_1 A Q_1 Q_2 \dots Q_n = \Delta$$

$$L^{-1} P A Q = \Delta$$

$$\boxed{P A Q = L \Delta}$$

*

Solving a system of eq's using G.E. with pivoting

$$Ax = b \text{, have } PAQ = LDU$$

$$PAx = Pb \rightarrow \text{column}$$

$$PAQx = PbQ, Q \text{ doesn't}$$

$$\begin{aligned} LDUx &= PbQ \quad \text{affect } x. \quad Ly = PbQ. \leftarrow \text{Solve for} \\ &\downarrow \\ y. & \quad Lx = y \quad \leftarrow \text{Solve for.} \end{aligned}$$

I] G.E. without pivoting

$$U = A, L = I.$$

col

→ for $k = 1$ to $m-1$:

row

→ for $j = k+1$ to m :

$$l_{jk} = u_{jk}$$

$$u_{kk},$$

for $i = k$ to m :

$$u_{j,i} = u_{j,i} - l_{jk} u_{ki}$$

II] G.E. with partial pivoting

$$U = A, L = I, P = I.$$

for $k = 1$ to $m-1$:

- select $i \geq k$ to maximize $|u_{ik}|$

- interchange rows u_i & u_k

for $j = 1$ to $k-1$
 $u_{kj} \leftrightarrow u_{ij}$

we only calculating floating point operations. (no. of)

2nd
(inner loop)

for $j = k+1$ to m :

$$l_{jk} = \frac{u_{jk}}{u_{kk}}$$

for $i = k$ to m :

$$u_{ji} = u_{ji} - l_{jk} u_{ki}$$

* Operation count

for each $k \leq i \leq m$,

the operation $u_{ji} = u_{ji} - l_{ji} u_{ii}$

takes 2 operations (1 multiplication
1 subtraction)

entries manipulated in

column $k = m-k = l$ (say).

division required in the k th column.

column = l .

1st column.

Total # of operations = $(m-1) + 2(m-1)$ for $(m-1)$ rows

2nd column

= $(m-2) + 2(m-2)$ operations
for $(m-2)$ rows.

$$\sum_{k=1}^{m-1} (m-k) + \sum_{k=1}^{m-1} 2(m-k)^2$$

:

$$\frac{5}{3} m^3 + \Theta(m^2)$$

If we use Gaussian elimination - $O(n^3)$
System of eqns ???

$\Rightarrow PAQ = LU$ Q is orthogonal and invertible
 $PA = LUQ^T$

$PAx = Pb$

$LUQ^T x = Pb$.

- Solve $Lz = Pb$ for z
- Solve $Uy = z$ for y
- Set $x = Q^T y$???

Notes

① $\det P = \pm 1$ ($\det Q$)

② $\det L = 1$

③ $\det U = \pm \det A$.

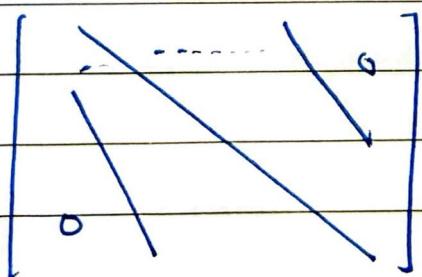
"product of
pivot

④ LU factorization of A (or PA or PAQ) exists irrespective of whether A is invertible or not.

⑤ The LU factorization of A exists if and only if
 $\Leftrightarrow \det \Delta_k \neq 0 \quad \forall 1 \leq k \leq n$,

where Δ_k = top left $k \times k$ submatrix of A .

* Band matrix



LU factorization preserves band structure.

i.e. both L & U of band matrices are band matrices too.

*

Cholesky factorization

Upsilonis

Let A be a symmetric 2×2 matrix

$$\begin{bmatrix} a & c \\ c & d \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ c/a & 1 \end{bmatrix} \begin{bmatrix} a & c \\ 0 & d - (\frac{c}{a})c \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 \\ c/a & 1 \end{bmatrix} \begin{bmatrix} a & 0 \\ 0 & d - (\frac{c}{a})c \end{bmatrix} \begin{bmatrix} 1 & c/a \\ 0 & 1 \end{bmatrix}$$

$L \quad D \quad L^T$

Theorem: If A is $\mathbb{R}^{n \times n}$ is symmetric and if
 Δ_k is non-singular $\forall 1 \leq k \leq n$, then \exists a unit
lower Δ^r matrix L and a diagonal matrix D s.t.

$$A = LDL^T$$

This factorisation is unique.

Proof:

$\det \Delta_k \neq 0 \quad \forall 1 \leq k \leq n$.

$\exists A$ has a LU factorization

let $A = LU$.

$$L^{-1} A L^{-T} = U L^{-T} = \text{Diagonal matrix.}$$

\therefore symmetric \therefore upper Δ^r matrix.

$$(L^{-1} A L^{-T})^T = L^{-T} A^T L^{-T}$$

$$L^{-T} A^T L^{-T} = D.$$

$$A = L D L^T$$

* Conditioning and Stability

* Conditioning of a problem

X : normed vector space of data.

Y : normed vector space of solⁿs.

$$f: X \rightarrow Y$$

Defⁿ: Absolute condition number

Let δx denote a small perturbation

$$\text{in } x \& \delta f = f(x + \delta x) - f(x)$$

Then the abs. condition number $\hat{K}(x)$ of f at x

$$\text{is defined as } \hat{K}(x) = \limsup_{\delta \rightarrow 0} \frac{\|\delta f\|}{\|\delta x\|} \begin{matrix} \leftarrow \text{abs. change} \\ \text{in soln} \end{matrix}$$

\rightarrow abs. change in data.

$$\approx \sup \frac{\|\delta f\|}{\|\delta x\|}$$

Defⁿ:

Relative condition number

$$K(x) := \sup_{\delta x} \frac{\|\delta f\| / \|f\|}{\|\delta x\| / \|x\|} \begin{matrix} \leftarrow \text{rel. change} \\ \text{in soln} \end{matrix}$$

\rightarrow rel. change
in data.

Ex. f is the problem of computing \sqrt{x} for $x > 0$.

$$f: x \mapsto \sqrt{x}.$$

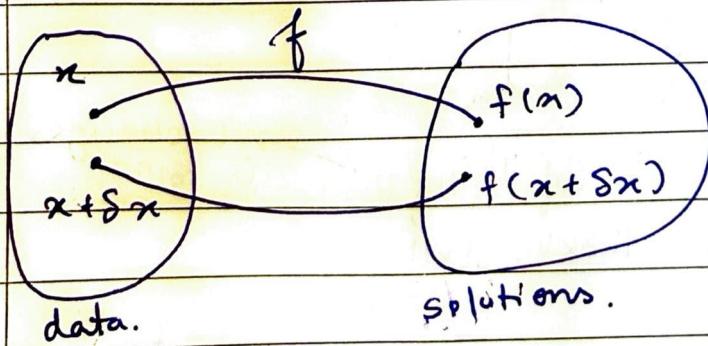
$$J(x) = \frac{1}{2\sqrt{x}}$$

$$K(x) \approx \kappa \quad K(x) = \frac{\|J(x)\|}{\|f'(x)\|/\|x\|} = \frac{\|y_2\sqrt{x}\|}{\|\sqrt{x}\|/\|x\|}$$

$= \frac{1}{2}$ (small conditioning no.
which says that this is
a well conditioned problem).

* Recall : conditioning \leftarrow abs. condition number κ
rel. condition number $K(x)$.

$$f: x \rightarrow \sqrt{x}, \quad K(x) = 1/2 \quad (f \text{ is well-conditional})$$



Ex: f: finding sq. rt. of quadratic polynomial.

$$\begin{array}{ll} \text{Input} & \text{Output} \\ (b, c) & \sqrt{b^2 - 4c} \\ & -\frac{b \pm \sqrt{b^2 - 4ac}}{2} \end{array}$$

$$k(x) \approx \frac{\|\delta f\| / \|f(x)\|}{\|\delta x\| / \|x\|}$$

$$k(x) \approx \sup_{\|x\| \rightarrow 0} \frac{\|\delta f\| / \|f(x)\|}{\|\delta(x)\| / \|x\|} = \frac{\|\delta f\| / \|\delta x\|}{\|\delta(x)\| / \|x\|}$$

if f is continuous

$$J_f(x) \\ \|\delta f(x)\| / \|x\|$$

$$J_f(x) = \left[\frac{-1 \pm \sqrt{b^2 - 4c}}{2} \right] \quad F \frac{\pm \sqrt{b^2 - 4c}}{\sqrt{b^2 - 4c}}$$

If we have a polynomial with repeated roots then there is small tolerance for perturbations in b and c .

- * If the polynomial has repeated roots $\|J(x)\|$ is ∞ in which case f is an ill-conditioned problem.

f: solving a system of eqns $\rightarrow Ax = b$.



given A, b

solve for x

given A, x

compute

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

b
(matrix-
vector)
multiplication

* Condition of matrix-vector multiplication

(for A)

2 problems

$$x \mapsto Ax (= b)$$

$$b \mapsto A^{-1}b (= x)$$

$$x \mapsto Ax \quad K(x) = \sup_{\delta x} \frac{\|A(x + \delta x) - Ax\|}{\|\delta x\|} / \|Ax\|$$

$$= \sup_{\delta x} \frac{\|A(\delta x)\|}{\|Ax\|} / \frac{\|\delta x\|}{\|x\|} \quad \because \sup_{\delta x} \frac{\|A\delta x\|}{\|\delta x\|} = \|A\|$$

$$= \|A\| \left(\frac{\|x\|}{\|Ax\|} \right)$$

$$\leq \|A\| \|A^{-1}\|$$

$$K(x) \leq \|A\| \|A^{-1}\| = K(A).$$

e.g.

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

depends on the
norm taken for
 $A \cdot A^{-1}$

* Condition of a system of eqns $Ax=b$ w.r.t. perturbations in A .

Theorem: Let b be fixed and consider the problem of solving $Ax=b$ w.r.t. perturbations in A .

Then the condition number of the problem is $K = K(A)$.

Proof: If A is perturbed by an infinitesimal amount δA then x must change to $x + \delta x$. So,

δA is a matrix.

$$(A + \delta A)(x + \delta x) = b.$$

$$\therefore \underbrace{Ax}_{b} + A\delta x + \delta Ax + \underbrace{\delta A\delta x}_{0} = b.$$

$$A\delta x + \delta Ax = 0.$$

$$\delta x = -A^{-1}\delta Ax.$$

$$\delta x = -A^{-1}\delta Ax.$$

$$\|\delta x\| \leq \|A^{-1}\delta Ax\|$$

$$\leq \|A^{-1}\| \|\delta A\| \cdot \|x\|$$

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \|\delta A\|$$

\Rightarrow rel. change in soln ($A \mapsto x$), b is fixed.

$$K = \sup_{\delta A} \frac{\|\delta x\| / \|x\|}{\|\delta A\| / \|A\|} \leq \sup_{\delta A} \frac{\|A^{-1}\| \|\delta A\|}{\|\delta A\| / \|A\|}$$

\downarrow rel. change in data.

$$\leq \|A^{-1}\| \|A\|$$

$$\leq K(A).$$

There is some A for which this bound is reached.

* It can be shown that δA and x can be chosen so that the bound is attained.

$$\therefore K = K(A)$$

* Properties of $K(A)$

- ① $K_p(A) = \frac{\|A\|}{\|A^{-1}\|}$ is norm dependent.
- ② For every matrix A $K(A) \geq 1$ ($\|A\| \|A^{-1}\| \geq \|A A^{-1}\| = 1 \|\|=1$)
- ③ $K(A) = K(A^{-1})$
- ④ $K(\alpha A) = K(A)$
- ⑤ If A is unitary or orthogonal
then $K_2(A) = 1$
(exercise) ??
- ⑥ $K_2(A)$ implies that $Ax=b$ where A is unitary or orthogonal
is a perfectly conditioned system.
→ most often 2 norm is used
- ⑦ By convention, $K(A) = \infty$ for a singular matrix A .

→ Does $\det A \approx 0$ indicate a large $K(A)$?
This is not necessarily true, for example.

e.g. (i) Let $B_n = \begin{bmatrix} 1 & -1 & -1 & \dots & -1 \\ & 1 & & & \\ & & 0 & \ddots & \\ & & & \ddots & 1 \end{bmatrix} \in \mathbb{R}^{n \times n}$

$$\det B_n = 1.$$

$$K(B_n) = n^{2^{n-1}}$$

(large,
badly
conditioned).

eg.
ii)

let $D_n = \text{diag}(10^{-1}, \dots, 10^{-1}) \in \mathbb{R}^{n \times n}$.

$$K_p(D_n) = 1 \quad \text{but} \quad \det D_n = 10^{-n} \quad (\approx 0)$$

Q. eg.

Effect of ~~size~~ → perturbations

$$\left(\begin{array}{cccc} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{array} \right) \left(\begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \end{array} \right) = \left(\begin{array}{c} 32 \\ 23 \\ 33 \\ 31 \end{array} \right) \quad \text{soln} \left(\begin{array}{c} 1 \\ 1 \\ 1 \\ 1 \end{array} \right)$$

Consider the perturbed data.

$$A(x + \delta x) = \left(\begin{array}{c} 32.1 \\ 22.9 \\ 33.1 \\ 30.9 \end{array} \right) \quad \text{soln} \left(\begin{array}{c} 9.2 \\ -12.6 \\ 4.5 \\ -1.1 \end{array} \right) \quad \text{cond} \quad \left(\begin{array}{c} 0.1 \\ -0.1 \\ 0.1 \\ -0.1 \end{array} \right)$$

Ex:

Consider perturbations in A

$$\left(\begin{array}{cccc} 10 & 7 & 8.1 & 7.2 \\ 7.08 & 5.04 & 6 & 5 \\ 8 & 5.98 & 9.89 & 9 \\ 6.99 & 4.99 & 9 & 9.98 \end{array} \right) \quad x = \left(\begin{array}{c} -81 \\ 137 \\ -34 \\ 22 \end{array} \right)$$



Compare

$$\frac{\| \delta A \| / \| A \|}{\| \delta x \| / \| x \|} \quad \times$$

$$\checkmark \quad \frac{\| \delta A \| / \| A \|}{\| \delta x \| / \| x \|}$$

Example in Last class

$$\text{Given } A = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}$$

conditioning of $Ax=b$.

- w.r.t perturbations in b

$$\kappa \approx \frac{\|Ax\|/\|x\|}{\|Ab\|/\|b\|} \approx \frac{8}{4200} \approx 1600$$

- w.r.t. perturbations in A

$$\kappa \approx \frac{\|Ax\|/\|x\|}{\|A\|/\|A\|} \approx \frac{160}{0.22/33} \approx \frac{160 \times 33}{0.22}$$

$$\delta A = \begin{pmatrix} 0 & 0 & 0.1 & 0.2 \\ 0.08 & 0.09 & 0 & 0 \\ 0 & -0.02 & -0.11 & 0 \\ -0.01 & -0.01 & 0 & -0.02 \end{pmatrix}$$

$$\boxed{\kappa_2(A) \approx 2984}$$

* Floating point representation

1985 - IEEE standard for floating point representation.

Kahan
& Stone --

754-2008
↓

754-2019.

$$\underbrace{}_{d_1} \underbrace{}_{d_2} \times 10^c \quad \boxed{\text{Base 10}}$$

e.g. -99.5 < C < 99.

$$\underbrace{}_{d_0} \cdot \underbrace{}_{d_1} \underbrace{}_{d_2} \dots \underbrace{}_{d_{p-1}} \times \beta^t$$

each $[0 \leq d_i \leq \beta - 1]$

$\underbrace{}_{p\text{-spaces}}$

$d_i \in \{0, 1\}$
if $\beta = 2$.

$$\rightarrow \beta = 2, p = 3, -1 \leq t \leq 2.$$

do $\neq 0$. | Normalized representation.

$$\underbrace{}_{d_0} \cdot \underbrace{}_{d_1} \underbrace{}_{d_2} \times 10^{-1}$$

$$\begin{array}{l} \underbrace{1}_{d_0} \cdot \underbrace{0}_{d_1} \underbrace{0}_{d_2} \times 2^0 \longleftrightarrow 1. \\ \underbrace{0}_{d_0} \cdot \underbrace{1}_{d_1} \underbrace{0}_{d_2} \times 2^1 \longleftrightarrow 1 \end{array} \quad \left. \begin{array}{l} \text{ambiguous} \\ \therefore d_0 \neq 0 \end{array} \right\}$$

* Construct a floating point system $\beta=2, p=3, -1 \leq t \leq 2$

$F \subseteq \mathbb{R}$

TF

TR

1.00×2^{-1}	\leftrightarrow	orthogonal ≈ 0.5
1.01×2^{-1}	\leftrightarrow	0.625
1.10×2^{-1}	\leftrightarrow	0.75
1.11×2^{-1}	\leftrightarrow	0.875

$\downarrow \times 2$

$$1.00 \times 2^0 \leftrightarrow 1$$

$$1.01 \times 2^0 \leftrightarrow 1.125$$

$$1.10 \times 2^0 \leftrightarrow 1.25$$

$$1.11 \times 2^0 \leftrightarrow 1.375$$

$\downarrow n$

$$1.00 \times 2^1 \leftrightarrow 2$$

$$1.01 \times 2^1 \leftrightarrow 2.125$$

$$1.10 \times 2^1 \leftrightarrow 2.25$$

$$1.11 \times 2^1 \leftrightarrow 2.375$$

$\downarrow y_2$

$$1.00 \times 2^2 \leftrightarrow 4$$

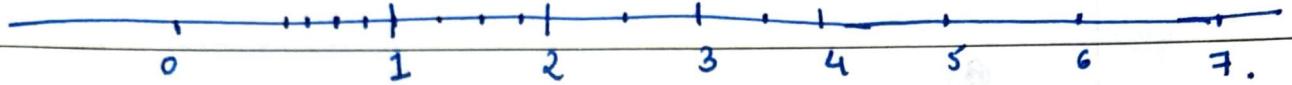
$$1.01 \times 2^2 \leftrightarrow 4.125$$

$$1.10 \times 2^2 \leftrightarrow 4.25$$

$$1.11 \times 2^2 \leftrightarrow 4.375$$

Non-uniform error

betw interval and outside interval



$$TF = \{0\} \cup \{d_0 \cdot d_1 d_2 \dots d_{p-1} \times \beta^t \mid \begin{array}{l} e_{\min} \leq t \leq e_{\max}, \\ 0 \leq d_i \leq \beta - 1, \\ d_0 \neq 0 \end{array}\}$$

density of representation is highest betw 0 & 1.

and as we move away to bigger no.s the representation becomes sparse.

x = actual real no. to be represented

x' = floating point representation of x .

$$\text{relative error representation} = \frac{|x - x'|}{|x|}$$

- Note that the floating pt. nos are not equally spaced
- Convention is 0 is represented by $[1.0 \times \beta^{e_{\min}-1}]$
- Rel representation when $x \in \mathbb{R}$ is represented by $x' \in TF$

$$\text{is defined as } \frac{|x - x'|}{|x|}$$

Q. What is the max rel. repr. error that can occur in a given floating pt. system?

eg. $x = 0.14159$ is represented by

$x = 3.14159$ is represented by 3.14×10^0

$$\text{then rel. error is represented by } = \frac{(3.14159 - 3.14)}{3.14159}$$

$$\approx \underline{\underline{0.0005}}$$

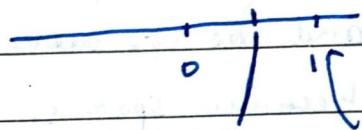
$$\text{eg. } \beta = 10, p = 5. \quad e_{\min} = -1$$

Q2.

$$9.\underline{9}\ \underline{9}\ \underline{9}\ \underline{9}\ \underline{9} \times 10^{-1}$$

$$\frac{9}{10} + \frac{9}{100} + \frac{9}{1000} + \frac{9}{10000} + \frac{9}{100000}$$

$$0.99999$$



$$I = 1.0000 \times 10^0$$

$$[0.99999] \quad 1.0001$$

$$1.0001 \times 10^0$$

$$1.0001$$

$$\boxed{1 + \beta^{1-p}}$$

↓
no. of bits
1.

$$\text{max representation error} = \left| \frac{1/000}{1.00005} \right| \left| \frac{1.00005 - 1}{1.00005} \right|$$

$$\text{now, } \frac{1 + \beta^{1-p}}{2} \rightsquigarrow 1$$

$$\text{max. representation error} = \left| \frac{1 + \beta^{1-p} - 1}{2} \right|$$

$$\left| \frac{1 + \beta^{1-p}}{2} \right|$$

Defⁿ:

The floating point representation of a real no. is of the form

$$\xrightarrow{\text{significand}} \underbrace{d.ddd}_{p\text{-digits}} \times \beta^e \xrightarrow{\text{base.}} \text{exponent}$$

P-precision

more precisely: the number represented by

$$\pm d.ddd \times \beta^e \text{ is } \boxed{\left(d_0 + d_1 \beta^{-1} + \dots + d_{p-1} \beta^{-(p-1)} \right) \times \beta^e}$$

each d_i is betⁿ $0 \leq d_i \leq \beta - 1$

- Let $e_{\min} \leq e \leq e_{\max}$, then # possible exponents = $e_{\max} - e_{\min} + 1$.
- # possible significands β^p .
- In general the floating pt. repr. of a real no. may not be unique. To calculate this, we consider "normalized" representation which $d_0 \neq 0$.

* Ulps (Units in the last place)

worst
floating
pt. no. / actual
number

x' x

$$\text{error} = |x - x'| \cdot \beta^{p-1} \text{ ulps (units in last place)}$$

If the floating pt. number $d.d\dots d \times \beta^e$ is used for rounding off a real no. z ,

$$\text{then the error is } \left| \frac{d.d\dots - z}{\beta^e} \right| \times \beta^{p-1} \text{ ulps.}$$

eq. $\beta = 10, p = 3$

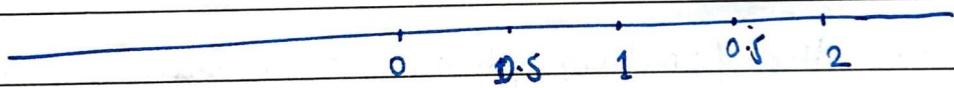
(i) Suppose $z = 0.0314 + 3.12 \times 10^{-2}$ is its fl. pt. representation.

$$\begin{aligned} \text{error} &= \left| \frac{3.12 - 0.0314}{10^{-2}} \right| \times 10^2 \\ &= 0.02 \times 10^2 \\ &= \underline{\underline{\alpha \text{ ulps.}}} \end{aligned}$$

(ii) $z = 0.0314159$

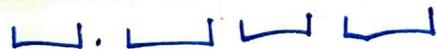
$$d.ddd - d \times \beta^e = 3.14 \times 10^{-2}$$

$$\begin{aligned} \text{error} &= \left| \frac{3.14 - 0.0314159}{10^{-2}} \right| \times 10^2 \\ &= 0.00159 \times 10^2 \\ &= 0.159 \text{ ulps.} \end{aligned}$$



for base β & precision p what is floating pt. no. that comes after 1.

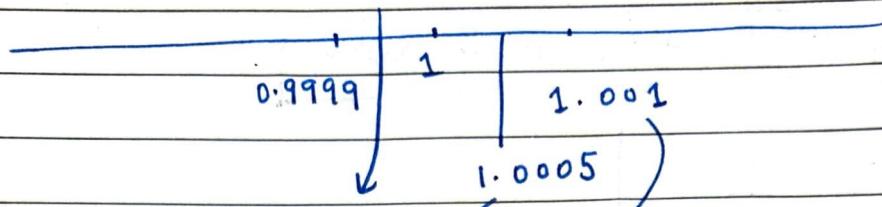
Base $\rightarrow 10$ precision $\rightarrow 4$



$$9.999 \times 10^{-1}$$

0.9999 → last real no.

last floating point number before 1 is $\underline{\underline{9.999 \times 10^{-1}}}$
first floating point number after 1 is $\underline{\underline{1.001 \times 10^0}}$



$$1 + \beta^{1-p} \times 10^e$$

$\frac{1 + \beta^{1-p} \times 10^e}{2}$

$$[1, 1 + \beta^{1-p}]$$

half half of this interval is $\boxed{\frac{1 + \beta^{1-p}}{2}}$

→ For base β , the possible abs. error

$|x - x'|$ can be as large as ???

$$\underbrace{0.0 \dots 0}_{p} \beta^1 \times \beta^e$$

where $\beta' = \beta/2$

this error is $\left(\frac{\beta}{2}\right) \beta^{-p} \times \beta^e$

$$\text{Relative error} = \frac{|x - x'|}{|x|}$$

$$\beta^e \leq x \leq \beta \times \beta^e$$

for a fixed e ,

$$10^e \leq 9.999 \times 10^e < 10 \times 10^e$$

$$\beta^e \leq |x| \leq \beta \cdot \beta^e$$

$$\frac{1}{\beta^e} \geq \frac{1}{|x|} \geq \frac{1}{\beta \cdot \beta^e}$$

$$\therefore \frac{|x - x'|}{|x|} \leq \frac{(\beta/2) \beta^{-p} \times \beta^e}{\beta^e} \leq \left(\frac{\beta}{2}\right) \beta^{-p}$$

* Rounding to 0 gives 100% error

$$\epsilon_{\text{mach}} := \frac{\beta^{1-p}}{2}$$

for a given machine with base β & precision p
the relative error is always bounded by ϵ_{mach} .

The max^m rel. repr. error in a floating pt. system with
base β & precision p is $\frac{\beta^{1-p}}{2} =: \epsilon_{\text{mach}}$

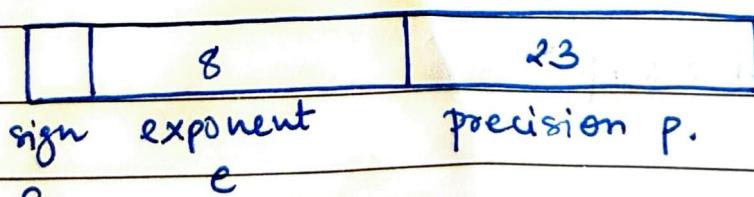
$$\underbrace{a + b}_{\substack{\downarrow \\ \text{rounded}}} = f(a+b)$$

* Last time: max rel. representation error is $\frac{1}{2} \beta^{1-p} =: \epsilon_{\text{mach}}$

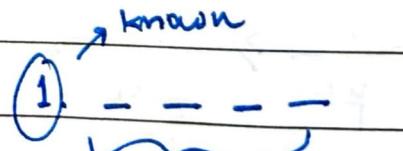
- IEEE single precision

$$-126 \leq e \leq 127$$

$\beta = 2$, $p = 23$, available length = 32 bits



$$(-1)^s \cdot 2^{e-127} \cdot (1+f)$$

① 
p places are allowed.

$$\epsilon_{\text{mach}} = \frac{1}{2} \beta^{1-p} = \frac{1}{2} \beta^{-24}$$

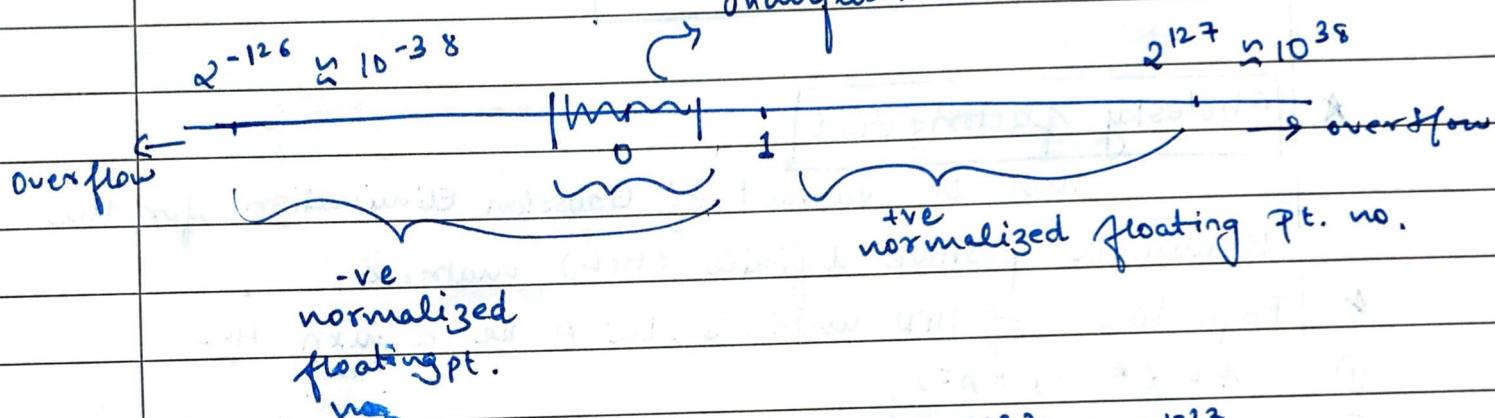
- IEEE double precision $B=2$, $p=52$, 64 bits length.

	11	52
sign	exponent	precision fraction/mantissa

$$e_{\min} = -1022 \quad e_{\max} = 1023.$$

$$\epsilon_{\text{mach}} = 2^{-53} \approx 10^{-16}$$

underflow/subnormal.



Range of normalized nos is 2^{-1022} to 2^{1023}
 $(\approx 10^{-308}$ to $10^{308}).$

① Fundamental property of floating pt. representation

$$\left| \frac{f(x) - x}{x} \right| \leq \epsilon_{\text{mach}}$$

for every $\forall x \in \mathbb{R}$, $\exists \epsilon$ with $|\epsilon| \leq \epsilon_{\text{mach}}$
such that $f(x) = x(1+\epsilon)$

$$\left| \frac{f(x) - x}{x} \right| = |\epsilon| \leq \epsilon_{\text{mach}}$$

② Let $x \star y$ denote the computed answer of the arithmetic operation \star (addition / subtraction / multiplication / division)

Fundamental property of floating pt. arithmetic -
 $\forall x, y \in \mathbb{R}, \exists \varepsilon \text{ with } |\varepsilon| \leq \varepsilon_{\text{mach}} \text{ s.t.}$

$$x \star y = xy(1 + \varepsilon)$$

* Cholesky factorization

This is variant of Gaussian Elimination for the Hermitian positive definite (HPD) matrices.

* Properties of HPD matrices - let A be a $m \times n$ HPD

$$\textcircled{1} \quad A = A^* \quad (A = A^T)$$

$$\textcircled{2} \quad x^* A y = y^* A x \quad (x^T A y = y^T A x) \quad \forall \text{ vectors } x, y.$$

$$\textcircled{3} \quad x^* A x \text{ is real} \quad \forall x \in \mathbb{C}^m \quad \text{moreover}$$

$$\text{PD} \Leftrightarrow x^* A x > 0 \quad \forall x \in \mathbb{C}^m$$

④ If X is a full rank $m \times n$ matrix

then $X^* A X$ is also HPD

$$\therefore (X^* A X)^* = X^* A X^{**} = X^* A X$$

$$\left\{ \begin{array}{l} \therefore x^* X^* A X x = (X x)^* A (X x) > 0 \\ \quad (X \text{-full rank} \Rightarrow X x \neq 0). \end{array} \right.$$

(5) As a consequence of (4)

① every principal submatrix of A is HPD

(ii) letting $x_i = \begin{bmatrix} 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix}$, $x_i^* A x_i = a_{ii} > 0$

\therefore all diagonal elements of A must be non-zero.

$$X = \begin{bmatrix} 1 & 0 & & \\ 0 & 1 & & \\ \vdots & & \ddots & \\ 0 & 0 & & 1 \end{bmatrix}_{n \times n} \quad \text{then } X^* A X = \text{top-left } k \times k \text{ submatrix.}$$

* Idea of Cholesky factorization

suppose $A = \begin{bmatrix} a_{11} & w^* \\ w & K \end{bmatrix} \xrightarrow{\text{1st row.}}$, note $a_{11} > 0$.

$$A \neq \begin{bmatrix} \sqrt{a_{11}} & 0 \\ 0 & R_1^* \end{bmatrix} \quad A \neq \text{f}$$

$$A = \begin{bmatrix} \sqrt{a_{11}} & 0 \\ w/\sqrt{a_{11}} & I \end{bmatrix} \begin{bmatrix} 1 & \cdots & 0 \\ \vdots & & \frac{K-w^*}{\sqrt{a_{11}}} \end{bmatrix} \begin{bmatrix} \sqrt{a_{11}} & w^*/\sqrt{a_{11}} \\ 0 & I \end{bmatrix}$$

$$A = R_1^* A_1 R_1$$

$$\begin{aligned}
 A &= R_1^* A_1 R_1 \\
 &= R_1^* (R_2^* (A_2) R_2) R_1 \\
 &= (R_1^* \dots R_m^*) A_m (R_m \dots R_1)
 \end{aligned}$$

$A = R^* R$
 ↓ ↓
 lower Upper triangular
 triangular.

* **Algorithm** $\Rightarrow a_{ii} \neq 0 \forall i$

eg.

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} = \begin{bmatrix} \bar{r}_{11} & 0 & 0 \\ \bar{r}_{12} & \bar{r}_{22} & 0 \\ \bar{r}_{13} & \bar{r}_{23} & \bar{r}_{33} \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ 0 & r_{22} & r_{23} \\ 0 & 0 & r_{33} \end{bmatrix}$$

$$\begin{bmatrix} \bar{r}_{11}^2 & \bar{r}_{11}\bar{r}_{12} & \bar{r}_{11}r_{13} \\ \bar{r}_{12}\bar{r}_{11} & \bar{r}_{12}^2 + \bar{r}_{22}^2 & \bar{r}_{12}r_{13} + \bar{r}_{22}r_{23} \\ \bar{r}_{13}\bar{r}_{11} & \bar{r}_{13}\bar{r}_{12} + \bar{r}_{23}\bar{r}_{22} & \bar{r}_{13}^2 + \bar{r}_{23}^2 + \bar{r}_{33}^2 \end{bmatrix}$$

$$r_{11} = \sqrt{a_{11}} \quad \& \text{ equality of matrices.}$$

Algorithm

for $k = 1$ to m

$$r_{kk} = \left(a_{kk} - \sum_{j=1}^{k-1} r_{kj}^2 \right) y_2$$

$$\frac{n(n-1)}{6}(2n-1)$$

for $i = k+1$ to m

$$r_{ik} = \left(a_{ik} - \sum_{j=1}^{k-1} r_{kj} a_{ij} \right) / r_{kk}$$

$$\frac{n^3}{3} - \frac{2n^2}{6} - \frac{n^2}{6}$$

Check operation count $\geq \frac{n^3}{3}$ flops.

- * Cholesky factorization is not same as LU factorization

Stability of an algorithm

The meaning of $\Theta(E_{\text{mach}})$.

the notation $\Psi(t) = \Theta(\Psi(t))$ means that \exists a positive constant C such that $\forall t$ close to some understood limit ($t \rightarrow 0$ or $t \rightarrow \infty$)

$$|\Psi(t)| \leq C |\Psi(t)| \quad (\text{Note that } c \text{ works uniformly for all } t).$$

Correct algorithm

Cholesky Algorithm

* for $k = 1$ to m :

$$r_{kk} = (a_{kk} - \sum_{j=1}^{k-1} r_{jk}^2)^{1/2}$$

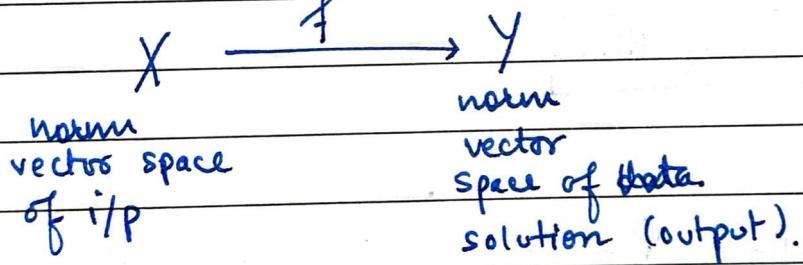
for $i = k+1$ to m :

$$r_{ik} = (a_{ik} - \sum_{j=1}^{k-1} r_{jk} r_{ji}) / r_{kk}$$

$\approx O(m^3/3)$

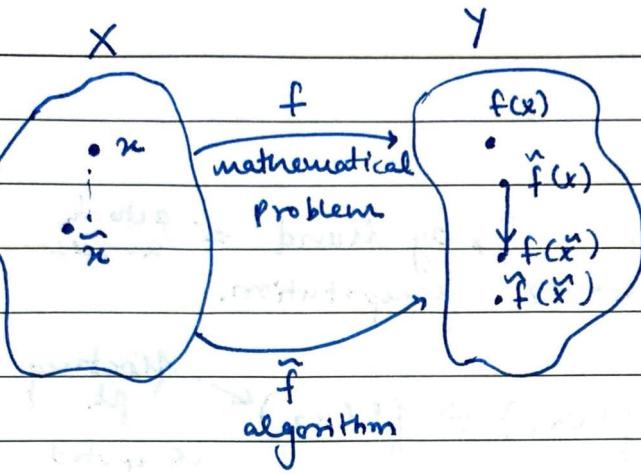
Stability of an algorithm

\hat{f} algorithm



Defn: An algorithm \hat{f} is said to be accurate if $\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|} = 0$

$$\boxed{\frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|} = \Theta(\epsilon_{mach})}$$



Defⁿ:

An algorithm \hat{f} is said to be (forward) stable if $\forall x \in X$,

$$\frac{\|\hat{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = O(\epsilon_{mach}), \text{ whenever}$$

$$\frac{\|x - \tilde{x}\|}{\|x\|} = O(\epsilon_{mach})$$

Nearly correct answer to a nearly correct question.

Defⁿ:

An algorithm \hat{f} is said to be backward stable if $\forall x \in X \quad \hat{f}(x) = f(\tilde{x})$ whenever $\frac{\|x - \tilde{x}\|}{\|x\|} = O(\epsilon_{mach})$.
for some \tilde{x} with

(Right answer to the nearly right question).

Remarks

- 1) Backward stability \Rightarrow forward stability

- 2) In general, backward stability analysis is easier to carry out.

Ex:

Subtraction

$$f(x_1, x_2) = x_1 - x_2 \xrightarrow{\text{by hand computation.}} \leftarrow \begin{array}{l} \text{actual} \\ \text{answer} \end{array}$$

$$\text{or} \\ f(x_1, x_2) = fl(x_1) \ominus fl(x_2) \xrightarrow{\substack{\text{floating pt.} \\ \text{computed answer.}}} \leftarrow$$

$fl(x_1, x_2)$

① By the fundamental ppty of fl. pt. representation.

$$fl(x_1) = x_1(1 + \varepsilon_1) \quad |\varepsilon_1| \leq \varepsilon_{\text{mach}}$$

$$fl(x_2) = x_2(1 + \varepsilon_2) \quad |\varepsilon_2| \leq \varepsilon_{\text{mach.}}$$

② By fundamental ppty of fl. pt. arithmetic

$$fl(x_1) \ominus fl(x_2) = (fl(x_1) - fl(x_2))(1 + \varepsilon_3) \quad \text{where } |\varepsilon_3| \leq \varepsilon_{\text{mach.}}$$

$$\begin{aligned} & f(x_1, x_2) \quad fl(x_1) \ominus fl(x_2) = x_1(1 + \varepsilon_1) - x_2(1 + \varepsilon_2) \\ & = x_1 - x_2 + (x_1 - x_2)(\varepsilon_1) \\ & = x_1 - x_2 + x_1\varepsilon_1 - x_2\varepsilon_2 \end{aligned}$$

$$\begin{aligned} & f(x_1, x_2) \Rightarrow f(x_1, x_2) = \tilde{x}_1 - \tilde{x}_2 \\ & = x_1(1 + \varepsilon_1) - x_2(1 + \varepsilon_2) \\ & = x_1 + x_1\varepsilon_1 - x_2 - x_2\varepsilon_2 \\ & = x_1 - x_2 + x_1\varepsilon_1 - x_2\varepsilon_2 \end{aligned}$$

$$f(x_1) \ominus f(x_2) = (x_1(1+\epsilon_1) - x_2(1+\epsilon_2))(1+\epsilon_3)$$

$$= x_1(1+\epsilon_1)(1+\epsilon_3) - x_2(1+\epsilon_2)(1+\epsilon_3)$$

$$\tilde{f}(x_1, x_2) = 1 + \epsilon_1\epsilon_3 + \epsilon_2\epsilon_3 + \epsilon_3^2$$

$$\leq 2\epsilon_{\text{mach}} + \epsilon_{\text{mach}}^2$$

$$\tilde{x}_1 = x_1(1+\epsilon_4) = x_1(1+\epsilon_4) - x_2(1+\epsilon_5) \quad \text{where } |\epsilon_4|, |\epsilon_5|$$

$$\tilde{x}_2 = x_2(1+\epsilon_5) = \tilde{x}_1 - \tilde{x}_2 = \tilde{f}(\tilde{x}_1, \tilde{x}_2) \leq 2\epsilon_{\text{mach}} + \epsilon_{\text{mach}}^2$$

i. Subtraction

is backward
stable.

$$\|x - \tilde{x}\| = O(\epsilon_{\text{mach}}) \quad \|x\|$$

② Outer product : stable but not backward stable.

$$f(x, y) = xy^* = A$$

$$(x + \delta x)(y + \delta y)^* = A + \delta A$$

Rank 1?
??

* Accuracy of a backward stable algorithm

Theorem Suppose a backward stable algorithm is applied on a computer (satisfying fundamental properties of floating pt. representation & floating pt. arithmetic) to solve a problem $f: x \rightarrow y$ with condition number K .

$$\text{Then } \frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} = O(K\epsilon_{\text{mach}})$$

functional stable algorithm.

$$\frac{\|f(x) - f(\tilde{x})\|}{\|f(x)\|} = O(K \epsilon_{mach})$$

Proof: defⁿ of backward stability : $\hat{f}(x) = f(\tilde{x})$ for some \tilde{x} s.t $\|\tilde{x} - x\| = O(\epsilon_{mach})$

$$K(x) = \sup_{\delta x \rightarrow 0} \left(\frac{\|\delta f\| / \|f\|}{\|\delta x\| / \|x\|} \right)$$

$$\Rightarrow K(x) \geq \frac{\|\hat{f}(x) - f(x)\| / \|f(x)\|}{\|x - \tilde{x}\| / \|x\|}$$

choose

$$\delta x = x - \tilde{x}$$

$$K(x) \geq \frac{\|\hat{f}(x) - f(x)\| / \|f(x)\|}{\|x - \tilde{x}\| / \|x\|}$$

$\approx O(K)$

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

$K(x)$

$$\Theta(K \epsilon_{mach}) = \frac{\|\hat{f}(x) - f(x)\|}{\|f(x)\|}$$

Thus, if an algorithm is backward stable, then the accuracy depends on the conditioning of the problem.

* Error Analysis for G.E. (See page 47-48).
algorithm for (GEPP).

$$|E_{jk}| = \left| \sum_{i=1}^j l_{ji} u_{ik} \cdot s_i \right| \leq \sum_{i=1}^j \|l_{ji}\| \|u_{ik}\| \cdot \underbrace{n\varepsilon}_{s_j \leq n\varepsilon} = n\varepsilon (\|L\| \cdot \|U\|)_{jk}$$

$$|E_{jk}| = n\varepsilon (\|L\| \cdot \|U\|)_{jk}$$

$$A = LU + E \quad \|E\| \leq n\varepsilon \|L\| \|U\| \quad \text{abs.}$$

for norms that do not depend on "signs of the matrix entries" (e.g. 1-norm, ∞ -norm, Frobenius norm, but not the 2-norm)

we may simplify this to -

$$\|E\| \leq n\varepsilon \|L\| \|U\| \quad \star\star$$

Remark: The above bound doesn't help much for larger n .

So in practice, error analysis is mostly computed on computed residual $r = Ax - b$. ($\tilde{x} = A^{-1}r + Ab$)

$$\therefore \|Sx\| = \|A^{-1}\| \|r\|$$

If GEPP is used $\|L\| \leq 1$

*

Projection matrices

A projection matrix is a matrix P which satisfies $P^2 = P$ (also called idempotent)

(In general, a linear map $T: V \rightarrow V$ is a projection if $T^2 = T$)

Note : ① If $v \in \text{range } P$, then $P(v) = v$

If $x = P(v)$,

$$P(x) = P(P(v)) = P^2(v) = P(v) = x$$

$$\therefore v = x$$

② If $v \notin \text{range } P$, then $P(v) - v \in \text{null}(P)$

$$\begin{aligned} P(P(v) - v) &= P^2(v) - P(v) \\ &= P(v) - P(v) = 0 \end{aligned}$$

③ If P is a projection then, $I - P$ is also a projection
Want $(I - P)^2 = (I - P)$

$$\begin{aligned} (I - P)^2 &= (I - P)(I - P)v \\ &\Rightarrow I(I - P)v - P(I - P)v \\ &= (I - P)(v - Pv) \\ &= v - Pv - Pv + P^2v \\ &= (I - P)v. \end{aligned}$$

④ $\text{range}(I - P) = \text{null}(P)$
(check.)

$$5) \text{ null } (I-P) = \text{range } (P)$$

(let $P = I - (I-P)$)

$$④ \text{ range } (I-P) = \text{null } (P)$$

(let $v \in (I-P)$)

$$\therefore (I-P)v = v.$$

$v - Pv = v.$ $v \neq 0.$

$$Pv = 0. \quad \therefore v \in \text{null space of } P.$$

$$v \notin (I-P)$$

$$\therefore (I-P)v = 0.$$

$$v - Pv = 0$$

$v = Pv.$ $\therefore v \in P$ v belongs to null space of $(I-P).$

$$⑥ \text{ range } (P) \cap \text{null space } (P) = \{0\}$$

The projection P separates the underlying vector space into two (disjoint) complementary subspaces

$$\boxed{\star \text{ V} = \text{range } P \oplus \text{null } P \star}$$

Note that : range (P) may not be orthogonal to the null space of P i.e. $\text{null } (P).$

Def": An orthogonal projection is a matrix P for which $\text{range } P \perp \text{null } P.$

P is an orthogonal projection if $\Leftrightarrow PAA^T = A$