

# Chapter 3. Direct Methods for Linear Systems.

## Section 1. Gauss Elimination Method.

$$\begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 3 \\ 3 \end{pmatrix} \rightarrow (2) - (1) \times \frac{1}{2}$$

$$\begin{pmatrix} 2 & 1 \\ 0 & \frac{3}{2} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 3 \\ \frac{3}{2} \end{pmatrix} \rightarrow \frac{3}{2}y = \frac{3}{2} \Rightarrow y = 1$$

$$2x + y = 3$$

$$\Rightarrow 2x = 3 - y = 2 \Rightarrow x = 1$$

$$A = (a_{ij})_{n \times n} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix} \in \mathbb{R}^n$$

Linear System:  $Ax = b$  (1).

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$$

Assume  $A$  is invertible

Idea: eliminate the entries in the lower part of the matrix.

For  $k = 1, 2, \dots, n-1$  do (as follows)

For  $i = k+1, k+2, \dots, n$  do  $\leftarrow$  let  $l_{ik} = a_{ik}/a_{kk}$

For  $j = k, k+1, \dots, n$  do  $b_i^{(k+1)} \leftarrow b_i^{(k)} - l_{ik}b_k^{(k)}$

$$a_{ij}^{(k+1)} \leftarrow a_{ij}^{(k)} - l_{ik} a_{kj}^{(k)}$$

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ & a_{22}^{(2)} & \dots & a_{2n}^{(2)} \\ & & a_{33}^{(3)} & \dots & a_{3n}^{(3)} \\ & & & \ddots & \ddots \\ & & & & a_{nn}^{(n)} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2^{(2)} \\ b_3^{(3)} \\ \vdots \\ b_n^{(n)} \end{pmatrix} \quad (2)$$

Assumption: The diagonal entries are always non-zero in the elimination process.  
 $a_{kk}^{(k)} \neq 0$ .

$$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad (X)$$

The upper right triangular system can be solved by the backward substitution method.

For  $k = n, n-1, \dots, 1$ , do

$$a_{kk}^{(k)} x_k = b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j$$

$$x_k = (b_k^{(k)} - \sum_{j=k+1}^n a_{kj}^{(k)} x_j) / a_{kk}^{(k)} \quad (*).$$

algorithm complexity.

$$O(n^3) + O(n^2).$$

## section 2. LU-decomposition Method.

Decompose  $A$  into the product of a lower triangular and an upper triangular matrix.

Denote the lower and upper triangular matrix by  $L$  and  $U$ , respectively.

$$A = LU \quad \text{with} \quad L = \begin{pmatrix} 1 & & & & \\ l_{21} & 1 & & & \\ l_{31} & l_{32} & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ l_{n1} & l_{n2} & \vdots & \ddots & 1 \end{pmatrix}$$

$$U = \begin{pmatrix} u_{11} & u_{12} & \cdots & \cdots & u_{1n} \\ & u_{22} & \cdots & \cdots & u_{2n} \\ & & \ddots & \ddots & \vdots \\ & & & \ddots & \vdots \\ & & & & u_{nn} \end{pmatrix}$$

$l_{ik}$

$$Ax = b \Rightarrow LUx = b.$$

$$\text{Let } Ux = y. \quad \Rightarrow \quad Ly = b.$$

step 1. solve  $Ly = b$ . by forward substitution.

$$\begin{pmatrix} 1 & & & & \\ l_{21} & 1 & & & \\ l_{31} & l_{32} & 1 & & \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \vdots & \ddots & 1 \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} \quad \checkmark$$

$$O(n^2)$$

Step 2, solve  $Ux = y$  by backward substitution.

$$\begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ & u_{22} & \dots & u_{2n} \\ & & \ddots & \\ & & & u_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}.$$

$O(n^2)$

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} = \begin{pmatrix} 1 & & & \\ l_{21} & 1 & & \\ l_{31} & l_{32} & 1 & \\ \vdots & \vdots & \vdots & \ddots & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ & u_{22} & \dots & u_{2n} \\ & & \ddots & \\ & & & u_{nn} \end{pmatrix}$$

$$a_{11} = u_{11} \quad a_{12} = u_{12} \quad a_{1j} = u_{1j} \quad j = 1, 2, \dots, n$$

$$l_{i1} u_{11} = a_{i1} \quad i = 2, 3, \dots, n \Rightarrow u_{i1} = a_{i1} / u_{11}$$

For  $k = 1, 2, \dots, n$ .

$$u_{kj} = ? \quad a_{kj} = \sum_{m=1}^k l_{km} u_{mj} \quad j = k, k+1, \dots, n$$

$$= u_{kj} + \sum_{m=1}^{k-1} l_{km} u_{mj}$$

$$u_{kj} = a_{kj} - \sum_{m=1}^{k-1} l_{km} u_{mj} \quad j = k, k+1, \dots, n$$

$$l_{ik} = ? \quad a_{ik} = \sum_{m=1}^k l_{im} u_{mk}, \quad i = k+1, \dots, n$$

$$l_{ik} \cdot u_{kk} = a_{ik} - \sum_{m=1}^{k-1} l_{im} u_{mk}$$

$i = k+1, k+2, \dots, n$

$$\Rightarrow l_{ik} = (a_{ik} - \sum_{m=1}^{k-1} l_{im} u_{mk}) / u_{kk}$$

Doolittle decomposition.

Section 3. QR decomposition

$$A = QR.$$

First, decompose a matrix  $A$  into the product of an orthogonal matrix  $Q$  and a right triangular matrix  $R$ .

$$Ax = b \Rightarrow QRx = b.$$

$$\text{Let } Rx = y, \quad Qy = b.$$

$$\Downarrow \\ y = Q^T b.$$

$$\Rightarrow Rx = Q^T b.$$

$$R = \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & \dots & r_{2n} \\ & & \ddots & \vdots \\ & & & r_{nn} \end{pmatrix}.$$

$$Q = (q_1, q_2, \dots, q_n)$$

$q_j$ :  $j^{\text{th}}$  column vector of  $Q$ .

$$q_j \in \mathbb{R}^n.$$

$$(q_1, q_2, \dots, q_n) \begin{pmatrix} r_{11} & r_{12} & \dots & r_{1n} \\ & r_{22} & \dots & r_{2n} \\ & & \ddots & \vdots \\ & & & r_{nn} \end{pmatrix} = (a_1, a_2, \dots, a_n)$$

$$q_i^T q_j = \begin{cases} 1 & j=i \\ 0 & j \neq i \end{cases}$$

$$\|q_j\|_2 = 1.$$

$a_j$ :  $j^{\text{th}}$  column vector of  $A$ .  
 $a_j \in \mathbb{R}^n$

$$\begin{cases} q_1 r_{11} = a_1 & (1) & \|r_{11}\| = \|a_1\|_2 \Rightarrow r_{11} = \|a_1\|_2 \\ q_1 r_{12} + q_2 r_{22} = a_2 & (2) & q_1 = \frac{a_1}{r_{11}} \\ q_1 r_{13} + q_2 r_{23} + q_3 r_{33} = a_3 \\ q_1 r_{14} + q_2 r_{24} + q_3 r_{34} + q_4 r_{44} = a_4 \\ \vdots \\ \sum_{i=1}^j q_i r_{ij} = a_j & (3) \end{cases}$$

Take inner product of (2) with  $q_1$

$$(q_1, q_1) r_{12} + \cancel{(q_1, q_2) r_{22}}^0 = (q_1, a_2) \Rightarrow r_{12} = (q_1, a_2).$$

inner product of (2) with  $q_2$ .

$$(q_2, q_1) r_{12} + (q_2, q_2) r_{22} = (q_2, a_2) \Rightarrow r_{22} = \frac{(q_2, a_2) - (q_2, q_1) r_{12}}{1}$$

$$\underline{q_2 r_{22} = a_2 - q_1 r_{12}}$$

$$r_{22} = \|a_2 - q_1 r_{12}\|_2 \quad \checkmark$$

$$q_2 = (a_2 - q_1 r_{12}) / r_{22}$$

Take inner product of (3) with  $f_m$ ,  $m=1, 2, \dots, j-1$ .

$$(f_m, \sum_{i=1}^j f_i r_{ij}) = (f_m, a_j).$$

$$(f_m, f_m) r_{mj} = (f_m, a_j) \Rightarrow r_{mj} = (f_m, a_j) / (f_m, f_m).$$

$m=1, 2, \dots, j-1.$

$$f_j r_{jj} = a_j - \sum_{i=1}^{j-1} f_i r_{ij}$$

$$r_{jj} = \|a_j - \sum_{i=1}^{j-1} f_i r_{ij}\|_2. \Rightarrow f_j = (a_j - \sum_{i=1}^{j-1} f_i r_{ij}) / r_{jj}$$

Gram-Schmidt orthogonalization

$(a_1, a_2, \dots, a_n)$

normalize  $a_1$ :  $f_1 \leftarrow a_1, \|f_1\|_2 = 1.$

$$a_2 - (a_2, f_1) f_1 = p_2, \quad a_1 = r_{11} f_1, \quad r_{11} = \|a_1\|_2$$

$$f_1 = a_1 / r_{11}$$

$$r_{22} = \|p_2\|_2, \quad r_{22} f_2 = p_2 \Rightarrow f_2 = p_2 / r_{22}$$

$$a_j - \sum_{i=1}^{j-1} (a_j, f_i) f_i = p_j, \quad j=1, 2, \dots, n.$$

$$r_{jj} = \|p_j\|_2, \quad r_{jj} f_j = p_j \Rightarrow f_j = p_j / r_{jj}.$$

QR method to solve linear system:  
 $AX = b$

$$\begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}$$

Householder matrix.

(reflection matrix.)

$$H = I - 2WW^T$$

$$\|W\|_2 = 1. \checkmark$$

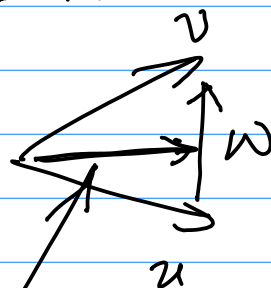
$$u, v \in \mathbb{R}^n$$

$$\|u\|_2 = \|v\|_2$$

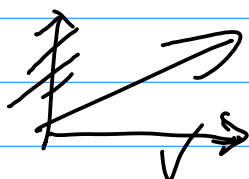
$$W = \frac{u-v}{\|u-v\|_2}$$

$$Hu = v$$

$$Hv = u.$$



$$(I - q_1 q_1^T) a_1 = \underline{a_1 - (q_1^T a_1) q_1}$$



$$(I - WW^T)v \Rightarrow$$

$$Hv = (I - 2WW^T)v = v - 2 \frac{(u-v)(u-v)^T}{\|u-v\|_2^2} v = u$$

$H$ : orthogonal matrix  
 $\downarrow$   
 symmetric

$$\begin{aligned} H^T H &= (I - 2WW^T)(I - 2WW^T) \\ &= I - 2WW^T - 2WW^T + 4 \underline{WW^T WW^T} = I \end{aligned}$$



$$\begin{pmatrix} u \\ a_{11} \\ a_{21} \\ \vdots \\ a_{n1} \end{pmatrix} \rightarrow \begin{pmatrix} v \\ r_{11} \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$$\|u\|_2 = |r_{11}| = \|u\|_2 = \|a_1\|_2$$

$$r_{11} = \pm \|a_1\|_2$$

$$\tilde{w}_1 = u - v = \begin{pmatrix} a_{11} - r_{11} \\ a_{21} \\ a_{31} \\ \vdots \\ a_{n1} \end{pmatrix}$$

$$w_1 = \frac{\tilde{w}_1}{\|\tilde{w}_1\|_2} = \frac{u - v}{\|u - v\|_2}$$

$$a_{11} - r_{11} = ?$$

$$H_1 = I - 2w_1 w_1^T$$

Next,

$$\tilde{a}_2 = \begin{pmatrix} a_{22} \\ a_{23} \\ \vdots \\ a_{2n} \end{pmatrix} \in \mathbb{R}^{n-1} \rightarrow \begin{pmatrix} r_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

$\|u_2\|_2 \qquad \qquad \|v_2\|_2$

$$|r_2| = \pm \|\tilde{a}_2\|_2$$

$$H_1 A x = H_1 b$$

$$w_2 = \frac{u_2 - v_2}{\|u_2 - v_2\|_2} \in \mathbb{R}^{n-1}$$

$$H_2 = I - 2w_2 w_2^T \in \mathbb{R}^{(n-1) \times (n-1)}$$

$$H_2 H_1 A x = H_2 H_1 b$$

↓

$$\begin{pmatrix} a_{33} \\ a_{43} \\ \vdots \\ a_{n3} \end{pmatrix} \in \mathbb{R}^{n-2} \rightarrow \begin{pmatrix} r_{33} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \in \mathbb{R}^{n-2}$$

$$H_{n-1} \cdots H_2 H_1 A x = H_{n-1} \cdots H_2 H_1 b$$

$$\| \quad \|$$

$$R x = H_{n-1} \cdots H_2 H_1 b$$

algorithm complexity  $O(n^3)$ .

## section 4. Stability Analysis

$$Ax = b. \quad (\text{original})$$

$$\tilde{A} \tilde{x} = \tilde{b} \quad (\text{practical})$$

step 1. only  $b$  is perturbed.

$$A \hat{x} = \hat{b}$$

$$\text{Let } \hat{b} = b + \delta b$$

$$\tilde{x} = x + \delta x.$$

perturbation  
 $\delta b \in \mathbb{R}^n.$

$\delta x$ : error.

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

$$Ax + A\delta x = b + \delta b$$

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

$\|\delta x\|$ : absolute error.

relative error:  $\frac{\|\delta x\|}{\|x\|}$

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

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

$$Ax = b.$$

$$\|b\| = \|Ax\| \leq \|A\| \cdot \|x\| \Rightarrow \|x\| \geq \|b\| / \|A\|.$$

$$\Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \frac{\|A^{-1}\| \cdot \|\delta b\|}{\|b\| / \|A\|} = \underbrace{\|A\| \cdot \|A^{-1}\|}_{\text{condition number of } A} \cdot \frac{\|\delta b\|}{\|b\|}$$

relative error

relative perturbation

condition number of  $A$

Norm of matrix  $A$ : related to vector norm

Definition of Vector Norm  $\| \cdot \|$

1).  $v \in \mathbb{R}^n$ .  $\|v\| \geq 0$   $\|v\| = 0$  iff  $v = 0$   
non-negativity.

2).  $\forall \lambda \in \mathbb{R}$ .  $\|\lambda v\| = |\lambda| \cdot \|v\|$   
homogeneity

3).  $\forall v, w \in \mathbb{R}^n$   $\|v + w\| \leq \|v\| + \|w\|$

Example:  $\|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$

$$\|v\|_\infty = \max_{1 \leq i \leq n} |v_i|$$

$$\|v\|_1 = \sum_{i=1}^n |v_i|$$

The norms are equivalent

(Norm Equivalence)

$$\| \cdot \|_a, \| \cdot \|_b$$

there exists two positive constants  $\mu, M > 0$

s.t.  $\mu \|v\|_b \leq \|v\|_a \leq M \|v\|_b \quad \forall v \in \mathbb{R}^n$

( $\mu, M$  are independent of  $v$ ).

$$(\|v\|_1, \|v\|_\infty)$$

$$(\|v\|_2, \|v\|_1)$$

$$(\|v\|_2, \|v\|_\infty)$$

$$\|v\|_2 \sim \|v\|_\infty$$

$$\|v\|_\infty = \max_{1 \leq i \leq n} |v_i| \leq \|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2} \leq \sqrt{n \max_{1 \leq i \leq n} v_i^2} = \sqrt{n} \|v\|_\infty$$

$\mu = 1 \quad M = \sqrt{n}$

$$\|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$$

$$\|v\|_1 = \sum_{i=1}^n |v_i|$$

$$\|v\|_{2,n} = \sqrt{\frac{1}{n} \sum_{i=1}^n v_i^2}$$

$$\|v\|_{1,n} = \frac{1}{n} \sum_{i=1}^n |v_i|$$

Matrix Norm :  $A \in \mathbb{R}^{n \times n}$

- 1).  $\|A\| \geq 0$ .  $\|A\| = 0$  iff  $A = 0$   
non-negativity.
- 2).  $\|\lambda A\| = |\lambda| \|A\|$ ,  $\lambda \in \mathbb{R}$ . homogeneity.
- 3).  $\|A+B\| \leq \|A\| + \|B\|$ . triangle inequality.
- 4).  $\|AB\| \leq \|A\| \|B\|$ .

Frobenius norm :  $A = (a_{ij})_{n \times n} \in \mathbb{R}^{n \times n}$

$$\|A\|_F = \sqrt{\sum_{i,j=1}^n a_{ij}^2}$$

Vector-norm induced norm of a matrix.  $A$

$$\|A\| = \sup_{\substack{v \in \mathbb{R}^n \\ v \neq 0}} \frac{\|Av\|}{\|v\|} \geq \frac{\|Aw\|}{\|w\|} \quad \forall w \in \mathbb{R}^n$$

Example :  $\|A\|_2 = \sup \frac{\|Av\|_2}{\|v\|_2}$

$$\|A\|_\infty = \sup \frac{\|Av\|_\infty}{\|v\|_\infty}$$

$$\|A\|_1 = \sup \frac{\|Av\|_1}{\|v\|_1}$$

Property :  $\|A\| \geq \frac{\|Aw\|}{\|w\|} \Rightarrow \|Aw\| \leq \|A\| \|w\|$

Remark: We can also use other vector norm to define an induced matrix norm.

$$\|v\|_p = \left( \sum_{i=1}^n |v_i|^p \right)^{\frac{1}{p}} \quad p \geq 1.$$

$$\|A\|_p = \sup \frac{\|Av\|_p}{\|v\|_p}.$$

Proposition : ①  $\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$

$$\begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{pmatrix}$$

$r_1, r_2, \dots, r_n \in \mathbb{R}^{1 \times n}$

$$\textcircled{2}. \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}| = \|A^T\|_\infty$$

$$\textcircled{3}. \|A\|_2 = [\rho(A^T A)]^{\frac{1}{2}}$$

Proof : ①.  $\|A\|_\infty = \sup_{\|v\|_\infty=1} \frac{\|Av\|_\infty}{\|v\|_\infty} = \sup_{\|v\|_\infty=1} \|Av\|_\infty$

$$= \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Let  $v = (v_1, v_2, \dots, v_n)^T$  be a normalized vector with  $\|v\|_\infty = \max_{1 \leq i \leq n} |v_i| = 1$ .

$$\begin{aligned} \|Av\|_\infty &= \max_{1 \leq i \leq n} \left| \sum_{j=1}^n a_{ij} v_j \right| \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| |v_j| \\ &\leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|. \Rightarrow \|A\|_\infty \leq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|. \end{aligned}$$

(1)

Let  $k$  be the row index so that

$$\sum_{j=1}^n |a_{kj}| = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

Choose a vector  $w = (w_1, w_2, \dots, w_n)^T \in \mathbb{R}^n$  s.t

$$w_j = \begin{cases} 1 & \text{if } a_{kj} \geq 0 \\ -1 & \text{otherwise} \end{cases} \quad \|w\|_\infty = 1.$$

$$\|Aw\|_\infty \geq |(Aw)_k| = \left| \sum_{j=1}^n a_{kj} w_j \right| = \sum_{j=1}^n |a_{kj}|$$

$$\wedge \quad \|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|$$

$$\Rightarrow \|A\|_\infty \geq \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (2)$$

$$(1), (2) \Rightarrow \|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|.$$

④

(3).  $\rho(B)$ : spectral radius

$$\rho(B) = \max_{1 \leq i \leq n} |\lambda_i(B)|$$

$$\|A\|_2 = [\rho(A^T A)]^{1/2}$$

$$\|A\|_2 = \sup \frac{\|Av\|_2}{\|v\|_2} \Rightarrow \|A\|_2^2 = \sup \frac{\|Av\|_2^2}{\|v\|_2^2}$$

$$= \sup \frac{(Av, Av)}{(v, v)} = \sup \frac{v^T A^T A v}{v^T v}$$

$A^T A$ : symmetric. non-negative definite

there exist a normalized mutually orthogonal basis, consisting of eigenvectors of  $A^T A$ .  
associated with eigenvalues

$$\Downarrow \\ (r_1, r_2, \dots, r_n) \in \mathbb{R}^n$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$$

$$r_i^T r_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

For any vector  $v$ , we have

$$v = c_1 r_1 + \dots + c_n r_n \quad c_i \in \mathbb{R}.$$

$$v^T v = \sum_{i=1}^n c_i^2$$

$$A^T A v = \sum_{i=1}^n \lambda_i c_i r_i \Rightarrow v^T A^T A v = \left( \sum_{i=1}^n c_i r_i \right)^T \sum_{i=1}^n \lambda_i c_i r_i$$

$$= \sum_{i=1}^n \lambda_i c_i^2$$

$$\frac{v^T A^T A v}{v^T v} = \frac{\sum_{i=1}^n \lambda_i c_i^2}{\sum_{i=1}^n c_i^2} \leq \frac{\lambda_1 \sum_{i=1}^n c_i^2}{\sum_{i=1}^n c_i^2} = \lambda_1 = \rho(A^T A).$$

The equality holds when  $c_2 = c_3 = \dots = c_n = 0$

$$\|A\|_2^2 = \sup_{\|v\| \neq 0} \frac{v^T A^T A v}{v^T v} = \lambda_1 = \rho(A^T A)$$

$$\Rightarrow \|A\|_2 = \rho(A^T A)^{1/2}$$

⊗

property: Let  $\|\cdot\|$  be an induced matrix norm.

if  $\|A\| < 1$  for some matrix  $A$ , then

$I+A$  is invertible and

$$\|(I+A)^{-1}\| \leq \frac{1}{1-\|A\|}.$$

Proof: Assume  $I+A$  is not invertible.

there exists a non-zero vector  $v \in \mathbb{R}^n$  s.t.

$$(I+A)v = 0$$

$$\Rightarrow v = -Av \Rightarrow \|v\| \leq \|A\| \|v\| \Rightarrow 1 \leq \|A\|$$

This is a contradiction

$$\text{Let } C = (I+A)^{-1}.$$

$$(I+A)C = I.$$

$$C = I - AC \quad \|C\| \leq \|I\| + \|A\| \|C\| = 1 + \|A\| \|C\|$$

$$\Rightarrow (1 - \|A\|) \|C\| \leq 1$$

$$\Rightarrow \|C\| \leq \frac{1}{1 - \|A\|} \quad \textcircled{H}$$

Proposit:  $\|A\| \|A^{-1}\|$ : condition number of  $A$   
 $\|A\| \|A^{-1}\| \geq 1$  in induced norm.

$$1 = \|I\| = \|AA^{-1}\| \leq \|A\| \|A^{-1}\|$$

Example of an ill-condition matrix

$$H = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{3} & \frac{1}{4} \\ \frac{1}{2} & \frac{1}{3} & \frac{1}{4} & \frac{1}{5} \\ \frac{1}{3} & \frac{1}{4} & \frac{1}{5} & \frac{1}{6} \\ \frac{1}{4} & \frac{1}{5} & \frac{1}{6} & \frac{1}{7} \end{pmatrix}$$

Hilbert matrix

$$x \in \mathbb{R}^n \quad Hx = b \quad \text{solve } Hy = b \text{ using LU or QR}$$

$$Hx \rightarrow b \quad \|x - y\| = ?$$

case II. matrix is perturbed

$$\tilde{A} = A + \delta A \quad \delta A \in \mathbb{R}^{n \times n}$$

$$\tilde{A} \tilde{x} = b$$

$$\Rightarrow (A + \delta A)(x + \delta x) = b \quad \delta x \in \mathbb{R}^n$$

$$\cancel{Ax} + A\delta x + \delta A x + \delta A \cdot \delta x = \cancel{b} \cdot 0$$

$$(A + \delta A) \delta x = -\delta A \cdot x$$



$$\delta x = -(A + \delta A)^{-1} \cdot \delta A \cdot x$$

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

$$\Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \|(A + \delta A)^{-1}\| \cdot \|\delta A\| \quad (3)$$

$$A + \delta A = A(I + A^{-1}\delta A)$$

$$(A + \delta A)^{-1} = (I + A^{-1}\delta A)^{-1} \cdot A^{-1}$$

$$(3) \Rightarrow \frac{\|\delta x\|}{\|x\|} \leq \|(I + A^{-1}\delta A)^{-1}\| \cdot \|A^{-1}\| \cdot \|\delta A\|$$

Assume  $\|A^{-1}\delta A\| < 1$  Assume  $\|A^{-1}\| \cdot \|\delta A\| < 1$

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

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

$$= \frac{1}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \cdot \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}$$

$$= \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \cdot \frac{\|\delta A\|}{\|A\|} = \frac{\text{Cond}(A)}{1 - \text{Cond}(A) \cdot \frac{\|\delta A\|}{\|A\|}} \cdot \frac{\|\delta A\|}{\|A\|}$$

The relative error also depends on the condition number of the matrix  $A$ .

case II. Both  $b$  and  $A$  are perturbed.

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

$$\Rightarrow \cancel{Ax} + \delta A \cdot x + (A + \delta A)\delta x = \cancel{b} + \delta b$$

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

$$\delta x = (A + \delta A)^{-1} [-\delta A \cdot x + \delta b]$$

$$\|\delta x\| \leq \|(A + \delta A)^{-1}\| \cdot (\|\delta A\| \cdot \|x\| + \|\delta b\|)$$

$$\frac{\|\delta x\|}{\|x\|} \leq \|(A + \delta A)^{-1}\| \cdot \left( \|\delta A\| + \frac{\|\delta b\|}{\|x\|} \right) \quad (5)$$

$$b = Ax \Rightarrow \|b\| \leq \|A\| \cdot \|x\| \Rightarrow \|x\| \geq \frac{\|b\|}{\|A\|}$$

$$(5) \leq \|(A + \delta A)^{-1}\| \cdot \left( \|\delta A\| + \|A\| \cdot \frac{\|\delta b\|}{\|b\|} \right)$$

$$= \|(A + \delta A)^{-1}\| \cdot \|A\| \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right)$$

$$\leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|A\| \cdot \frac{\|\delta A\|}{\|A\|}} \cdot \left( \frac{\|\delta A\|}{\|A\|} + \frac{\|\delta b\|}{\|b\|} \right).$$

Remark: 1. When the matrix  $A$  is symmetric and positive definite, the LU decomposition may take the form  $LI^T = A$ .

with  $L = \begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & & \ddots & \\ l_{n1} & l_{n2} & & l_{nn} \end{pmatrix}$

$$\begin{pmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ \vdots & & \ddots & \\ l_{n1} & l_{n2} & & l_{nn} \end{pmatrix} \begin{pmatrix} l_{11} & l_{21} & l_{31} & \\ & l_{22} & l_{32} & \\ & & l_{33} & \\ & & & \ddots & \\ & & & & l_{nn} \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \dots \\ a_{21} & a_{22} & \\ a_{31} & & \\ \vdots & & \\ a_{n1} & & \end{pmatrix}$$

$$l_{11}^2 = a_{11} \Rightarrow l_{11} = \sqrt{a_{11}}$$

$$l_{i1} l_{11} = a_{i1} \quad i = 1, 2, \dots, n$$

$$\Rightarrow l_{i1} = a_{i1} / l_{11}$$

$$l_{21}^2 + l_{22}^2 = a_{22} \Rightarrow l_{22} = \sqrt{a_{22} - l_{21}^2}$$

...

The method is called the Cholesky decomposition

2. When the matrix  $A$  is tri-diagonal

$$A = \begin{pmatrix} a_1 & c_1 & & \\ b_1 & a_2 & c_2 & \\ & b_2 & a_3 & c_{n-1} \\ & & \ddots & \ddots & c_n \\ & & & b_{n-1} & a_n \end{pmatrix} = LU.$$

$$L = \begin{pmatrix} 1 & & & \\ \beta_1 & 1 & & \\ & \beta_2 & 1 & \\ & & \ddots & \ddots \\ & & & \beta_{n-1} & 1 \end{pmatrix} \quad U = \begin{pmatrix} \gamma_1 & \delta_1 & & \\ & \gamma_2 & \delta_2 & \\ & & \gamma_3 & \ddots \\ & & & \ddots & \gamma_{n-1} & \delta_n \end{pmatrix}$$

$$\gamma_1 = a_1 \quad \beta_1 \gamma_1 = b_1 \Rightarrow \beta_1 = b_1 / a_1$$

$$\gamma_1 = c_1 \quad \beta_1 \gamma_1 + \alpha_2 = a_2 \Rightarrow \alpha_2 = a_2 - \beta_1 c_1.$$

$$\beta_2 \alpha_2 = b_2 \Rightarrow \beta_2 = b_2 / \alpha_2.$$

...

$$\boxed{LUx = f} \Rightarrow \begin{cases} Ly = f \\ Ux = y \end{cases}$$

The algorithm complexity is  $O(n)$ .