

$$\begin{bmatrix} Q_{11} & \dots & Q_{13} \\ \vdots & & \vdots \\ Q_{31} & \dots & Q_{33} \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ \vdots \\ x_3 \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_3 \end{bmatrix}$$

$Q_{11}x_1 + Q_{12}x_2 +$
 $+ Q_{13}x_3 = b_1$
 coordinates

Basis: a triple of vectors s.t. we can write every vector b as a linear combination of them

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_m \end{bmatrix} \quad w = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix} \quad v \cdot w \equiv v^T w$$

$v_1 \cdot w_1 + \dots + v_m \cdot w_m$

$$v \cdot w = \|v\| \cdot \|w\| \cdot \cos \theta$$

Orthogonal if $v \cdot w = 0$

Mapp $x \rightarrow Ax$ geometric transformation: rotation, scaling, mirroring.

$$A \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} x_1 + x_2 \\ x_2 \end{bmatrix}$$



$$I = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad I \cdot v = v \text{ identity matrix}$$

$v \rightarrow A(Bv)$ composition is a transformation that corresponds to the matrix product AB

$$v \rightarrow A(Bv) = (AB)v$$

Matrix product:

$$\begin{array}{c|c} i & A \\ \hline \end{array} \cdot \begin{array}{c|c} j & B \\ \hline \end{array} = \begin{array}{c|c} i & C_{ij} \\ \hline \end{array}$$

$m \times n$ $n \times p$

must be equal

$$C_{ij} = a_{i1} \cdot b_{1j} + a_{i2} \cdot b_{2j} + \dots + a_{in} \cdot b_{nj}$$

vectors = $m \times 1$ matrices

Matrix algebra:

- $A(B+C) = AB + AC$ $(AB)C = A(BC)$
- $(A+B)^2 = A^2 + B^2 + \underline{AB + BA}$ $AB \neq BA$

$$AB = 0 \rightarrow B = 0 \text{ or } A = 0$$

$$AB = AC \rightarrow B = C$$

Rank of a matrix A : is the smallest number such that all the columns of A can be written as linear combinations of n vectors v_1, v_2, \dots, v_n

Block operation:

$$\begin{matrix} m_1 & m_2 \\ m_2 & \\ m_3 & \end{matrix} \left[\begin{array}{c|c} A_{11} & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right] \cdot \begin{matrix} n_1 & n_2 & n_3 \\ n_2 & \\ n_3 & \end{matrix} \left[\begin{array}{c|c|c} P_1 & P_2 & P_3 \\ \hline \hline B_{11} & B_{12} & \\ B_{21} & B_{22} & \\ \hline \hline C_{11} & C_{12} & C_{13} \end{array} \right] = \begin{matrix} n_1 & n_2 & n_3 \\ n_2 & \\ n_3 & \end{matrix} \left[\begin{array}{c|c|c} P_1 & P_2 & P_3 \\ \hline \hline C_{11} & C_{12} & C_{13} \end{array} \right]$$

$$C_{12} = A_{11} \cdot B_{112} + A_{12} \cdot B_{212}$$

Linear systems: given $b \in \mathbb{R}^m$, $A \in \mathbb{R}^{m \times n}$, which are the coordinates, I need to write b as a linear combination of the columns of A .

$$b = \begin{bmatrix} 4 \\ 4 \\ 0 \end{bmatrix} \quad A = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad Ax = \begin{bmatrix} 2 \\ 0 \\ 0 \end{bmatrix} x_1 + \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} x_2 + \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} x_3 = b$$

If solution not available, we can get close to b whilky

$$\min \|Ax - b\|$$

If A is square ($n \times n$) and invertible: $x = A^{-1}b$, $A, b \in \mathbb{R}^n$

$$A^{-1} \cdot A = A \cdot A^{-1} = I = \begin{bmatrix} 1 & & & \\ & \ddots & & 0 \\ 0 & & \ddots & \\ & & & 1 \end{bmatrix}$$

Orthogonality: $U \in \mathbb{R}^{n \times n}$ is orthogonal if $U^T U = I$ / $U U^T = I$

$$U^{-1} = U^T$$

If U is orthogonal, then $\forall v \in \mathbb{R}^n \|Uv\| = \|v\|$

$$\|v\|^2 = v^T v = v^T \cdot I \cdot v = v^T (U^T U) v = (Uv)^T (Uv) = \|Uv\|^2$$

$$\|x\|_A = x^T A x$$

Def: $\{u_1, u_2, \dots, u_n\}$ is orthonormal if $u_i^T u_j = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$

$$I = U^T U$$

$$\begin{bmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{bmatrix} = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_n^T \end{bmatrix} \cdot \begin{bmatrix} u_1 | u_2 | \dots | u_n \end{bmatrix} = \begin{bmatrix} u_1^T u_1 & u_1^T u_2 & \dots \\ u_2^T u_1 & u_2^T u_2 & \dots \\ \vdots & \vdots & \ddots \\ u_n^T u_1 & u_n^T u_2 & \dots & u_n^T u_n \end{bmatrix}$$

Def: $\mathbf{U}_1, \mathbf{U}_2$ orthogonal $\rightarrow \mathbf{U}_1 \mathbf{U}_2$ is orthogonal

$$(\mathbf{U}_1 \cdot \mathbf{U}_2)^\top (\mathbf{U}_1 \cdot \mathbf{U}_2) = \mathbf{U}_2^\top \cdot \mathbf{U}_1^\top \cdot \mathbf{U}_1 \cdot \mathbf{U}_2 = \mathbf{U}_2^\top \cdot \mathbf{U}_2 = \mathbb{I}$$

\mathbb{I} \mathbb{I}

Eigenvalues / Eigenvectors

Given \mathbf{A} square matrix $\mathbf{A} \in \mathbb{R}^{m \times m}$ if $\mathbf{A}\mathbf{v} = \mathbf{v}\lambda$ for a certain vector $\mathbf{v} \in \mathbb{R}^m$, $\lambda \in \mathbb{R}$, then we say that \mathbf{v} is an eigenvector and λ is an eigenvalue of \mathbf{A} .

For almost all matrices:

$$\mathbf{A} = \mathbf{V} \Lambda \mathbf{V}^{-1} = \underbrace{\left[\mathbf{v}_1 \mid \mathbf{v}_2 \mid \dots \mid \mathbf{v}_m \right]}_{\text{eigenvectors of } \mathbf{A}} \underbrace{\begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_m & \end{bmatrix}}_{\text{eigenvalues of } \mathbf{A}} \underbrace{\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}}_{=}$$

$$\mathbf{A} \mathbf{v}_1 = \mathbf{v}_1 \lambda_1, \quad \mathbf{A} \mathbf{v}_2 = \mathbf{v}_2 \lambda_2, \dots, \quad \mathbf{A} \mathbf{v}_m = \mathbf{v}_m \lambda_m$$

$$= \mathbf{v}_1 \lambda_1 w_1 + \mathbf{v}_2 \lambda_2 w_2 + \dots + \mathbf{v}_m \lambda_m w_m$$

$$\underbrace{\boxed{} \quad \boxed{} \quad \boxed{}} + \boxed{} \quad \boxed{} + \boxed{} \quad \boxed{} \quad \text{sum of } m \text{ non-zero terms}$$

$$[\mathbf{V}, \Lambda] = \text{eig}(\mathbf{A}) \quad \text{cost } O(m^2)$$

Eigenvectors are less unique: if v eigenvector, λv is also an eigenvector (with the same λ)

If v_1, v_2 eigenvectors with $\lambda_1 = \lambda_2 = \lambda$, then $\alpha v_1 + \beta v_2$ is also an eigenvector with eigenvalue λ

$$A(v_1\alpha + v_2\beta) = A v_1 \alpha + A v_2 \beta = v_1 \lambda \alpha + v_2 \lambda \beta = (v_1\alpha + v_2\beta)\lambda$$

The (Spectral theorem) If $A = A^T$ (symmetric), then

$A = V \Lambda V^{-1}$ always exists, the eigenvalues λ_i are real,

V can be taken orthogonal

Th: let A be symmetric with minimum eigenvalue λ_{\min} and maximum λ_{\max} . Then, $\lambda_{\min} \|x\|^2 \leq x^T A x \leq \lambda_{\max} \|x\|^2$

Proof: let us first do the case $A = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_m \end{bmatrix}$, then

$$\begin{bmatrix} x_1 & x_2 & \dots & x_m \end{bmatrix} \begin{bmatrix} \lambda_1 & & & \\ & \ddots & & \\ & & \lambda_m & \\ & & & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \begin{bmatrix} x_1 & \dots & x_m \end{bmatrix} \begin{bmatrix} x_1 & & \\ & \ddots & \\ & & x_m \end{bmatrix}$$

$$\lambda_{\min} (x_1^2 + \dots + x_m^2) \leq \lambda_1 x_1^2 + \lambda_2 x_2^2 + \dots + \lambda_m x_m^2 \leq \lambda_{\max} (x_1^2 + \dots + x_m^2)$$

Def: if $\lambda_i \geq 0$ for each eigenvalue λ_i of \mathbf{Q} , then

$x^T \mathbf{Q} x \geq 0$. In this case we call \mathbf{Q} positive semidefinite.

Def: if $\lambda_i > 0$ for each eigenvalue λ_i of \mathbf{Q} , then

$x^T \mathbf{Q} x > 0$ for each vector $x \neq 0$. In this case, we call \mathbf{Q} positive definite.

\downarrow
 \mathbf{Q} is invertible

Th: let $A \in \mathbb{R}^{m \times n}$. Then, $A^T A$ is symmetric, positive semidefinite matrix. It is also valid for $A^T A^T$.

Proof: $(A^T A)^T = A^T \cdot (A^T)^T = A^T A$, as the product is sym.

for each $x \in \mathbb{R}^m$

$$x^T (A^T A) x = (Ax)^T (Ax) = \|Ax\|^2 \geq 0$$

Singular value decomposition $A \in \mathbb{R}^{m \times n}$, each square matrix

can be written as:

$$A = U \Sigma V^T = \begin{bmatrix} u_1 | u_2 | \dots | u_m \end{bmatrix} \cdot \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & 0 \\ 0 & & & \sigma_m \end{bmatrix} \cdot \begin{bmatrix} v_1^T \\ v_2^T \\ \vdots \\ v_m^T \end{bmatrix}$$

with U orthogonal, Σ diagonal, V orthogonal

and $\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_m \geq 0$ singular values

$$= \mu_1 \sigma_1 v_1^T + \mu_2 \sigma_2 v_2^T + \dots + \mu_m \sigma_m v_m^T$$

$$\boxed{} \quad \boxed{} \quad \boxed{} + \boxed{} \quad \boxed{} \quad \boxed{} + \dots + \boxed{} \quad \boxed{} \quad \boxed{}$$

σ_i unique for each matrix, μ_i and v_i may be not.

Th: If A rectangular $\in \mathbb{R}^{m \times n}$, it can always be decomposed as $A = U \Sigma V^T$ $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ orthogonal,

$\Sigma \in \mathbb{R}^{m \times n}$ diagonal (i.e. $\Sigma_{ij} = 0$ whenever $i \neq j$)

$m > n$ (A tall thin)

$$A = \begin{matrix} U \\ m \times m \end{matrix} \quad \begin{matrix} \Sigma \\ m \times m \end{matrix} \quad \begin{matrix} V^T \\ m \times n \end{matrix}$$

$n > m$ (A short fat)

$$A = \begin{matrix} U \\ m \times m \end{matrix} \quad \begin{matrix} \Sigma \\ m \times m \end{matrix} \quad \begin{matrix} V^T \\ m \times n \end{matrix}$$

$$\sigma_1 \geq \sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_{\min(m, n)} \geq 0$$

$$A = \left[\mu_1 \left[\mu_2 \left(\dots \left[\mu_n \right] \right) \cdot \begin{bmatrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{bmatrix} \right] \cdot \begin{bmatrix} v_1^T \\ v_2^T \\ v_3^T \end{bmatrix} \right]$$

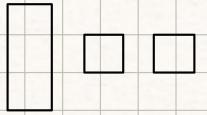
if $m > n$, $\mu_{n+1} \dots \mu_m$ match up with zeros and so not

appear in the product. I can write

$$U = \begin{bmatrix} U_1 & | & U_2 \\ m & & m-m \end{bmatrix} \quad \Sigma = \begin{bmatrix} \Sigma \\ 0 \end{bmatrix}^m \quad A = U \Sigma V^T = U, \Sigma, V^T$$

$m \times m \quad m \times m \quad m \times m$

$$U \Sigma = U_1 \Sigma_1 + \cancel{U_2 \cdot 0}$$



Properties of SVD

- $\text{rank}(A) = \text{number of non-zero } \sigma_i$'s:

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_m = 0$$

$$A = \mu_1 \sigma_1 V_1^T + \dots + \mu_r \sigma_r V_r^T + \cancel{\mu_{r+1} \sigma_{r+1} V_{r+1}^T \dots}$$

$\text{Imag}(A) = \text{span} \{ \mu_1, \mu_2, \dots, \mu_r \}$. Vectors are linearly independent column vector of A , and r is the rank of A , which is the dimension of the image.

$$\text{Ker}(A) = \{ z \in \mathbb{R}^m \text{ s.t. } A z = 0 \}$$

When V_K , with $K > r$ $A V_K = \mu_1 \sigma_1 V_1^T V_K + \dots + \mu_r \sigma_r V_r^T V_K = 0$

Hence $V_{r+1}, V_{r+2}, \dots, V_m$ are in $\text{Ker}(A)$.

For fact, $\text{Ker}(A) = \text{span} \{ V_{r+1}, V_{r+2}, \dots, V_m \}$

V is orthogonal

$A = U \Sigma V^T$ suppose invertible, SVD of A^{-1} :

$$A^{-1} = (V^T)^{-1} \Sigma^{-1} U^{-1} = V \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & 0 \\ 0 & & \frac{1}{\sigma_m} \end{bmatrix} U^T = \begin{bmatrix} U_1 & | & U_2 & | & \dots & | & U_m \end{bmatrix} \Sigma^{-1} \begin{bmatrix} M_1 \\ \vdots \\ M_m \end{bmatrix}^T$$

$$= U_1 \frac{1}{\sigma_1} M_1^T + U_2 \frac{1}{\sigma_2} M_2^T + \dots + U_m \frac{1}{\sigma_m} M_m^T$$

$$= \begin{bmatrix} U_m & | & U_{m-1} & | & \dots & | & U_1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma_m} & & \\ & \ddots & & \\ & & \frac{1}{\sigma_1} \end{bmatrix} \begin{bmatrix} M_m^T \\ \frac{M_{m-1}^T}{M_m^T} \\ \vdots \\ \frac{M_1^T}{M_2^T} \end{bmatrix}$$

SVD of A^{-1}

orth.

diag.

orth.

Let $A = U \Sigma V^T$, eigenvalue decomposition of $A^T A$:

$$A^T A = (U \Sigma V^T)^T (U \Sigma V^T) = V \Sigma^T \underbrace{U^T U}_{I} \Sigma V^T$$

$$\Sigma^T \Sigma = \begin{bmatrix} \sigma_1 & \dots & 0 & | & 0 \\ & \ddots & & | & \\ 0 & & \sigma_m & | & \end{bmatrix} \begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & & \\ & & \sigma_m \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & & \\ 0 & & \sigma_m^2 \end{bmatrix}$$

$$A^T \cdot A = V \begin{bmatrix} \sigma_1^2 & & 0 \\ & \ddots & & \\ 0 & & \sigma_m^2 \end{bmatrix} V^T$$

Matrix norm: $\|x\| = (\mathbf{x}^\top \mathbf{x})^{1/2}$

Def: given any norm $\|\cdot\|$ on vectors, we can define a norm on matrices as:

$$\|M\| = \max_{\substack{v \in \mathbb{R}^{n \times m} \\ v \neq 0}} \frac{\|Mv\|}{\|v\|} = \max_{\substack{u \in \mathbb{R}^m \\ \|u\|=1}} \|Mu\|$$

$v = \lambda u$
↑ scale norm 1

- $\|Mv\| \leq \|M\| \cdot \|v\|$
- $\|A\| \geq 0$ equality only for $A=0$
- $\|\lambda A\| = |\lambda| \cdot \|A\|$
- $\|A+B\| \leq \|A\| + \|B\|$
- $\|AB\| \leq \|A\| \cdot \|B\|$
- $\|A\omega\| \leq \|A\| \cdot \|\omega\|$

If we have $\|\cdot\| = \|\cdot\|_2$, then: $\|UA\|_2 = \|A\|_2$

$\|A\omega\|_2 = \|A\|_2$ with ω orthogonal

If $A = U\Sigma V^\top$, then $\|A\|_2 = \|\Sigma\|_2 = \|\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix}\|_2 =$

$\max(\sigma_1, \dots, \sigma_m) = \sigma_1$ why? $\left\| \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_m \end{bmatrix} \begin{bmatrix} u_1 \\ \vdots \\ u_m \end{bmatrix} \right\|_2 = \left\| \begin{bmatrix} \sigma_1 u_1 \\ \vdots \\ \sigma_m u_m \end{bmatrix} \right\|_2 =$

$$= \sqrt{b_1^2 \mu_1^2 + \dots + b_m^2 \mu_m^2} \leq \sqrt{\underbrace{b_1^2 (\mu_1^2 + \dots + \mu_m^2)}_1} = b_1$$

Frobenius norm: $\|A\|_F = \left(\sum_{i,j} e_{i,j}^2 \right)^{1/2}$

satisfies all properties of norms, including $\|AU\|_F = \|UA\|_F = \|A\|_F$

But $\|\mathbb{I}\|_2 = 1 \quad \|\mathbb{I}\|_F = \sqrt{m}$

Th: for both the Frobenius and the ∞ -norm, the matrix X

that achieves $\min_{X \in \mathbb{R}^{m \times m}} \|A - X\|$ is $X_k = \mu_1 b_1 v_1^T + \dots + \mu_k b_k v_k^T$
 $\text{rank } X \leq k$

where $U, \Sigma, V = \text{SVD}(A)$

$$\rightarrow X_k = \begin{bmatrix} \mu_1 | \mu_2 | \dots | \mu_k \end{bmatrix} \begin{bmatrix} b_1 & b_2 & \dots & b_k \\ \vdots & \ddots & & 0 \\ 0 & 0 & \dots & b_k \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_k^T \end{bmatrix} =$$

$$= \begin{bmatrix} \mu_1 | \mu_2 | \dots | \mu_m \end{bmatrix} \begin{bmatrix} b_1 & & 0 \\ \dots & b_k & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} v_1^T \\ \vdots \\ v_m^T \end{bmatrix}$$

$$\|A - X\| = \left\| U \left(\begin{bmatrix} b_1 & 0 \\ \vdots & b_k \\ 0 & b_m \end{bmatrix} - \begin{bmatrix} b_1 & & 0 \\ \dots & b_k & 0 \\ 0 & 0 & \dots & 0 \end{bmatrix} \right) V^T \right\| = \left\| \begin{bmatrix} 0 & & & & 0 \\ \dots & 0 & b_{k+1} & 0 \\ 0 & 0 & \dots & b_m \end{bmatrix} \right\|$$

$$\text{error of rank } k-1 \text{ approx } \|A - X_1\|_F^2 = b_1^2 + b_2^2 + \dots + b_m^2$$

if $b_1 \gg b_2, b_3, \dots$ rank-1 approximation is a very good one

Principal component analysis:

$x_1, x_2, \dots, x_n \in \mathbb{R}^m$ data (length- m vectors of features)

$$A = \underbrace{\begin{bmatrix} x_1 & | & x_2 & | & \dots & | & x_m \end{bmatrix}}_{\text{examples}} \} \text{features}$$

① Removing the mean $\mu = \frac{1}{n} (x_1 + \dots + x_n)$ $\hat{x}_i = x_i - \mu \quad i=1..n$

② Consider SVD of $\hat{A} = [\hat{x}_1 | \dots | \hat{x}_m]$ and interpret u_i, v_i as different features of the data

OR

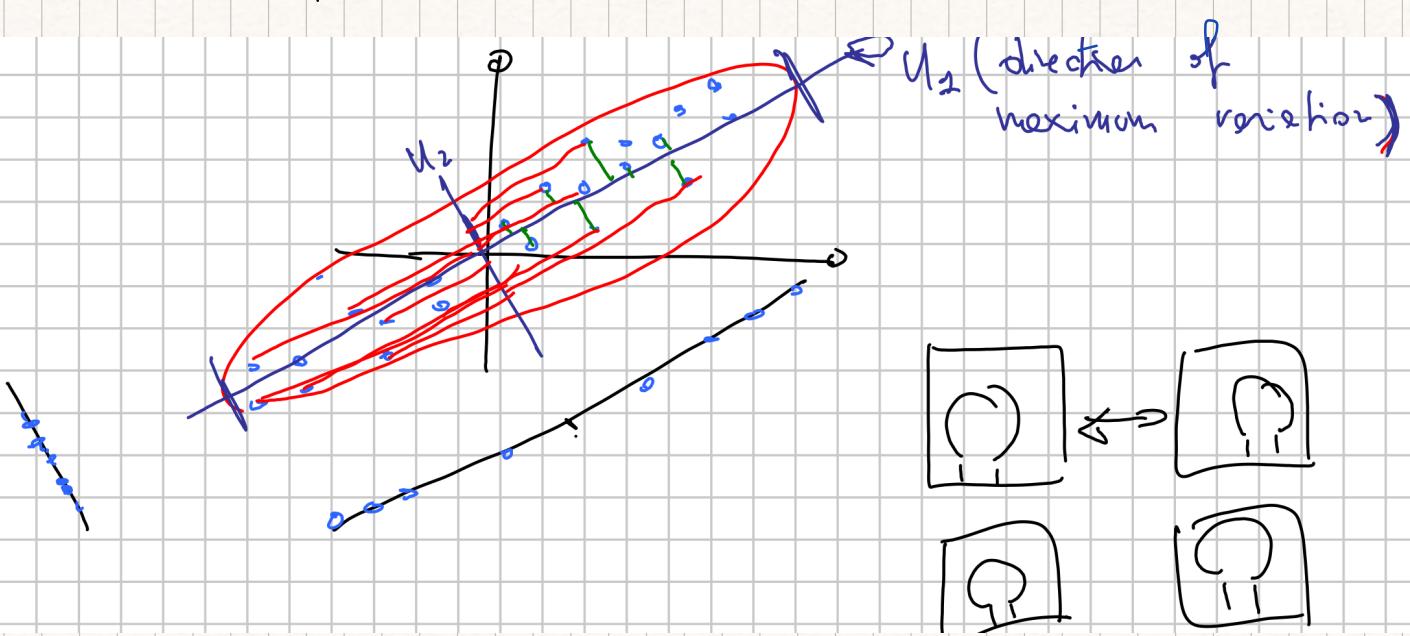
② as told by statistician: build covariance matrix

$$C = \frac{1}{n} (\hat{x}_1 \hat{x}_1^\top + \hat{x}_2 \hat{x}_2^\top + \dots + \hat{x}_m \hat{x}_m^\top) \quad C \text{ is positive definite}$$

PCA builds a basis that lets us write:

$$\hat{x}_j = x_j - \mu = \mu_1 \lambda_{1j} + \mu_2 \lambda_{2j} + \dots + \mu_m \lambda_{mj} \quad j=1, 2, \dots, m$$

where \hat{x}_j vary the most in the direction u_1 ($\sum \lambda_{1j}^2$ is max.)



Then how we get the λ_{ij} ?

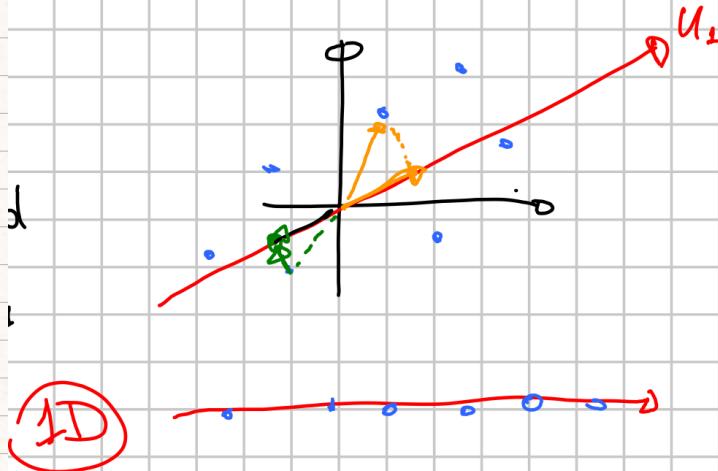
$$\hat{x}_j = \hat{A} \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \frac{\mu_1 b_1 v_1^T e_j}{\lambda_{1j}} + \frac{\mu_2 b_2 v_2^T e_j}{\lambda_{2j}} + \dots + \frac{\mu_m b_m v_m^T e_j}{\lambda_{mj}}$$

$$L = \sum V^T = U^T \hat{A} \quad \hat{x}_j = U L_j \quad L_j = U^T \hat{x}_j$$

Dimensionality reduction / plotting with PCA

Idea: direction of maximum variation gives me a "better plot."

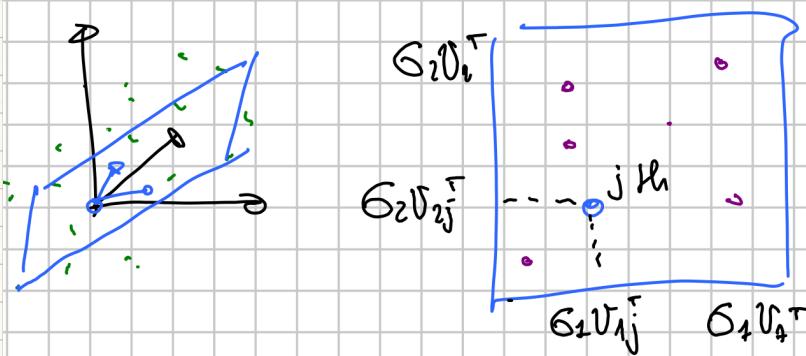
①D plotting \leftrightarrow rank-1 approximation : all of my data point are replaced by multiples of the same vector μ_1



$$A \approx X_1 = \mu_1 b_1 v_1^T$$

(2D) : plotting \leftrightarrow rank-2 approximation : all data point replaced by linear combination of 2 vectors

$$A \approx X_2 = \mu_1 b_1 v_1^T + \mu_2 b_2 v_2^T$$



Limitations:

- Euclidean distance only
- no labeling of μ_i 's.
- "flattening" the data: uses only matrix structure rather than the 4D-structure of any data
- we are lucky that rank-k approximation of a 2D-matrix $\min \|A - X\|$ has an easy solution in $\|\cdot\|_F, \|\cdot\|_2$, $\text{rank } X \leq k$

Least squares problem: Best approximation problem: given vectors $\alpha_1, \alpha_2, \dots, \alpha_m \in \mathbb{R}^m$ and a "target" vector $b \in \mathbb{R}^m$, find the coefficients that represent a linear combination of the α_i 's that goes as close as possible to b

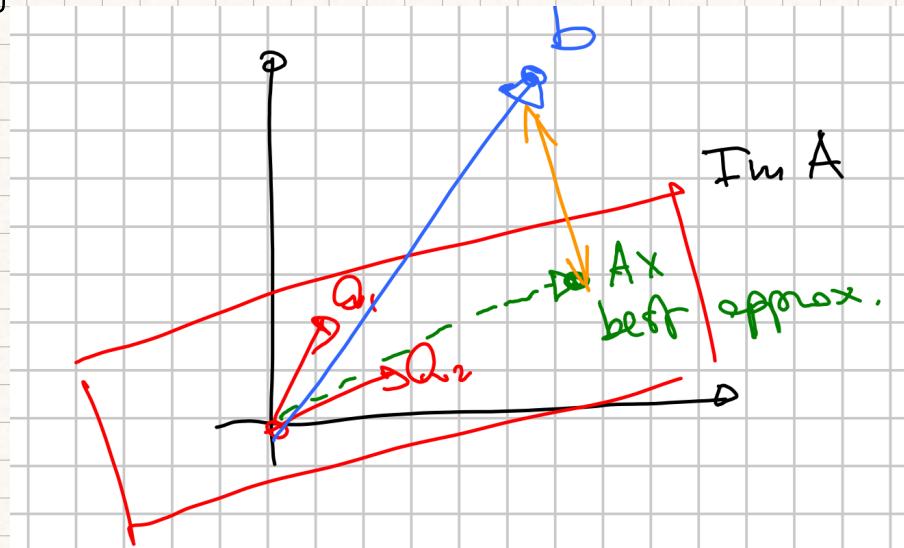
$$\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_m x_m = Ax$$

"

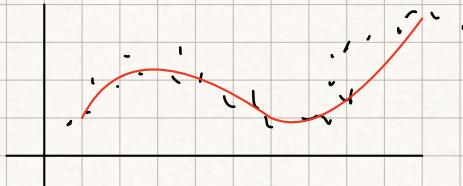
$$\begin{bmatrix} \alpha_1 | \alpha_2 | \dots | \alpha_m \end{bmatrix} \in \mathbb{R}^{m \times m} \quad x \in \mathbb{R}^m$$

We know from linear algebra that if A were invertible then there is one solution that gets exactly b

But in general $m \neq n$ and no exact solution $\min \|Ax - b\|_2$



Polynomial fitting: given point $(x_1, y_1), (x_2, y_2) \dots (x_m, y_m)$



Let's find the degree - 3 polynomial

$$y = p(x) = ax^3 + bx^2 + cx + d \quad a, b, c, d \text{ unknown}$$

such that $p(x_i) \approx y_i$ min $\sum_{i=1}^n (p(x_i) - y_i)^2$

$$\min_{a, b, c, d} \| \begin{bmatrix} x_1^3 & x_1^2 & x_1 & 1 \\ x_2^3 & x_2^2 & x_2 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ x_m^3 & x_m^2 & x_m & 1 \end{bmatrix} \cdot \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix} - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \|^2 = \| \begin{bmatrix} p(x_1) - y_1 \\ p(x_2) - y_2 \\ \vdots \\ p(x_m) - y_m \end{bmatrix} \|^2$$

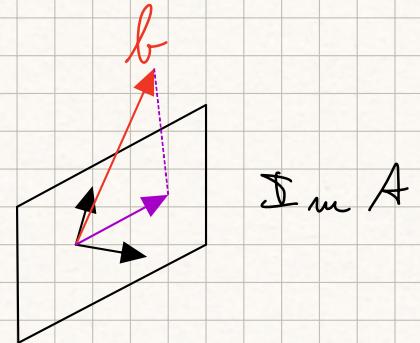
A x b

Least square problem

$$\min_{x \in \mathbb{R}^n} \|Ax - b\|_2$$

$$A \in \mathbb{R}^{m \times n}$$

$$b \in \mathbb{R}^m$$



$$A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 1 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad b = \begin{bmatrix} 4 \\ 4 \\ 0 \\ 1 \end{bmatrix}$$

We have multiple solutions:

$$x_1 = \begin{bmatrix} 4 \\ 2 \\ 0 \end{bmatrix} \quad x_2 = \begin{bmatrix} 0 \\ 0 \\ 4 \end{bmatrix} \quad x_3 = \begin{bmatrix} 2 \\ 1 \\ 2 \end{bmatrix}$$

$$Ax_1 = Ax_2 = Ax_3 = \begin{bmatrix} 4 \\ 4 \\ 0 \\ 1 \end{bmatrix}$$

this because exist $z \neq 0$ s.t. $Az = 0$

$$z = \begin{bmatrix} 2 \\ 1 \\ -2 \end{bmatrix} \quad \text{and multiple } \ker A = \text{span}(z)$$

Def: $A \in \mathbb{R}^{m \times n}$ has full column rank if $\text{rank}(A) = m$
 (i.e. $\ker A = \{0\}$, there are no $z \neq 0$ s.t. $Az = 0$)

The min $\|Ax - b\|$ has multiple solutions whenever A does not have full column rank.

Theorem: A has full column rank if and only if $A^T A$ is positive definite. (Recall: $A^T A$ is p. d. if $z^T A^T A z > 0; \forall z \neq 0$)

Proof: A has full col. rank $\Leftrightarrow A z \neq 0$ for all $z \neq 0$
 $\Leftrightarrow \|Az\| \neq 0 \quad \forall z \neq 0 \Leftrightarrow (Az)^T (Az) = \|Az\|^2 \neq 0 \quad \forall z \neq 0 \Leftrightarrow$
 $\Leftrightarrow z^T A^T A z \neq 0 \quad \forall z \neq 0 \Leftrightarrow A^T A \text{ positive defined} \quad \square$

We assume that we have problems with full-col. rank

$$\min_{x \in \mathbb{R}^n} \|Ax - b\| \Leftrightarrow \min_{x \in \mathbb{R}^n} \frac{1}{2} \|Ax - b\|^2$$

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} (Ax - b)^T (Ax - b) = \min_{x \in \mathbb{R}^n} \frac{1}{2} (x^T A^T A x - x^T A^T b - b^T A x + b^T b)$$

equal

because $(x^T A^T b)^T = b^T A x$

$$\text{min}_{x \in \mathbb{R}^m} \frac{1}{2} x^T A^T A x - b^T A x + b^T b$$

$$\frac{1}{2} x^T Q x + q^T x + c \quad \text{quadratic function}$$

$$\nabla f = A^T A x - A^T b$$

$\nabla^2 f = A^T A$ positive definite matrix.

To find a minimum, set gradient to 0

$$x \text{ minimum} \iff A^T A x - A^T b = 0 \iff x = (A^T A)^{-1} A^T b$$

$$A^T A x = A^T b$$

the solution is $x = (A^T A)^{-1} A^T b$, whenever A has full column rank.

Remark: the solution is $x = A^+ b$, where $A^+ = (A^T A)^{-1} A^T$ is called pseudoinverse.

Computational cost of solving the problem with $x = (A^T A)^{-1} A^T b$:

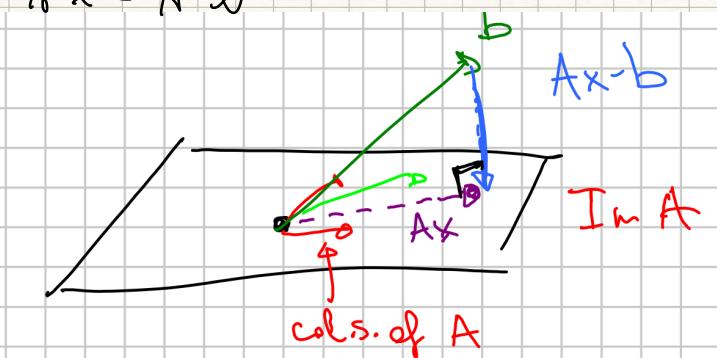
1 Compute $A^T \cdot A$ $\leftarrow 2mn^2$ flops
 $m \times n \quad m \times n$

2 Compute $A^T \cdot b$ $\leftarrow 2mn$
 $n \times m \quad n \times 1$

3 Solve $(A^T \cdot A) x = (A^T \cdot b)$ $\leftarrow \frac{2}{3} n^3$ (Gaussian elimination,

(LU feet)

In memory: if you can't solve $Ax = b$, multiply both sides by A^T : $A^T A x = A^T b$



$Ax - b$ is orthogonal to all vectors in $\text{Im } A$, in particular to the columns of A $A = [v_1 | v_2 | \dots | v_m]$

$$v_i^T (Ax - b) = 0 \quad \forall i$$

$$\begin{array}{c} \uparrow \\ \left[\begin{array}{c} v_1^T \\ \vdots \\ v_m^T \end{array} \right] \end{array} \left\{ (Ax - b) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \right. \quad A^T (Ax - b) = 0$$

Method of normal equations: $x = (A^T A)^{-1} (A^T b)$

Pseudoinverse: $A^+ = (A^T A)^{-1} A^T$ $A \in \mathbb{R}^{m \times n}$ full. column rank.

- Size $(A^+) = \text{size } (A^T)$ $m \geq n$
- If A orthogonal, $A^+ = A^T$, because $A^T A = I$
- $A^+ A = (A^T A)^{-1} A^T A = I$ $\boxed{} \boxed{} = \square$

$$\bullet A \cdot A^T \neq I$$

Lemme: for every vector $v \in \mathbb{R}^m$, the matrix

$$H = I_m - \frac{2}{v^T v} v v^T$$

$$v^T v = \|v\|^2 \text{ scalar prod.}$$

$$v v^T = \boxed{} \quad \boxed{} = \boxed{}$$

is orthogonal and symmetric.

Matrices of the form (*) are called Householder reflectors.

$$\text{Proof: } H^T = \left(I - \frac{2}{v^T v} v v^T \right)^T = I^T - \left(\frac{2}{v^T v} v v^T \right)^T = I - \frac{2}{v^T v} v v^T = H$$

$$H^T H = H^2 = \left(I - \frac{2}{v^T v} v v^T \right)^2 = I - I \frac{2}{v^T v} v v^T - \frac{2}{v^T v} v v^T \cdot I + \frac{4}{(v^T v)^2} v v^T v^T = I$$

$$\text{Equivalent form: } I - \frac{2}{\|v\|^2} v v^T = I - 2\mu\mu^T, \text{ where } \mu = \frac{1}{\|v\|} \cdot v$$

If H is a Householder reflector, and w is a vector, then we can compute Hw with cost $O(n)$ rather than $O(n^2)$.

$$Hw = (I - 2\mu\mu^T)w = w - 2\mu(\mu^T w)$$

1. compute scalar product $\lambda = \mu^T w \quad O(n) \quad 2n$
2. compute linear combination $w - \mu(\lambda\mu) \quad O(n) \quad 2n$

Lemma: let $x, y \in \mathbb{R}^n$ with $\|x\| = \|y\|$. If one chooses $v = x - y$, then $H = I_m - \frac{2}{\|v\|^2} v v^T$ is such that $Hx = y$

Th: for every $A \in \mathbb{R}^{m \times m}$, there exist $Q \in \mathbb{R}^{m \times m}$ orthogonal, $R \in \mathbb{R}^{m \times m}$ upper triangular such that $A = Q \cdot R$

upper triangular: $R_{ij} = 0$ if $i > j$

$$\begin{matrix} & \diagdown \\ 0 & \diagup \end{matrix}$$

$$\begin{matrix} & \diagdown \\ 0 & \diagup \end{matrix}$$

We have already shown this for $n=1$; given any $A \in \mathbb{R}^{m \times n}$ there is H Householder reflector s.t. $H \cdot A = \begin{bmatrix} S \\ 0 \\ \vdots \end{bmatrix}$

$$A = H \begin{bmatrix} S \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (H^{-1} = H^T = H, \text{ so } H^2 = I)$$

QR decomposition algorithm:

Step 1: take H_1 , Householder reflector such that

$$H_1 \cdot \begin{bmatrix} A_{1,1} \\ \vdots \\ A_{m,1} \end{bmatrix} = \begin{bmatrix} S \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

↑ first column of A

$$H_1 A = \begin{bmatrix} S_1 & * & \dots & * \\ 0 & \ddots & & \\ \vdots & & \ddots & \\ 0 & \ddots & \ddots & \ddots \end{bmatrix}$$

Step 2: Not $H_2 H_1 A = \begin{bmatrix} * & S_2 & * & \dots & * \\ \vdots & 0 & \vdots & \ddots & \vdots \\ * & 0 & * & \dots & * \end{bmatrix}$, as it would destroy zeros in the first column

so ...

$$Q_2 = \left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & H_2 \end{array} \right] \quad Q_2 A = \left[\begin{array}{c} v_1 \\ \vdots \\ v_m \end{array} \right] = \left[\begin{array}{c} v_1 \\ \hline H_2 \cdot \left[\begin{array}{c} v_2 \\ \vdots \\ v_m \end{array} \right] \end{array} \right] \quad Q_2 (H_2 A) = \left[\begin{array}{cccc} s_1 & * & \dots & * \\ 0 & s_2 & \dots & * \\ \vdots & 0 & \ddots & \vdots \\ 0 & 0 & \dots & s_n \end{array} \right]$$

$H_2 \in \mathbb{R}^{(m-1) \times (m-1)}$ where H_2 is chosen s.t. $(H_2 A)$ is mapped to $\begin{bmatrix} s_2 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

Step 3: $Q_1 = H_1$

$$Q_2 Q_1 A = \left[\begin{array}{ccccc} s_1 & x & x & \dots & x \\ 0 & s_2 & x & \dots & \vdots \\ \vdots & 0 & x & \dots & \vdots \\ 0 & 0 & x & \dots & x \end{array} \right] \quad H_3 \in \mathbb{R}^{(m-2) \times (m-2)}$$

that maps $(Q_2 Q_1 A)_3$ to $\begin{bmatrix} s_3 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

$$Q_3 = \left[\begin{array}{c|c} 1 & 0 \\ \hline 0 & H_3 \\ 0 & \vdots \\ 0 & H_3 \end{array} \right] \text{ w that}$$

$$Q_3 (Q_2 Q_1 A) = \left[\begin{array}{ccccc} s_1 & x & \dots & x & \} \text{ unchanged} \\ 0 & s_2 & x & \dots & x \\ \vdots & 0 & s_3 & \dots & x \\ 0 & 0 & 0 & \ddots & x \\ 0 & 0 & 0 & \ddots & x \end{array} \right] \} \text{ multiplied each column by } H_3$$

We can continue until the end, when (after n steps)

$$Q_m Q_{m-1} \dots Q_2 Q_1 A = \left[\begin{array}{ccccc} x & \dots & x & \dots & x \\ 0 & x & \dots & x & \dots \\ \vdots & 0 & \ddots & \ddots & x \\ 0 & 0 & \ddots & \ddots & x \end{array} \right] = R \text{ upper triangular}$$

orthogonal

$$Q_{k_i} = \left[\begin{array}{c|c} I_{k-1} & 0 \\ \hline 0 & R_k \\ k-1 & m-k+1 \end{array} \right]_{m-k+1}$$

$$\begin{bmatrix} \vdots & \cdots & 0 \\ 0 & \cdots & 1 \\ 0 & - & 0 \end{bmatrix}$$

Multiply each row on the left by $Q_1 Q_2 Q_3 \dots Q_{m-1}$,

use that $Q_k \cdot Q_{k_i} = I$ to get $Q_1 Q_2 \dots Q_{m-1} Q_m Q_{m-1} \dots Q_2 Q_1 A = Q_1 Q_2 \dots Q_{m-1} R$

□

orthogonal triangular

Problem : Cubic cost at each iteration $O(m^2 n)$



total cost is a fourth power $O(m^2 n^2)$

We need to use fast householder arithmetic:

$$HM = (I - 2uu^T)M = M - (2u)(u^T M) \leftarrow \text{high optimization}$$

~~PXP PXP~~

$O(p^2 q)$

$O(pq)$

$O(pq)$

bring the cost down

from quadratic to cubic

$$\text{Minor optimization: } H_i A = \begin{bmatrix} s_1 x & \cdots & x \\ 0 & \ddots & \vdots \\ \vdots & x & \cdots & x \end{bmatrix}$$

If we know that the first column of $H_{k_i} \cdot R(k_i: \text{end}, k_i: \text{end})$ is equal to $\begin{bmatrix} s_k \\ \vdots \\ 0 \end{bmatrix}$, we do not need to recompute it.

2 If $\mu = u$, then Q_m does very little, I can drop it.

for tell thin A

$$A = QR \leftarrow \begin{bmatrix} Q_0 & | & Q_C \\ \hline m & & m-m \end{bmatrix} \left[\begin{array}{c|c} R_0 & \\ \hline 0 & m-m \end{array} \right] = Q_0 R_0$$

$m \quad Q_0 \quad n$
 $n \quad \cancel{\square} \quad n$

So we can write $A = Q_0 R_0$ with thinner matrices $Q_0 \in \mathbb{R}^{m \times m}$
 $R_0 \in \mathbb{R}^{m \times m}$

thin QR factorization: returns only Q_0, R_0 , and has
 cheaper cost $O(mn^2)$. How to return only Q_0 without
 forming the whole Q ?

$$Q_0 = [\mathbb{I} - 2\mu_1 \mu_1^\top] \cdot \begin{bmatrix} 1 & 0 \\ 0 & \mathbb{I} - 2\mu_2 \mu_2^\top \end{bmatrix} \cdot \begin{bmatrix} \mathbb{I}_2 & \\ \hline & \mathbb{I} - 2\mu_3 \mu_3^\top \end{bmatrix} \cdots \begin{bmatrix} \mathbb{I}_{m-1} & \\ \hline & \mathbb{I} - 2\mu_m \mu_m^\top \end{bmatrix}$$

Can be represented only via $\mu_1, \mu_2, \dots, \mu_m \rightarrow O(mn)$

If you want to compute e.g. products $Q \cdot w, Q^\top w$, then
 you can use $\textcircled{*}$ plus hot Householder products to
 compute products with the same time complexity.

→ Don't return Q or Q_0 , just return $\mu_1, \mu_2, \dots, \mu_m$

• You can use this trick also to compute $Q_0 = Q \begin{bmatrix} I_m \\ 0 \end{bmatrix} m$

$$Q_0 = \left[I - 2\mu \mu_1^\top \right] \dots \left[\begin{array}{c|c} F_{m-1} & 0 \\ \hline 0 & I - 2\mu_m \mu_m^\top \end{array} \right] \begin{bmatrix} I_m \\ 0 \end{bmatrix}$$



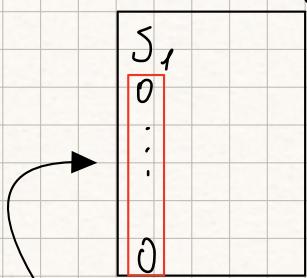
Compute multiplying to the right

Computational cost:

thin QR: $O(mn^2) \rightarrow mn^2 - \frac{2}{3}m^3 + O(mn)$

- square: $\frac{4}{3}m^3$ as twice as much as LU factorization
- full-thin $m \gg n \quad 2mn^2$.

Optimal storage:



$$\mu_1 \in \mathbb{R}^{m \times 1}$$

$$\mu_1 = \frac{v_1}{\|v_1\|}$$

up to normalization

recycle this position to store w_1

$$A = Q R = Q_0 R_0$$

$$Q = \begin{bmatrix} Q_0 & Q_C \end{bmatrix} \in \mathbb{R}^{m \times m} \quad \text{orthogonal}$$

$$R = \begin{bmatrix} R_0 \\ 0 \end{bmatrix} \in \mathbb{R}^{m \times m}$$

R_0 upper triangular

How to use the QR factorization to solve L.S. prob.

Find $x \in \mathbb{R}^{n \times m}$ that minimizes $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$

$$\|Ax - b\| = \|Q_R x - b\| = \|Q^T(Q_R x - b)\| \geq \|Rx - Q^T b\| =$$

Orth. matrix preserve the 2-norm

$$= \left\| \begin{bmatrix} R_0 \\ 0 \end{bmatrix} x - \begin{bmatrix} Q_0^T \\ Q_1^T \end{bmatrix} b \right\| = \left\| \begin{bmatrix} R_0 x \\ 0 \end{bmatrix} - \begin{bmatrix} Q_0^T b \\ Q_1^T b \end{bmatrix} \right\| = \left\| \begin{bmatrix} R_0 x - Q_0^T b \\ Q_1^T b \end{bmatrix} \right\|$$

We make the norm as small as possible choosing x as:

Second block: does not contain x , its entries are always going to be in the sum of squares.

First block: becomes 0 if x solves $R_0 x = Q_0^T b$ (square linear system with matrix R_0). If R_0 is invertible, this is always possible.

min $\|Ax - b\|$ equals to $\|Q_1^T b\|$, and is achieved by

$$x = R_0^{-1} Q_0^T b$$

When R_0 is invertible \rightarrow Lemma: R_0 is invertible i.f.f. A has full column rank.

Proof: A has full col. rank $\iff A^T A$ is invertible

$\leftrightarrow R^T Q^T QR$ is invertible $\leftrightarrow [R_0^T \ 0] \begin{bmatrix} R_0 \\ 0 \end{bmatrix} = R_0^T R_0$ is invertible $\leftrightarrow R_0$ is invertible.

Second algorithm to solve O.L.S. prob. $\min \|Ax - b\|_2$,

with $A \in \mathbb{R}^{m \times n}$ with full column rank:

1 Compute thin QR factorization of A $\mathcal{O}(mn^2)$ if $m \gg n$

2 form $C = Q_0^T b$ $\mathcal{O}(mn)$

3 Solve lin. system $R_0 x = C$ $\mathcal{O}(n^2)$

$$\begin{bmatrix} R_{11} & \dots & R_{1n} \\ 0 & \ddots & \vdots \\ 0 & \dots & R_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} c_1 \\ \vdots \\ c_n \end{bmatrix}$$

Solving least squares problem with SVD

$$\min_{x \in \mathbb{R}^m} \|Ax - b\|_2, \quad A \in \mathbb{R}^{m \times n}, \quad m \geq n, \quad b \in \mathbb{R}^m$$

$$A = U S V^T$$

$$U \in \mathbb{R}^{m \times m}$$

$$U = \begin{bmatrix} U_0 & U_{0c} \end{bmatrix}_{m \times m}$$

$$S = \begin{bmatrix} S_0 \\ 0 \end{bmatrix}_{m \times n}$$

$$S = \begin{pmatrix} * & & \\ & \ddots & \\ & & * \end{pmatrix} \in \mathbb{R}^{m \times n}$$

$$V = R^{m \times n}$$

$$A = U_0 S_0 V^T$$

$$V^T x = y$$

$$\|Ax - b\| = \|USV^T x - b\| = \|U^T(U SV^T x - b)\| = \|S\bar{y} - U^T b\| =$$

$$= \left\| \begin{bmatrix} 6_1 & \dots & 6_n \\ 0 & \ddots & 0 \\ \vdots & & \vdots \\ 0 & & 0 \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} - \begin{bmatrix} u_1^T b \\ \vdots \\ u_m^T b \end{bmatrix} \right\| = \left\| \begin{bmatrix} 6_1 y_1 \\ \vdots \\ 6_m y_m \\ 0 \end{bmatrix} - \begin{bmatrix} u_1^T b \\ \vdots \\ u_m^T b \end{bmatrix} \right\|$$

We can choose y s.t. the first m rows are equal to 0

$$y_i := \frac{u_i^T b}{6_i}$$

$$x = V y = \begin{bmatrix} v_1 & | & v_2 & | & \dots & | & v_m \end{bmatrix} \cdot \begin{bmatrix} \frac{u_1^T b}{6_1} \\ \vdots \\ \frac{u_m^T b}{6_m} \end{bmatrix} = v_1 \frac{u_1^T b}{6_1} + v_2 \frac{u_2^T b}{6_2} + \dots + v_m \frac{u_m^T b}{6_m}$$

well regular
value count more

$$x = V \begin{bmatrix} 6_1 & 0 & \dots & 0 \\ 0 & \ddots & & 0 \\ & & \ddots & 0 \\ & & & 6_m \end{bmatrix} \begin{bmatrix} u_1^T b \\ \vdots \\ u_m^T b \end{bmatrix} = \underbrace{VS_0^{-1}U_0^T}_A b$$

The value of the minimum is $\left\| \begin{bmatrix} 0 \\ \vdots \\ u_{m+1}^T b \\ \vdots \\ u_m^T b \end{bmatrix} \right\| = \|U_0^T b\|$

$6_i \neq 0$ for all $i = 1, 2, \dots, n$ i.f.f. A has full column rank. Indeed,

A full col. rank $\leftrightarrow A^T A$ is well regular

$$A^+ A = (U S V^T)(U S V^T) = U S^T V^T U S V^T = V \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_m^2 \end{bmatrix} V^T$$

which is invertible i.f.f. $\sigma_i \neq 0 \forall i$.

Also SVD reveal how far A is from a matrix without full column rank.

X doesn't have a full col. rank \leftrightarrow rank X is $n-1$ or lower

$\min \|A - x\| = \sigma_m$ the minimum value is obtained with truncated SVD.

σ_m is the smallest singular value, w A has full col. rank
 $\leftrightarrow \sigma_m \neq 0$.

With SVD we can solve problem in which A doesn't have a full rank. $\|Ax - b\| = \sqrt{\left(\begin{array}{c} \sigma_1 y_1 - u_1^T b \\ \vdots \\ \sigma_m y_m - u_m^T b \end{array} \right)^T \left(\begin{array}{c} \sigma_1 y_1 - u_1^T b \\ \vdots \\ \sigma_m y_m - u_m^T b \end{array} \right)}$

$$\|Ax - b\| = \sqrt{\left(\begin{array}{c} \sigma_1 y_1 - u_1^T b \\ \vdots \\ \sigma_m y_m - u_m^T b \\ 0 \\ \vdots \\ 0 \end{array} \right)^T \left(\begin{array}{c} \sigma_1 y_1 - u_1^T b \\ \vdots \\ \sigma_m y_m - u_m^T b \\ 0 \\ \vdots \\ 0 \end{array} \right)}$$

If $\sigma_{n+1} = \sigma_{n+2} = \dots = \sigma_m = 0$ i.e. the rank of A is n , the $y_{n+1} \dots y_m$ do not effect the value of the norm. We can choose them arbitrarily
 $\rightarrow \infty$ solutions

Among the infinite solutions, one can identify the one with the smallest norm.

$$\min_{x \text{ solves}} \|x\| = \min \|Vx\| = \min \|y\|$$

the optimum is

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} \frac{u_1^T b}{6_1} \\ \frac{u_2^T b}{6_2} \\ \vdots \\ 0 \end{bmatrix}$$

} forced so that y is a solution
} to get minimum norm

Solution of rank-deficient least square problem

$$\begin{aligned} x &= V y = v_1 \frac{u_1^T b}{6_1} + \dots + v_m \frac{u_m^T b}{6_m} \\ &= \begin{bmatrix} v_1 & | & v_2 & | & \dots & | & v_m \end{bmatrix} \begin{bmatrix} 1/6_1 & & & & & & \\ & \ddots & & & & & \\ & & 1/6_m & & & & \\ & & & \ddots & & & \\ & & & & \ddots & & \\ & & & & & \ddots & \\ & & & & & & 1/6_m \end{bmatrix} \begin{bmatrix} u_1^T b \\ \vdots \\ u_m^T b \end{bmatrix} \\ &= V \begin{bmatrix} 1/6_1 & & & & & & \\ & \ddots & & & & & \\ & & 1/6_m & & & & \\ & & & \ddots & & & \\ & & & & \ddots & & \\ & & & & & \ddots & \\ & & & & & & 0 \end{bmatrix} U_0^T b \end{aligned}$$

A^+

Def: pseudoinverse of a matrix A without full. col. rank is

$$A^+ = V \begin{bmatrix} 1/6_1 & & & & & & \\ & \ddots & & & & & \\ & & 1/6_m & & & & \\ & & & \ddots & & & \\ & & & & \ddots & & \\ & & & & & \ddots & \\ & & & & & & 0 \end{bmatrix} U_0^T$$

with this definition, the vector $x = A^+ b$ is always a

solution of $\min \|Ax - b\|$, and if there is an infinite number of them this formula gives the one with the smallest norm.

Truncating: in SVD many singular values (and their u_i, v_i) count less, as they are responsible for smaller changes in norm.

Ridge regression is alternative to truncation

$$\min_{x} \|Ax - b\|^2 + \lambda^2 \|x\|^2 \quad \lambda > 0$$

Trick:

$$\|Ax - b\|^2 + \lambda^2 \|x\|^2 = \left\| \begin{bmatrix} A \\ \lambda I_m \end{bmatrix}_m \cdot x - \begin{bmatrix} b \\ 0 \end{bmatrix}_m \right\|^2 = \left\| \begin{bmatrix} Ax - b \\ \lambda x \end{bmatrix} \right\|^2$$

solution of the ridge-regression problem:

$$x_{\text{ridge}} = \begin{bmatrix} A \\ \lambda I_m \end{bmatrix}^+ \begin{bmatrix} b \\ 0 \end{bmatrix} = \left(\begin{bmatrix} A^\top & \lambda I_m \end{bmatrix} \begin{bmatrix} A \\ \lambda I_m \end{bmatrix} \right)^{-1} \begin{bmatrix} A^\top & \lambda I_m \end{bmatrix} \begin{bmatrix} b \\ 0 \end{bmatrix}$$

$$= (A^\top A + \lambda^2 I)^{-1} (A^\top b)$$

additional term with respect to $x = (A^\top A)^{-1} A^\top b$
always positive definite

Alternative expression with SVD $A = U \Sigma V^\top$

$$x_{\text{ridge}} = V \begin{bmatrix} \frac{6_1}{6_i^2 + 2^2} & \dots & \frac{6_m}{6_m^2 + 2^2} \end{bmatrix} U^T b$$

The quantity $\frac{1}{6_i}$, which appears in the analogous expression for x is replaced by $\frac{6_i}{6_i^2 + 2^2}$

$$\frac{6_i}{6_i^2 + 2^2} \begin{cases} \frac{1}{6_i} & \text{if } 6_i >> 2 \\ 0 & \text{if } 6_i \ll 2 \end{cases}$$

This approximation truncating singular values that are smaller than ≈ 2

Recap: 3 algorithms for LS problems:

Normal equations

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

$$m > n \quad \approx m n^2$$

$$m \leq n \quad \approx 4/3 m^3$$

QR

$$x = Q_s^{-1} Q_s^T b$$

$$\approx 2 m n^2$$

$$\approx 4/3 m^3$$

SVD

$$x = \underbrace{U S_0^{-1} U^T}_{\text{Mseudoinverse } A^+} b$$

$$2 m n^2$$

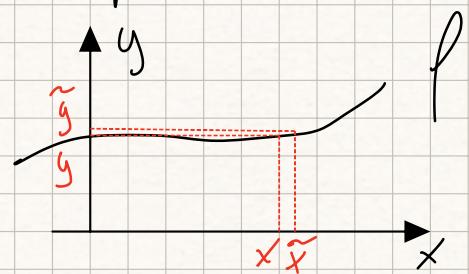
$$1/3 m^3$$

Condition number: a way to measure the sensitivity of a problem for small perturbations of the inputs

given $f: \mathbb{R} \rightarrow \mathbb{R}$

$$y = f(x)$$

$$\tilde{y} = f(\tilde{x})$$

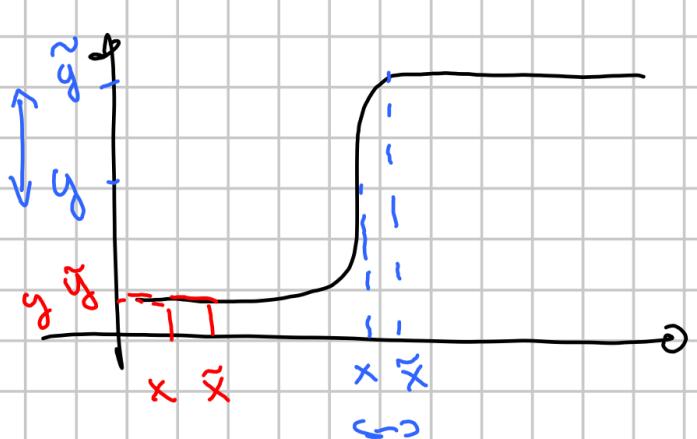


$$f(\tilde{x}) = f(x) + f'(x)(\tilde{x} - x) + O((\tilde{x} - x)^2)$$

$\underbrace{f(\tilde{x})}_{\tilde{y}}$ $\underbrace{f(x)}_{y}$

Relative error difference:

$$\frac{|\tilde{y} - y|}{|y|} = \frac{|f'(\tilde{x})(\tilde{x} - x)|}{|f(x)|} + O((\tilde{x} - x)^2) = \frac{|f'(\tilde{x})| \cdot |\tilde{x} - x|}{|f(x)|} \cdot \frac{|\tilde{x} - x|}{|\tilde{x}|} + O(\cdot)$$

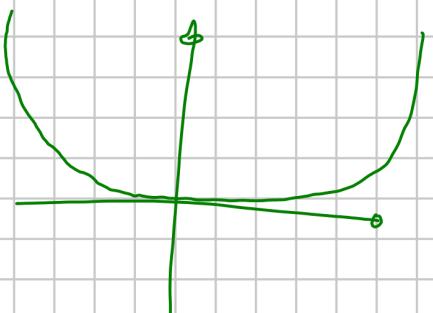
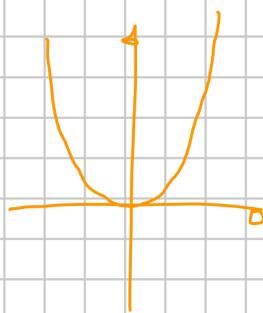


condition number
of f in point x rel. error
on input

Formal definition: condition number of $f: \mathbb{R} \rightarrow \mathbb{R}$

$$\lim_{|\tilde{x} - x| \rightarrow 0} \frac{|f(\tilde{x}) - f(x)|}{|f(x)|} \cdot \frac{|\tilde{x} - x|}{|x|}$$

For function with more inputs or outputs, the answer may also depend on the direction:



We can define a maximum over all directions

Def:

$$K f, x = \limsup_{\|x^2 - x\| \rightarrow 0}$$

$$\frac{\|f(x^2) - f(x)\|}{\|f(x)\|}$$

$$\frac{\|f(x^2) - f(x)\|}{\|x^2 - x\|}$$

$$\frac{\|f(x^2) - f(x)\|}{\|x\|}$$

In terms of the Jacobian (matrix of derivatives)

$$K f, x = \frac{\|\nabla f(x)\| \cdot \|x\|}{\|f(x)\|}$$

The condition number is a warning sign: it tells you whether your problem is very sensitive to input perturbations.

Conditioning of linear algebra problems:

solving linear systems $A \in \mathbb{R}^{n \times n}$ invertible, $b \in \mathbb{R}^n$, $Ax = b$

Let \tilde{b} a perturbation of b

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

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa \frac{\|\tilde{b} - b\|}{\|b\|}$$

$Ax = b$ exact output

Theorem: this condition number is $\kappa(A) = \|A\| \cdot \|A^{-1}\|$

$$\text{Proof: } \|\tilde{x} - x\| = \|A^{-1}\tilde{b} - A^{-1}b\| = \|A^{-1}(A^{-1}\tilde{b} - A^{-1}b)\| \leq \|A^{-1}\| \cdot \|\tilde{b} - b\|$$

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

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \frac{\|A^{-1}\| \|\tilde{b} - b\|}{\|b\|} = \boxed{\|A^{-1}\| \cdot \|A\|} \cdot \frac{\|\tilde{b} - b\|}{\|b\|}$$

conditional number

of solving linear equations w.r.t. b

Def: $\kappa(A) = \|A\| \cdot \|A^{-1}\|$ is called the condition number of A

Condition number and singular value

$$\|A\| = 6_1$$

$$\|A^{-1}\| = \|(U S V^T)^{-1}\| = \|(V^T)^{-1} S^{-1} U^{-1}\| = \|V \begin{bmatrix} 1/\sigma_1 & & \\ & \ddots & \\ & & 1/\sigma_n \end{bmatrix} U^T\|$$

Based on SVD

$$= \left\| \begin{bmatrix} v_m & \dots & v_1 \end{bmatrix} \begin{bmatrix} 1/\sigma_n & & \\ & \ddots & \\ & & 1/\sigma_1 \end{bmatrix} \begin{bmatrix} u_1 & \dots & u_m \end{bmatrix}^T \right\| = \frac{1}{\sigma_m}$$

SVD of A^{-1}

$$k(A) = \|A\| \cdot \|A^{-1}\| = \frac{6_1}{6_m} \quad \text{formula for } k(A) \text{ in terms of SVD}$$

Condition number and distance to singularity:

$$\min_{x \text{ singular}} \|A - x\| = 6_m$$

$$\min_{x \text{ singular}} \frac{\|A - x\|}{\|A\|} = \frac{6_m}{6_1} = \frac{1}{k(A)}$$

Condition number of solving linear systems w.r.t. perturbations of A : $Ax = b$ $\tilde{A}\tilde{x} = \tilde{b}$

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \underbrace{\|A\| \cdot \|A^{-1}\|}_{k(A)} \cdot \frac{\|\tilde{A} - A\|}{\|A\|} + O(\|\tilde{A} - A\|^2)$$

Condition number of least square problems.

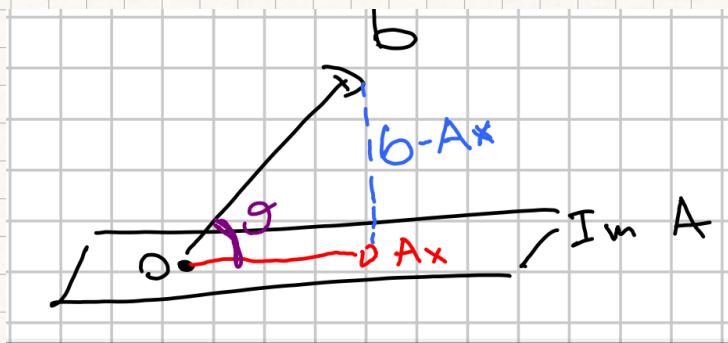
Consider the L.S. problem $\min \|Ax - b\|$, $A \in \mathbb{R}^{m \times n}$ full column rank $b \in \mathbb{R}^m$. The condition number w.r.t. input b is:

$$k_{b \rightarrow x} \leq \frac{k(A)}{\cos \theta}$$

and w.r.t. input A

$$K_A \rightarrow x \leq k(A) + k(A)^2 \tan \theta$$

where $k(A) = \frac{6_1}{6_n}$, θ is the angle s.t. $\cos \theta = \frac{\|A\|}{\|B\|}$



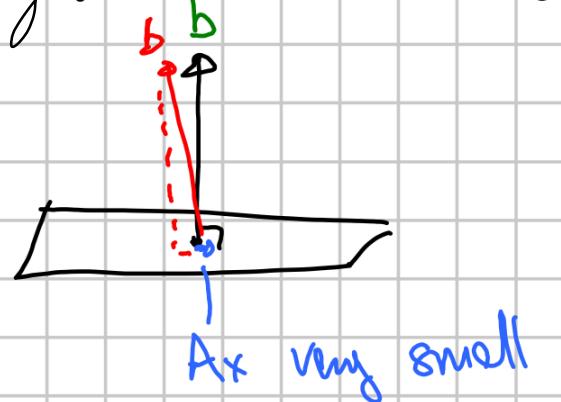
θ = angle between $\mathbb{L}_m A$ and b

$$\|b - Ax\| = \|Q_c^\top b\| = \|U_c^\top b\|$$

$$\|Ax\| = \|Q_0^\top b\| = \|U_0^\top b\|$$

interesting case:

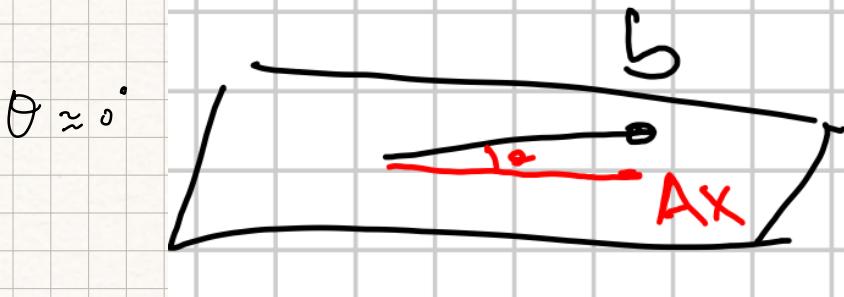
$$\theta \approx 90^\circ$$



$$K_B \rightarrow x = \frac{k(A)}{\cos \theta} \quad \text{very large}$$

$$K_A \rightarrow x = k(A) + \tan \theta \cdot k(A)^2 \quad \text{very large}$$

small change in b or A \rightarrow large change in x



$$b - Ax \quad \text{very small}$$

$$K_B \rightarrow x = K(A)$$

$$K_A \rightarrow x = K(A)$$

behave like a linear system, $k(A)$

$$0^\circ < \theta < 30^\circ$$

$$K_B \rightarrow x = \frac{k(A)}{\cos(\theta)}$$

$$K_A \rightarrow x = k(A) + \tan \theta \cdot k(A)^2$$

more sensitive to perturbation

of A then a linear system

Note that $\frac{1}{\kappa(A)} \geq \text{rel. distance between } A \text{ and the closest matrix}$

without full column rank.

Sensitivity: how does the solution change w.r.t. small input changes → condition number

Stability: how well does an algorithm compute the solution to the problem

Stability of numerical algorithms

If your input is not exact floating point numbers, then you cannot compute the exact solution.

Real input x

F.P. approximation \tilde{x}

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \mu \approx 10^{-16}$$

or $\frac{\|\tilde{x} - x\|}{\|x\|} \leq O(\mu, n) \mu \text{ for vectors / matrices}$

You can't compute $f(x)$, only $f(\tilde{x})$

$$\frac{\|f(\tilde{x}) - f(x)\|}{\|f(x)\|} \leq \kappa_{f,x} \frac{\|\tilde{x} - x\|}{\|x\|} \quad \text{"intrinsic error"}$$

To assess error of an algorithm, the only recourse is keeping track of each floating point error

$$a + b \longrightarrow (a+b)(1+\epsilon) \quad |\epsilon| \leq u$$

machine precision

Trick: backward stability analysis; for some algorithms you can write the result \tilde{y} as an exact $f(\tilde{a}, \tilde{b})$ with \tilde{a}, \tilde{b} "close" to a, b .

$$\tilde{y} = \underbrace{a_1 b_1 + a_2 b_2 + a_3 b_3}_y + a_1 b_1 (\epsilon_1 + \epsilon_4 + \epsilon_5) + a_2 b_2 (\epsilon_2 + \epsilon_4 + \epsilon_5) + a_3 b_3 (\epsilon_3 + \epsilon_5)$$

$$= a_1 \hat{b}_1 + a_2 \hat{b}_2 + a_3 \hat{b}_3$$

$$\hat{b}_1 = b_1 (1 + \epsilon_1 + \epsilon_4 + \epsilon_5)$$

$$\hat{b}_2 = b_2 (1 + \epsilon_2 + \epsilon_4 + \epsilon_5)$$

$$\hat{b}_3 = b_3 (1 + \epsilon_3 + \epsilon_5)$$

$$\|\hat{b} - b\|_2 = \left\| \begin{bmatrix} b_1 (\epsilon_1 + \epsilon_4 + \epsilon_5) \\ b_2 (\epsilon_2 + \epsilon_4 + \epsilon_5) \\ b_3 (\epsilon_3 + \epsilon_5) \end{bmatrix} \right\|_2 \leq 3u \|b\|_2$$

This tells us immediately an accuracy bound on y :

$$\frac{|\tilde{y} - y|}{|y|} \leq u_{f,b} \frac{\|\hat{b} - b\|}{\|b\|} = u_{f,b} \cdot 3u$$

\Rightarrow the result obtained is just a vector 3 away from the (unavoidable) error coming from machine-precision approximation of u

$$\frac{|\tilde{y} - y|}{|y|} \leq u_{f,b} \frac{\|\hat{b} - b\|}{\|b\|}$$

$\underbrace{}_u$

We say that an algorithm (or a computation) is backward stable if the computed output $\tilde{y} = f(\tilde{x})$ is exactly equal to $f(x)$, where $\frac{\|x' - x\|}{\|x\|} = O(n)$. In this case, the algorithm is as accurate as it could get given the intrinsic error.

Algorithmic error \approx intrinsic error

⚠️ This $O(n)$ can include small-degree polynomials of the dimensions, because typically these are overestimates, and because they are typically much smaller than possible condition numbers.

This trick doesn't work for every algorithm,

Sometimes it's possible to prove stability "a posteriori": after you have computed a solution, you can show it solves a nearby problem.

ex: residual test for square linear systems

$A \in \mathbb{R}^{n \times n}$, $b \in \mathbb{R}^n$ find $x \in \mathbb{R}^n$ that solves $Ax = b$

Suppose your algorithm computes \tilde{x} (approx. solution)

residual: $r = A\tilde{x} - b$

Th: the error satisfies

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa(A) \cdot \frac{\|r\|}{\|b\|}$$

↓
 $\|A\| \cdot \|A^{-1}\|$

Proof: $r = A\tilde{x} - b \iff A\tilde{x} = \underbrace{b + r}_{\hat{b}}$

\tilde{x} is the exact solution of modified inputs A, \hat{b} . Then, the cond. number bound tells us that:

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa(A) \cdot \frac{\|\hat{b} - b\|}{\|b\|} = \kappa(A) \cdot \frac{\|r\|}{\|b\|}$$

→ solution with small residual are accurate (up to $\kappa(A)$)

problem $a = Ax - b$ is effected by an error itself. Then error is bounded by $\|A\| \cdot \|x\| \cdot O(\mu)$

→ a is typically accurate at least as order of magnitude.

Conversely, backward stable solution → small residual:

if \tilde{x} solves $A\tilde{x} = \hat{b}$ with $\frac{\|\hat{b} - b\|}{\|b\|} = O(\mu)$

then $A\tilde{x} - b = \hat{b} - b$ small w.r.t. $\|b\|$

$$\tilde{A} = A + \epsilon \quad \frac{\|\epsilon\|}{\|A\|} = O(n)$$

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

$$\tilde{b} = b + f \quad \frac{\|f\|}{\|b\|} = O(n)$$

$$\text{then } (A + \epsilon)\tilde{x} = b + f \iff A\tilde{x} - b = f - \epsilon\tilde{x}$$

$$\|A\tilde{x} - b\| \leq \|f\| + \|\epsilon\| \cdot \|\tilde{x}\| = O(n) \cdot (\|f\| + \|A\| \cdot \|\tilde{x}\|)$$

QR factorization: $A = QR$ if $\|A - \tilde{Q}\tilde{R}\|$ is small, then you have computed an exact solution of $qr(A + \epsilon)$ with $\tilde{\epsilon} = \tilde{Q}\tilde{R} - A$, hence

$$\frac{\|\tilde{R} - R\|}{\|R\|} \leq K_{qr,R} \cdot \frac{\|\epsilon\|}{\|A\|}$$

Small residual implies

Algorithm as accurate as possible

Residual tests for LS problems: $r = A\tilde{x} - b$ is not necessarily small

$$\min \frac{1}{2} \|A\tilde{x} - b\|^2 =$$

$$\min \frac{1}{2} \tilde{x}^T A^T A \tilde{x} - b^T A \tilde{x} + \text{const} \rightarrow \text{small gradient}$$

$$\nabla_f(\tilde{x}) = A^T A \tilde{x} - A^T b \quad (\text{residual of normal equations})$$

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

Residual bound on normal equations:

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa(A^T A) \cdot \frac{\|A^T A \tilde{x} - A^T b\|}{\|A^T b\|}$$

$$A^T A = V \begin{bmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \ddots & \\ & & & \sigma_n^2 \end{bmatrix} V^T \quad \leftarrow \text{SVD of } A^T A \quad \kappa(A^T A) = \frac{\sigma_1^2}{\sigma_n^2} = \kappa(A)^2$$

Even if $\frac{\|A^T A \tilde{x} - A^T b\|}{\|A^T b\|} = O(\mu)$, we get an error of the order of

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa(A^T A) \cdot O(\mu), \text{ but } \kappa_{LS, A} \text{ and } \kappa_{LS, B} \text{ can be better numbers}$$

then $\kappa(A^T A) \approx \kappa(A)^2$

Th: Suppose $A = Q_0 R_0$ is a thin QR factorization exactly computed!

Let $r_0 = Q_0^T (A \tilde{x} - b)$. Then \tilde{x} is the exact solution of

$$\min \|A \tilde{x} - (b + Q_0 r_0)\|.$$

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa_{LS, B} \cdot \frac{\|Q_0 r_0\|}{\|b\|} = \kappa_{LS, B} \cdot \frac{\|R_0 r_0\|}{\|b\|}$$

↑ orthogonal

true cond. number of the problem

For good (backward stable) algorithms, relative residuals \approx machine precision and this let you prove that

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq \kappa_{\text{prob.}} O(\mu)$$

best one can do because of intrinsic error.

"*a priori*" backward stability for alg. that are based on orthogonal transformations

Stability of matrix multiplication

$$y = \mathbf{c}^T b \quad \tilde{y} = \tilde{\mathbf{c}}^T \tilde{\mathbf{b}}$$

$$|\tilde{y} - y| \leq O(n) \cdot |\mathbf{c}| \cdot |b|$$

$$C = A \cdot B \quad A, B, C \in \mathbb{R}^{n \times n} \quad \tilde{C} = \tilde{A} \odot \tilde{B}$$

$$|\tilde{C}_{ij} - C_{ij}| \leq O(n) \cdot |A_{i:}| \cdot |B_{:j}|$$

This gives rise to a normwise bound:

$$\|\tilde{C} - C\| \leq O(n) \|A\| \cdot \|B\|$$

↑
depends on the chosen norm

If one of the matrices is orthogonal, then:

$$\|\tilde{C} - C\| \leq O(n) \|A\| \cdot \|B\| \leq O(n) \cdot \|B\|$$

$$\tilde{C} = C + \tilde{\epsilon} \quad \|\tilde{\epsilon}\| \leq O(n) \cdot \|B\| \quad \leftarrow \text{forward stability bound}$$

$$\|C\| = \|A \cdot B\| = \|B\| \text{ if } A \text{ orthogonal}$$

$$\frac{\|\Gamma\|}{\|C\|} \leq O(n)$$

Backward stability of:

$$\tilde{C} = C + \tilde{\Theta} = A \tilde{B}$$

$$\tilde{B} = B + \tilde{F}$$

$$\frac{\|\tilde{F}\|}{\|B\|} = O(n)$$

$$\tilde{C} = C + \tilde{\Theta} = AB + \tilde{\Theta} = A \left(B + \underbrace{A^T \tilde{\Theta}}_F \right)$$

$$\|\tilde{F}\| = \|A^T \tilde{\Theta}\| = (\|\tilde{\Theta}\| \leq O(n) \cdot \|B\|)$$

↑
orthogonal

What we proved: multiplication with the standard alg., by an orthogonal matrix is backward stable.

Remark: if A were not orthogonal, the same computation gives

$$\tilde{C} = A \odot B = A(B + \tilde{F}) \quad \|\tilde{F}\| \leq O(n) \cdot \|B\| \cdot \kappa(A)$$

If H is a Householder reflector, then the product is backward stable

$$\tilde{C} = H \odot B = H(B + \tilde{F})$$

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

$$\|\tilde{F}\| = O(n) \|B\|$$

$$HB = B - 2u(u^T B)$$

Proof via $\tilde{C} = C + \tilde{\Theta}$, $(\|\tilde{\Theta}\| = \|B\| O(n))$ by keeping track of $(1+\epsilon_i)$ errors

Th: $[\tilde{Q} \quad \tilde{R}] = \varphi_2(A)$ computed via Bleuselwolker reflectors returning matrices such that $\tilde{Q} \tilde{R} = A + F$, $\|F\| = O(n) \|A\|$.
 (\tilde{Q} represented implicitly via reflectors v_1, u_2, \dots, v_m)

Proof: exact qr steps

$$\underbrace{\begin{bmatrix} I_{v_{k-1}} \\ H(u_k) \end{bmatrix}}_{Q_k} \underbrace{\begin{bmatrix} x & \cdots & \cdots & x \\ 0 & \ddots & & \vdots \\ \vdots & & x & \ddots \\ 0 & 0 & x & \ddots & x \end{bmatrix}}_{V_{k-1} \quad A_{k-1}} = \underbrace{\begin{bmatrix} x & \cdots & \cdots & x \\ 0 & x & \cdots & \vdots \\ 0 & x & \ddots & \vdots \\ 0 & 0 & 0 & x \end{bmatrix}}_{A_k}$$

$$\underbrace{\begin{bmatrix} I_{v_{k-1}} \\ H(\tilde{u}_k) \end{bmatrix}}_{\tilde{Q}_k} \circ \tilde{A}_{k-1} = \tilde{A}_k \quad \tilde{Q}_k \text{ is exactly orthogonal}$$

$$\tilde{A}_k = \tilde{Q}_k \tilde{A}_{k-1} + \tilde{F} = \tilde{Q}_{v_k} (\tilde{A}_{k-1} + \tilde{F}_{v_k})$$

$$\|\tilde{F}_{k-1}\| = O(n) \|\tilde{A}_{k-1}\| = O(n) \|A\|$$

$$\tilde{A}_1 = \tilde{Q}_1 (A + F_0)$$

$$\tilde{A}_2 = \tilde{Q}_2 (\tilde{A}_1 + F_1) = \tilde{Q}_2 \tilde{Q}_1 (A + F_0 + \tilde{Q}_1 F_1)$$

$$\tilde{A}_3 = \tilde{Q}_3 (\tilde{A}_2 + \tilde{F}_2) = \tilde{Q}_3 \tilde{Q}_2 \tilde{Q}_1 (A + F_2 + \tilde{Q}_1^T F + \tilde{Q}_2^T \tilde{Q}_1 F_2)$$

All errors written as multiplication
of the matrix A

$$\tilde{R} = \tilde{Q}_m \tilde{Q}_{m-1} \dots \tilde{Q}_1 (A + F_2 + \dots + \tilde{Q}_{m-1} \dots Q_1 F_{m-1})$$

m errors terms, each
of norm $O(n) \|A\|$

$$\tilde{Q} = \tilde{Q}_m \dots \tilde{Q}_1 (A + F) \quad \|F\| (= O(n) \|A\|)$$

→ exact computed \tilde{R} (and \tilde{x}) are the exact result of
 $qr(\hat{A}) \quad \hat{A} \approx A + F \quad \|F\| = O(n) \|A\| \quad \blacksquare$

In addition, one can also prove that backward substitution

$$\begin{bmatrix} x & \vdots & x \\ 0 & \ddots & \vdots \\ \vdots & \ddots & x \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

is backward stable: \tilde{x} satisfies $\underbrace{(R+F)x}_{\tilde{R}} = b$

All steps to compute $x = qr$ are backward stable

QR computes \tilde{x} that solves exactly $\tilde{x} = \min_{\tilde{x}} \|\hat{A}\tilde{x} - \hat{b}\|$ and will always produce small $e_0 = \Omega^T (A\tilde{x} - b)$, and is as accurate

as it could be $\frac{\|(\tilde{x} - x)\|}{\|x\|} = O(u) \cdot (K_{CS,A} + K_{CS,B}).$

The new results holds for $x \rightarrow \text{null}$:

All have backward errors of order $\|A\| = \|S\|_1 \text{ times } O(u)$

$\rightarrow \tilde{x}$ solves $\min_x \|A^T \tilde{x} - b\|$ exactly

This does not hold for normal equations:

$$M = A^T A \quad \text{non-backward stable steps} \rightarrow \text{large error}$$

$$c = A^T b \quad \text{the error } \frac{\|\tilde{x} - x\|}{\|x\|} \text{ is always going to be of}$$

$$x_{\text{null}} = M^{-1} c \quad \text{magnitude at least } K(1) O(u) = k(A)^2 O(u)$$

$k(A)^2$ can be much larger than $K_{CS,A} + K_{CS,B}$

	NE	QR	SVD
$M \approx n$	$\frac{4}{3}n^3$	$\frac{4}{3}n^3$	$\approx 13n^3$
$M >> n$	Mn^2	$2Mn^2$	$2Mn^2$
	\uparrow	\uparrow	\uparrow

unstable, esp.
when Θ is small

backward
stable

gives info on
distance from
instability,
works better on
almost-singular problems

Matlab:

$$x = A \setminus b \text{ uses QR.}$$

LU factorization: it's a factorization

$$A = L U$$

$A \in \mathbb{R}^{n \times n}$

lower triangular
with 1 on the diagonal

upper triangular

To be: follow the plan of QR factorization, but use lower triangular matrices instead of Householder reflector.

This produces essentially Gaussian elimination:

$$A = \begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \end{bmatrix}$$

$$\text{row } i = \text{row } i - \frac{A_{i1}}{A_{11}} (\text{row } 1)$$

eliminate entries by
subtracting $\text{row } 1$ to $\text{row } i = 2, 3, \dots, n$

a multiple of $\text{row } 1$

$$A_2 = \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -a_2 & 1 & 0 \\ -a_3 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} x & x & x \\ x & x & x \\ x & x & x \end{bmatrix}$$

$$l_{ij} = \frac{A_{ij}}{A_{ii}}$$

multiplying on the left by $L_1 \leftrightarrow$ one step of G.E.

second step of G.E.

$$\begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -b_3 & 1 \end{bmatrix} \cdot \begin{bmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \end{bmatrix}$$

$$l_{ij} = \frac{(A_2)_{ij}}{(A_2)_{22}} \quad A_2$$

Successive steps : multiplying on the left by matrices of the form

$$L_k = \begin{bmatrix} 1 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}.$$

non zero entries only in the k-th column

$$U = L_3 A_3 = L_3 L_2 A_2 = \underbrace{L_3 L_2 L_1}_{} A$$

lower tr. with ones on diagonal
invertible

$$L_1^{-1} L_2^{-1} L_3^{-1} U = A$$

lower triangular upper triangular
matrix

Update at each step:

$$U_{k+1:n} \cdot l_{k+1:n} = U_{k+1:n} - \begin{bmatrix} L_{k+1,k} \\ \vdots \\ L_{m,k} \end{bmatrix} \cdot [U_k, U_{k+1}, \dots, U_{k,m}]$$

$(m-k) \times (m-k)$
 $(m-k)^2$
subtractions

$(m-k)^2$ products
rank-1 matrix

$$\text{Total cost: } 2(m-1)^2 + 2(m-2)^2 + \dots + 2 \cdot 2^2 + 2 \cdot 1^2 = 2 \sum_{k=1}^{m-1} k^2 = \frac{2}{3} m^3 + O(m^2)$$

(half as much as QR)

If we never encounter divisions by zero, this produces a factorization $A = LU$ for every $A \in \mathbb{R}^{m \times m}$.

We can use this factorization to solve linear systems:

$$Ax = b \rightarrow LUx = b \xrightarrow{y} \begin{cases} Ux = y \\ Ly = b \end{cases} \quad \text{This form can be solved}$$

With substitutions:

- ① Compute $A = LU$

② solve $Lg = b$ for $g \in \mathbb{R}^n$

③ solve $Ux = g$ for $x \in \mathbb{R}^n$

equivalent to writing

$$x = A^{-1}b = (L U)^{-1}b = U^{-1}(L^{-1}b)$$

(U) factorization is NOT stable

$$\begin{bmatrix} 10^{-30} & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & -1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 10^{30} & 1 & 0 \\ 10^{30} & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} 10^{-30} & 1 & 1 \\ 0 & -10^{30} & -10^{30} \\ 0 & -10^{30} & -1-10^{30} \end{bmatrix}$$

\downarrow
becomes exactly singular
when you approximate it
with floating point numbers

Stability analysis:

computed \tilde{L}, \tilde{U} are s.t. $\tilde{L} \tilde{U} = A + \Theta$ $\|\Theta\| = O(n) \underbrace{\|L\| \cdot \|U\|}_{\text{may both be much larger than } \|A\|}$

Fix: pivoting: exchange rows at each time step to ensure pivots are as much large as possible

Step 1: Swap rows to bring the largest value in the 1st column to (1,1) $A = \begin{bmatrix} x & x & x \\ \text{red box} & x & x \\ x & x & x \end{bmatrix}$, then proceed with

step 1 of G.E. LU factorization with pivoting is stable most of the time!

Pivoting ensures $\|L_{ij}\| \leq 1$

However, $\|U\|/\|A\|$ can still grow exponentially $\sim 2^m$

LU with pivoting is still the go-to algorithm to solve square linear problems. With QR, instead, $\|Q\|=1$, $\|R\|=\frac{\|A\|}{\|A\|}$ but it costs 2x times.

Note: we cannot use LU for least-squares problems because multiplying by L scales every row.

LU factorization + sparse matrices:

fill in: entries in A that were zero get filled in.

Direct algorithms for symmetric linear systems

Symmetric matrix; can we take advantage of symmetry

$$Mx = b \quad , \quad M = M^T$$

$$L_1 M \quad L_1 = \begin{bmatrix} 1 & & & \\ x & \ddots & & \\ & \ddots & \ddots & \\ & & & 1 \end{bmatrix}$$

trick: to keep symmetry, instead of multiplying $L_1 M$, compute

$$L_1 M L_1^T \quad (L_1 M L_1^T) = L_1 M L_1 = \begin{bmatrix} x & 0 & 0 & 0 \\ 0 & x & - & x \\ \vdots & x & \ddots & \vdots \\ 0 & x & - & x \end{bmatrix} \quad \text{by symmetry, there are zeros}$$

Similarly, after more steps,

$$L_2 (L_1 M L_1^T) L_2^T = \begin{bmatrix} x & 0 & \cdots & - & \cdots & 0 \\ 0 & x & \cdots & 0 & \cdots & - \\ \vdots & 0 & x & \cdots & - & x \\ 0 & 0 & \vdots & \vdots & \vdots & x \end{bmatrix}$$

$$L_{m-1} L_{m-2} \cdots L_2 L_1 M L_1^T L_2^T \cdots L_{m-2}^T L_{m-1}^T = 0 \quad \text{diagonal}$$

$$M = L_1^{-1} L_2^{-1} \cdots L_{m-1}^{-1} D (L_{m-1}^T)^{-1} (L_{m-2}^T)^{-1} \cdots (L_2^T)^{-1} (L_1^T)^{-1} = L D L^T$$

$$\begin{bmatrix} 1 & & & \\ x & \ddots & & 0 \\ \vdots & \vdots & \ddots & \\ x & x & \ddots & x \end{bmatrix}$$

$$L^T$$

lower tr.
with 1 on
diagonal

Algoletes to get the next matrix:

$$(M_K)_{K:m, K:m} = (M_{K-1})_{K:m, K:m} - \frac{f}{(M_{K-1})_{K-1, K-1}} \cdot (M_{K-1})_{K, K:m}$$

↑ ↑ ↓

3 symmetric matrices

$$M = \begin{bmatrix} M_{11} & M_{12} & M_{13} & M_{14} \\ M_{21} & \boxed{M_{22}} & M_{23} & M_{24} \\ M_{31} & M_{32} & \boxed{M_{33}} & M_{34} \\ M_{41} & M_{42} & M_{43} & \boxed{M_{44}} \end{bmatrix}$$

$$M^{(2)}_{2:4, 2:4} \leftarrow M^{(1)}_{2:4, 2:4} - \begin{bmatrix} \frac{M_{21}}{M_{11}} \\ \frac{M_{31}}{M_{11}} \\ \frac{M_{41}}{M_{11}} \end{bmatrix} \begin{bmatrix} M_{12}^{(1)} & M_{13}^{(1)} & M_{14}^{(1)} \end{bmatrix}$$

Since all matrices are symmetric, one can compute the lower triangle, and get the upper triangle from symmetry. Computational cost solved!

$$\frac{2}{3} n^3 \text{ LU}$$

$$\frac{1}{3} n^3 \text{ LDU}^\top$$

We need column pivoting for stability. In order to keep symmetry one must swap columns accordingly:

swap row 1-3. $\rightarrow \downarrow \quad \downarrow$

and column 1-3. $\rightarrow \quad \times \quad \times$

Th: if M is positive definite, then

- no zero pivots
- factorization is stable

To prove it, we need two easy properties of P.D. matrices:

1) $M \succ 0 \rightarrow M_{kk} > 0 \quad \forall k = 1, 2, \dots, n$

$$M_{kk} = [0 \dots 0 \underset{k}{\uparrow} 0 \dots 0] \in \mathbb{R}^n \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} = e_k^T M e_k > 0$$

Recap: $M \succ 0$ means $x^T M x > 0 \quad \forall x \neq 0$

2) $M \succ 0 \rightarrow B M B^T \succ 0 \quad \forall B \text{ invertible} \quad M, B \in \mathbb{R}^{n \times n}$

$$x \neq 0 \quad x^T (B M B^T) x = y^T M y > 0$$

\uparrow
 $y = B^T x \neq 0$

Properties 1+2 ensure that all matrices encountered during LDL are P.D., and all pivots are > 0 .

In particular, $D_{kk} > 0 \quad \forall k$. The LDL^T factorization for a P.D. matrix can be written in 2 different forms:

$$D = \sqrt{d_{11}} d_{22} \quad \rightarrow \quad \sqrt{\sqrt{d_{11}} \sqrt{d_{22}}} \quad \rightarrow \quad \sqrt{\sqrt{d_{11}} \sqrt{d_{22}}} \quad \rightarrow$$

$$\left[\begin{array}{c} \vdots \\ d_{nn} \end{array} \right] = \left[\begin{array}{c} \vdots \\ \sqrt{d_{nn}} \end{array} \right] \cdot \left[\begin{array}{c} \vdots \\ \sqrt{d_{22}} \\ \vdots \\ \sqrt{d_{nn}} \end{array} \right]$$

$$M = LDL^T = \frac{\left(D^{1/2} \right)^{1/2} L^T}{R^T} R = R^T R \quad R \text{ upper triangular}$$

Th: Every P. D. $M \succ 0$ can be written as $M = R^T R$, where $R^T \in \mathbb{R}^{n \times n}$ is upper triangular. **Cholesky factorization**

Recall: from LS problem: $A \in \mathbb{R}^{m \times n}$, $A = Q \circ R_0 \rightarrow A^T A = R_0^T R_0$

$$G_i(R) = G_i(A^T A) \rightarrow \|R\| = \|A\|^{1/2} \quad \text{norm of } R \text{ doesn't grow!}$$

We can use LDL^T and $R^T R$ to solve linear systems!

$$Mx = b \quad M = LDL^T$$

$$x = M^{-1} b = (L^T)^{-1} D^{-1} L^{-1} b$$

$$\begin{cases} y = L^{-1} b \rightarrow Ly = b \\ z = D^{-1} y \rightarrow Dy = y \\ x = (L^T)^{-1} z \rightarrow L^T x = z \end{cases}$$

$\frac{1}{3} n^3$ for the factorization +
 $O(n^2)$ for the substitution

How to choose an algorithm to solve $Ax = b$?

A pos. def \rightarrow use Cholesky / LDL^T

$$\left. \right\} \frac{1}{3}n^3$$

A symmetric \rightarrow use $LDL^T +$ pivoting

A non-symmetric \rightarrow use LU

$$\left. \right\} \frac{2}{3}n^3$$

If A is sparse, one can use sparse versions of these 3 algorithms