

# Review: matrices & decompositions

## Outer and inner products

Inner product of  $a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \in \mathbb{R}^n$  and  $b = \begin{pmatrix} b_1 \\ \vdots \\ b_n \end{pmatrix} \in \mathbb{R}^n$  is

$$a^T b = \sum_{i=1}^n a_i b_i = \langle a, b \rangle$$

Recall,  $a$  and  $b$  are orthogonal if  $a^T b = 0$

Outer product of  $a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \in \mathbb{R}^n$  and  $b = \begin{pmatrix} b_1 \\ \vdots \\ b_m \end{pmatrix} \in \mathbb{R}^m$  is

$$a \cdot b^T = \begin{pmatrix} a_1 b_1 & \dots & a_1 b_m \\ \vdots & & \vdots \\ a_n b_1 & \dots & a_n b_m \end{pmatrix} \in \mathbb{R}^{n \times m}$$

$$\frac{a^T}{b} = \text{scalar}$$

$$|a - \overline{b^T}| = \boxed{\text{matrix}}$$

Inner product

vs

Outer product

Given  $A \in \mathbb{R}^{n \times p}$  and  $B \in \mathbb{R}^{p \times m}$  there are two ways to view the product  $AB$ .

"inner product view"

$$A = \begin{pmatrix} -\alpha_1^T - \\ \vdots \\ -\alpha_n^T - \end{pmatrix} \quad B = \begin{pmatrix} | & | \\ b_1 & \dots & b_m \\ | & | \end{pmatrix}$$

$$AB = \begin{pmatrix} \alpha_1^T b_1 & \dots & \alpha_1^T b_m \\ \vdots & \ddots & \vdots \\ \alpha_n^T b_1 & \dots & \alpha_n^T b_m \end{pmatrix}$$

"outer product view"

$$A = \begin{pmatrix} | & | \\ a_1 & \dots & a_p \\ | & | \end{pmatrix} \quad B = \begin{pmatrix} -b_1^T - \\ \vdots \\ -b_p^T - \end{pmatrix}$$

$$AB = \sum_{j=1}^p a_j \cdot b_j^T = \underbrace{\boxed{n \times m}}_{p \text{ matrices}} + \dots + \underbrace{\boxed{n \times m}}$$

## Rank factorization

Given  $A \in \mathbb{R}^{n \times p}$  with  $\text{rank}(A) = r$ , one can find

- full column rank  $C \in \mathbb{R}^{n \times r}$
- full column rank  $F \in \mathbb{R}^{p \times r}$

such that  $A = C F^T = \sum_{j=1}^r c_j f_j^T$

$$A = \boxed{C} \quad = \quad \boxed{c_1} \quad \boxed{f_1^T} \quad = \quad \underbrace{\boxed{c_1 f_1^T} + \dots + \boxed{c_r f_r^T}}_{r \text{ matrices of rank 1}}$$

## QR decomposition

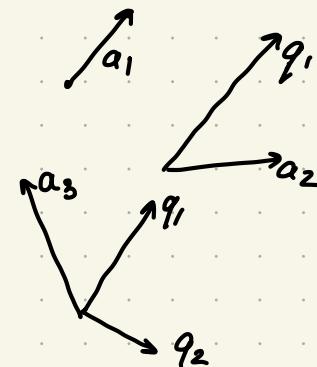
Given  $A \in \mathbb{R}^{n \times p}$  it can be decomposed as  $A = QR$

- where
- $Q$  is column-orthogonal, i.e.  $Q^T Q = I$
  - $R$  is upper-triangular

$$\begin{array}{c|c|c} \begin{array}{|c|c|} \hline 1 & 1 \\ \hline q_1 \dots q_p \\ \hline 1 & 1 \\ \hline \end{array} & = & \begin{array}{c|c} 1 & 1 \\ q_1 \dots q_p \\ 1 & 1 \\ \hline Q & \end{array} & \begin{array}{c|c} 1 & 1 \\ 0 & \ddots \\ 0 & 0 \\ \hline R & \end{array} \end{array}$$

To find  $Q$  and  $R$  Gram-Schmidt procedure is used.

- $e_1 = a_1 \rightarrow q_1 = \frac{e_1}{\|e_1\|}$
- $e_2 = a_2 - \sum_{k=1}^{k-1} \langle a_2, q_i \rangle q_i \rightarrow q_2 = \frac{e_2}{\|e_2\|}$
- $e_k = a_k - \sum_{i=1}^{k-1} \langle a_k, q_i \rangle q_i \rightarrow q_k = \frac{e_k}{\|e_k\|}$



What if  $\text{rank}(A) = r < p$ ? This will be reflected in R.

$$\begin{array}{c|c} \begin{array}{|c|c|} \hline & | \\ \hline a_1 & \dots & a_p \\ \hline & | \\ \hline \end{array} & = \end{array} \begin{array}{c|c} \begin{array}{|c|c|c|} \hline & | & | \\ \hline q_1 & \dots & q_m & q_p \\ \hline & | & | \\ \hline \end{array} & \left. \begin{array}{c} \\ \\ \\ \end{array} \right\} r \\ \underbrace{\qquad\qquad}_{\text{basis for}} \\ \text{the column} \\ \text{space of } A \end{array}$$

Example: linear regression

$$\hat{\beta} = (X^T X)^{-1} X^T y = (R^T R)^{-1} R^T Q^T y = R^{-1} Q^T y$$

$R \hat{\beta} = Q^T y$  can be solved via back solving (fast!)

## Eigen (spectral) decomposition

Consider  $A \in \mathbb{R}^{n \times n}$ , assume  $A$  is symmetric.

Then  $A = U \Lambda U^T$  where

- $U \in \mathbb{R}^{n \times n}$  is orthogonal, i.e.  $U^T U = U U^T = I$

- $\Lambda = (\lambda_1, \dots, \lambda_n)$  is diagonal

$U = (U_1, \dots, U_n)$ ,  $U_i$  are eigenvectors

$\lambda_i$  are corresponding eigenvalues

Sometimes, it is useful to view ED as

$$A = \sum_{i=1}^n \lambda_i U_i U_i^T$$

## Properties:

- $A u_j = \lambda_j u_j$ , i.e. mapping  $A$  doesn't rotate  $u_j$ .  
|  $A = \sum_{i=1}^n \lambda_i u_i u_i^T \Rightarrow A u_j = \sum_{i=1}^n \lambda_i u_i u_i^T u_j = \lambda_j u_j$
- Eigen vectors are determined up to sign  
|  $u_i \rightarrow -u_i$
- The rank of  $A$  is the number of non-zero  $\lambda_i$ .
- If  $V = \sum_{i=1}^n d_i u_i$  then  $A V = \sum_{i=1}^n \lambda_i d_i u_i$   
|  $A V = \sum_{i=1}^n d_i A u_i = \sum_{i=1}^n d_i \lambda_i u_i$  ↑  
Scaling factors

## Definite Matrices

A symmetric matrix  $A \in \mathbb{R}^{n \times n}$  is positive-definite (PD)  
if  $v^T A v > 0$  for all  $v \in \mathbb{R}^n$   $v \neq 0$

It is positive semi-definite (PSD)  
if  $v^T A v \geq 0$  for all  $v \in \mathbb{R}^n$

$$\text{PSD: } A \succeq 0$$

$$\text{PD: } A > 0$$

## Properties:

$$\bullet A \geq 0 \Leftrightarrow \lambda_i \geq 0$$

$$\bullet A > 0 \Leftrightarrow \lambda_i > 0$$

$| A > 0 \Rightarrow v^T A v > 0 \quad \forall v \neq 0,$  let  $v = u_i;$

$$u_i^T A u_i = \lambda_i \|u_i\|^2 = \lambda_i > 0$$

$$\lambda_i > 0 \Rightarrow v^T A v = (u_i^T v)^T \Lambda (u_i^T v) = \sum_{i=1}^n \lambda_i (y_i)^2$$

$\bullet$  ED can be used to compute powers of  $A$   
as  $A^k = U \Lambda^k U^T$

$$| A^k = \underbrace{U \Lambda}_{\Lambda^k} \underbrace{U^T}_{\Lambda^k} \cancel{\Lambda} \cancel{U^T} = U \Lambda^k U^T$$

$\bullet$  for  $A > 0$  it holds  $A^{-1} = U \Lambda^{-1} U^T$

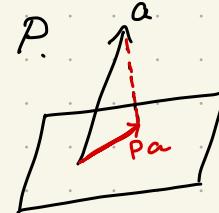
for  $A \geq 0$  it is true  $A^{1/2} = U \Lambda^{1/2} U^T$

$$| A^{1/2} \cdot A^{1/2} = U \Lambda^{1/2} \cancel{U^T} \cancel{\Lambda} \Lambda^{1/2} U^T = U \Lambda U^T = A$$

## Projection

Matrix  $P \in \mathbb{R}^{n \times n}$  is a projection matrix if  $P^2 = P$ .

Orthogonal projection is when  $P = P^T = P^2$ .



Properties: if  $P$  is orthogonal projection then

- Eigenvalues are 0 or 1

$$P = U \Lambda U^T, \quad P^2 = \sqrt{\Lambda} \Lambda^2 \sqrt{\Lambda}^T = P = \sqrt{\Lambda} \Lambda^T$$

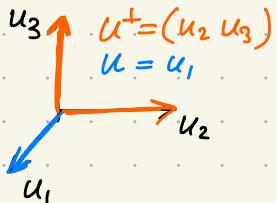
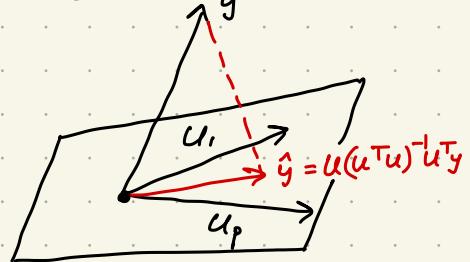
- It is PSD.

- Diagonal elements are bounded  $0 \leq p_{ii} \leq 1$

$$P = P^2 = P^T P \Rightarrow p_{ii} = \sum_{k=1}^n p_{ki}^2 \geq p_{ii}^2 \geq 0$$

## Examples

- Given  $U \in \mathbb{R}^{n \times p}$ ,  $P = U(U^T U)^{-1} U^T$  is orthogonal projection onto the column space of  $U$ 
  - $y \in \mathbb{R}^n \quad P_U(y) = U\beta$  and
  - minimize  $\|y - U\beta\|^2$  — regression
- If  $U$  is column orthogonal  $P = UU^T$
- $P = I - UU^T$  is projection onto the column space of  $U^\perp$ 
  - $\tilde{U} = (U, U^\perp) \Rightarrow \tilde{U} \tilde{U}^T = UU^T + U^\perp U^\perp^T = I$



$U_{\perp}$  is orthogonal complement of the basis  $U$  if

- $U_{\perp} \in \mathbb{R}^{n \times (n-p)}$  is column orthogonal, i.e.  $U_{\perp}^T U_{\perp} = I$
- $\tilde{U} = (U \ U_{\perp}) \in \mathbb{R}^{n \times n}$  is orthogonal, i.e.  $\tilde{U} \tilde{U}^T = \tilde{U}^T \tilde{U} = I$

This also implies that

- $U_{\perp}$  is column orthogonal and  $U^T U_{\perp} = 0_{p \times (n-p)}$

## Column centering

Given sample matrix  $X \in \mathbb{R}^{n \times p}$ ,  $X = \begin{pmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{pmatrix}$

- the **Sample mean** is  $\bar{x} = \frac{x^T 1}{n} = \frac{1}{n} \sum_{i=1}^n x_i$

- the column centered version is  $\tilde{x} = C \cdot x$ ,

where  $C = I - \frac{1 \cdot 1^T}{n} \in \mathbb{R}^{n \times n}$  is **centering operator**

$$C x = x - \frac{1^T}{n} x = x - 1 \cdot \bar{x}^T = x - \begin{pmatrix} \bar{x}^T \\ \vdots \\ \bar{x}^T \end{pmatrix}$$

$C$  is Orthogonal projection Operator

(Onto what space?)

|  $U = \frac{1}{\sqrt{n}} \cdot 1_n \in \mathbb{R}^n$  then  $U^T U = 1$ ,  $C$  is projection on  $U^\perp$

## Sample covariance matrix

Given  $X = (f_1, \dots, f_p) \in \mathbb{R}^{n \times p}$  sample covariance matrix is  $S \in \mathbb{R}^{p \times p}$  with  $S_{ij} = \text{cov}(f_i, f_j)$ .

If  $X = \begin{pmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{pmatrix}$  and  $C = I - \frac{11^T}{n}$ , then

$$S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = \frac{X^T C X}{n-1}$$

$$| X^T C X = X^T C^T C X = \tilde{X}^T \tilde{X}$$

If  $X$  is centered then  $S = \frac{X^T X}{n-1}$

## Properties:

- Sample covariance matrix is PSD
- |  $V^T S V = \frac{1}{n-1} V^T X^T X V = \frac{\|XV\|^2}{n-1}$
- In fact,  $V^T S V = \text{Var}(X_V)$
- In regression  $\hat{\beta} = S_x^{-1} S_{xy}$  where  
 $S_x \in \mathbb{R}^{P \times P}$  is sample covariance for  $X$  and  
 $S_{xy} \in \mathbb{R}^{P \times 1}$  is sample covariance between  $X$  and  $y$

$$\hat{\beta} = \left( \frac{X^T X}{n-1} \right)^{-1} \frac{X^T y}{n-1}$$

## Singular Value Decomposition (SVD)

Given  $A \in \mathbb{R}^{n \times p}$ , SVD expresses  $A$  as

$$A = UDV^T \text{ where}$$

- $U \in \mathbb{R}^{n \times n}$  and  $V \in \mathbb{R}^{p \times p}$  are orthogonal
- $D \in \mathbb{R}^{n \times p}$  diagonal with non-negative entries

sorted in decreasing order.

$n > p$

$$A = \begin{matrix} | & | \\ u_1, \dots, u_n \\ | & | \end{matrix} \begin{matrix} d_1 & & & \\ & \ddots & & \\ & & d_p & \\ & & & \vdots \\ & & & v_1^T \\ & & & \vdots \\ & & & v_p^T \end{matrix}$$

left singular vectors      singular values      right singular vectors

$$\begin{array}{|l} U^T U = U U^T = I \\ V^T V = V V^T = I \\ d_1 \geq \dots \geq d_p \geq 0 \end{array}$$

$n > p$

$$\boxed{A} = \begin{matrix} | \\ | \\ | \\ | \end{matrix} \boxed{u_1, \dots, u_n} \boxed{d_1 \dots d_p} \boxed{-v_1^T \dots -v_p^T}$$

Compact form for SVD:

$$\boxed{A} = \begin{matrix} | \\ | \\ | \\ | \end{matrix} \boxed{1 \dots} \boxed{d_1 \dots d_p} \boxed{-v_1^T \dots -v_p^T}$$

↑                           ↑  
diagonal                   orthogonal  
column orthogonal

Note that  $A = \sum_{i=1}^{\min(n,p)} d_i \cdot \underbrace{u_i \cdot v_i^T}_{\text{rank-1 matrices ranked by "importance"}}$

## Properties:

- $A V_i = d_i u_i$  and  $A^T u_i = d_i v_i$  for  $i \leq \min(n, p)$
- |  $A V_j = \sum_{i=1}^p d_i u_i v_i^T v_j = d_j c_j \|v_j\|^2$
- Every  $A$  has SVD.
- $D$  is unique,  $U$  and  $V$  are not.

For example, if we have repeated singular values there are many ways to compute  $U$  and  $V$ .

$$D = \begin{pmatrix} d & & \\ & \ddots & \\ & & d \end{pmatrix} = d I \text{ then } UDV^T = \underbrace{UQ}_U D Q^T V^T = \tilde{U} D \tilde{V}^T$$

$$D = \begin{pmatrix} d & & & \\ & d & & \\ & & \ddots & \\ & & & d_p \end{pmatrix} \quad U = \left( \underbrace{\begin{matrix} | & | \\ u_1 & \dots & u_k \\ | & | \end{matrix}}_Q \dots \right) \quad V = \left( \underbrace{\begin{matrix} | & | \\ v_1 & \dots & v_k \\ | & | \end{matrix}}_Q \dots \right)$$

- If  $\text{rank}(A) = r$ , then  $d_{r+1} = \dots = d_p = 0$

## Link between SVD and ED

- If  $A$  is PSD, the decompositions are the same.  
What if  $A$  is symmetric but not PSD?

$$\begin{aligned} A &= U \Lambda U^T \text{ if } d_i \geq 0 \text{ then } U = U, D = \Lambda, V = U \\ \text{If } d_i < 0 \text{ then } u_i &= u_i, d_i = -d_i, v_i = -u_i \end{aligned}$$

Given  $A \in \mathbb{R}^{n \times p}$  with  $n > p$

- $v_i$  are eigenvectors for  $A^T A$ ,  $d_i^2$  are eigenvalues

$$A = \underbrace{U \Lambda V^T}_{p \times p} \quad A^T A = V \Lambda V^T \cancel{U \Lambda U^T} = V \Lambda^2 V^T$$

- $u_i$  are eigenvectors for  $A A^T$ ,  $d_i^2$  are eigenvalues  
(plus some extra zeroes)

$$A A^T = \underbrace{U \Lambda^2 U^T}_{n \times p} \quad \underbrace{\Lambda^2}_{p \times p}$$

Example  $X = UDV^T$  then ridge regression solution

$$\hat{y}_\lambda = X \hat{\beta}_\lambda = X (X^T X + \lambda I)^{-1} X^T y = \sum_{i=1}^{\min(p, n)} \frac{d_i^2}{d_i^2 + \lambda} \underbrace{\langle u_i, y \rangle}_{\text{scaledown factor}} u_i$$

length of projection of  $y$  onto  $u_i$ .

$$\hat{y}_\lambda = UDV^T (V D^2 V^T + \lambda I)^{-1} V D U^T y =$$

$$= UDV^T \cancel{N(D^2 + \lambda I)^{-1} N} V D U^T y =$$

$$= U \frac{\lambda^2}{D^2 + \lambda I} U^T y = \sum_{i=1}^p \frac{d_i^2}{d_i^2 + \lambda} u_i u_i^T y$$

## Matrix and Vector norms

Given vector  $v \in \mathbb{R}^n$

- $\ell_1$  norm  $\|v\|_1 = \sum_{i=1}^n |v_i|$

- $\ell_2$  norm  $\|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2}$  (often  $\|v\|$ )

Given matrix  $A \in \mathbb{R}^{n \times p}$

- Frobenius norm  $\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p A_{ij}^2 = \text{tr}(A^T A) = \sum_{i=1}^{\min(n,p)} d_i^2(A)$

Recall that for  $A \in \mathbb{R}^{n \times n}$  trace is  $\text{tr}(A) = \sum_{i=1}^n a_{ii}$

### Properties

- $\text{tr}(A) = \text{tr}(A^T)$
- $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$
- Nuclear norm  $\|A\|_* = \sum_{i=1}^{\min(n,p)} d_i(A)$

## Gradients

If  $f: \mathbb{R}^n \rightarrow \mathbb{R}$ , so it takes  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n$  and maps it to a number, then

$$\nabla_x f(x) = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix} \in \mathbb{R}^n$$

If  $f: \mathbb{R}^{n \times p} \rightarrow \mathbb{R}$ , so it takes  $A = \begin{pmatrix} A_{11} \dots A_{1p} \\ \vdots \\ A_{n1} \dots A_{np} \end{pmatrix} \in \mathbb{R}^{n \times p}$  and maps it to a number, then

$$\nabla_A f(A) = \begin{pmatrix} \frac{\partial f(A)}{\partial A_{11}} & \dots & \frac{\partial f(A)}{\partial A_{1p}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{n1}} & \dots & \frac{\partial f(A)}{\partial A_{np}} \end{pmatrix} \in \mathbb{R}^{n \times p}$$

## Examples:

- $S \in \mathbb{R}^{n \times n}$ ,  $S = S^T$ ,  $x \in \mathbb{R}^n$

if  $f(x) = x^T S x$  then  $\nabla_x f(x) = 2 S x$

$$x^T S x = \sum_{i=1}^n \sum_{j=1}^n S_{ij} x_i x_j \quad \frac{\partial f(x)}{\partial x_i} = 2 \sum_{j=1}^n S_{ij} x_j = 2 S_i^T x$$

$$\nabla_x f(x) = \begin{pmatrix} 2 S_1^T x \\ \vdots \\ 2 S_n^T x \end{pmatrix} = 2 S x$$

- $A, B \in \mathbb{R}^{n \times p}$

if  $f(A) = \|A\|_F^2$  then  $\nabla_A f(A) = 2 A$

$$\|A\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p A_{ij}^2 \quad \frac{\partial f(A)}{\partial A_{ij}} = 2 A_{ij} \quad \nabla_A f(A) = 2 \begin{pmatrix} A_{11} & \dots & A_{1p} \\ \vdots & \ddots & \vdots \\ A_{nn} & \dots & A_{np} \end{pmatrix} = 2A$$

if  $f(A) = \text{tr}(AB^T)$  then  $\nabla_A f(A) = B$

$$\text{tr}(AB^T) = \sum_{i=1}^n \sum_{j=1}^p A_{ij} B_{ij} \quad \frac{\partial f(A)}{\partial A_{ij}} = B_{ij}$$

## Recap from Lecture 1

Rank factorization:  $\text{rank}(A) = r \Rightarrow C$  and  $F$  are full col. rank

$$\begin{matrix} A &= & C &\cdot F^T \\ \boxed{n \times p} &= & \boxed{n \times r} &\cdot \boxed{r \times p} \end{matrix}$$

---

QR decomposition:  $Q^T Q = I$ ,  $R$  is upper triangular

$$\begin{matrix} A &= & Q &\cdot R \\ \boxed{n \times p} &= & \boxed{n \times p} &\cdot \boxed{p \times p} \end{matrix}$$

---

SVD:  $U^T U = I$ ,  $D = \text{diag}(d_1, \dots, d_p)$ ,  $d_1 \geq \dots \geq d_p \geq 0$ ,  $V^T V = V V^T = I$

$$\begin{matrix} A &= & U &\cdot D &\cdot V^T \\ \boxed{n \times p} &= & \boxed{n \times p} &\cdot \boxed{p \times p} &\cdot \boxed{p \times p} \end{matrix}$$

**ED:**  $A^T = A \Rightarrow U^T U = U U^T = I$ ,  $\Lambda = \text{diag}(d_1, \dots, d_n)$

$$\begin{array}{c} A = U \cdot \Lambda \cdot U^T \\ \boxed{n \times n} = \boxed{n \times n} \quad \boxed{n \times n} \quad \boxed{n \times n} \end{array}$$

**Property:**  $A = U D V^T \Rightarrow A A^T = \underbrace{U D^2 U^T}_{\text{e.vec e.val}} \text{ and } A^T A = \underbrace{V D^2 V^T}_{\text{e.vec e.val}}$

**Orthogonal projection:**  $P^2 = P = P^T$

- onto  $U = (\vec{u}_1, \dots, \vec{u}_p) \Rightarrow P = U(U^T U)^{-1} U^T$
- onto  $U$  if  $U^T U = I \Rightarrow P = U U^T$
- onto  $U^\perp \Rightarrow P = I - U U^T$

**Example:**  $C = I - \frac{11^T}{n}$  centering operator

**PSD:** Symmetric  $A \geq 0$  if  $V^T A V \geq 0$  for any  $V$ .  
equivalently,  $\lambda_i(A) \geq 0$

**Examples:** orth. projection, sample covariance