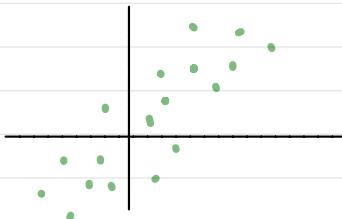


①

# Singular Value Decomposition

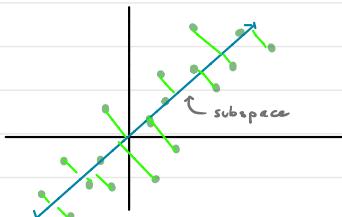
Principal axes: let  $y_1, \dots, y_n \in \mathbb{R}^p$ .



The data are  $p$ -dimensional, but the point cloud may be "close" to a lower-dimensional (linear) space.

Goal: ① Find the "best" low dimensional subspace  
 ② Find the approximation  $\hat{y}_i$  of  $y_i$  for  $i=1\dots,n$ .

Intuitively:



Start with question ② first.

Let  $L \subset \mathbb{R}^p$  be an  $r$ -dimensional subspace

Let  $B = [b_1 \cdots b_r]$  be an orthonormal basis for  $L$

This means ①  $\forall x \in L \exists a \in \mathbb{R}^r : x = a.b_1 + \dots + a_r.b_r$   $\begin{cases} B \text{ is a} \\ \text{basis} \end{cases}$

②  $b_j^T b_j = 1, b_j^T b_n = 0, B^T B = I_r \quad \begin{cases} B \text{ is orthonorm} \\ B \in \mathcal{D}_{r,p} \end{cases}$

z

Task: for  $y_i \in \mathbb{R}^p$ , find  $\hat{y}_i \in L$  to minimize  $\|y_i - \hat{y}_i\|^2$ .

Solution:  $\hat{y}_i \in L \Leftrightarrow \hat{y}_i = Ba$  some  $a \in \mathbb{R}^r$

$$\begin{aligned}\|y_i - \hat{y}_i\|^2 &= \|y_i - Ba\|^2 = (y_i - Ba)^T (y_i - Ba) \\ &= y_i^T y_i - 2 y_i^T Ba + a^T B^T B a \\ &= a^T a - 2 a^T B^T y_i + c\end{aligned}$$

Minimize in  $a$  with calculus:  $\nabla_a \|y_i - Ba\|^2 = 2a - 2B^T y_i$ :

$$\text{optimal } a \text{ satisfies } a = B^T y_i$$

(Take 2nd deriv to show a min)

so  $\hat{y}_i = Ba = BB^T y_i$  ( $BB^T$  is the proj. matrix onto  $L$ , the linear space spanned by the r columns of  $B$ )  
 (compare to  $\hat{y}$  in regression)

Matrix form:  $Y = \begin{pmatrix} y_1^T \\ \vdots \\ y_n^T \end{pmatrix} \in \mathbb{R}^{n \times p}$

$$\begin{aligned}\hat{Y} &= \begin{pmatrix} y_1^T B B^T \\ \vdots \\ y_n^T B B^T \end{pmatrix} = \begin{pmatrix} y_1^T B \\ \vdots \\ y_n^T B \end{pmatrix} B^T \in \mathbb{R}^{n \times p} \\ &\quad n \times p \qquad r \times p\end{aligned}$$

$$= Y B B^T$$

(3)

Now solve problem ① Find optimal subspace / optimal  $B$

Task: minimize  $\|Y - YBB^T\|^2$  in  $B \in \mathbb{R}^{P \times r}$ ;  $B^T B = I_r$

Simpler task:  $\min_{b \in \mathbb{R}^P} \|Y - Ybb^T\|^2$  in  $b \in \mathbb{R}^P$ ;  $b^T b = 1$

$$\begin{aligned}\text{Solution: } \|Y - Ybb^T\|^2 &= \|Y(I - bb^T)\|^2 \\ &= \text{tr}(Y(I - bb^T)(I - bb^T)Y^T) \\ &= \text{tr}(Y(I - bb^T)Y^T) \\ &= \text{tr}(YY^T) - \text{tr}(Ybb^TY^T) \\ &= \|Y\|^2 - b^T Y^T Y b\end{aligned}$$

Minimizing in  $b$  means maximizing  $b^T Y^T Y b$  in  $b$ .

Try taking a derivative:  $\frac{\partial}{\partial b} b^T Y^T Y b = 2 Y^T Y b = 0 \Leftrightarrow b = 0$

$$\text{forget } b^T b = 1: \quad \frac{\partial}{\partial b} \left[ b^T Y^T Y b - \lambda(b^T b - 1) \right] = 0$$

$$\Leftrightarrow 2 Y^T Y b - \lambda 2 b = 0 \Leftrightarrow Y^T Y b = \lambda b$$

$\Rightarrow$  critical point when  $b$  is an eigenvector of  $Y^T Y$ .

$\Rightarrow$  critical points of  $\text{RSS}(b) = \|Y - Ybb^T\|^2$  are  $\lambda$ -vecs of  $Y^T Y$

$$\begin{cases} \text{tr}(S) = \sum s_{ii} \\ \text{tr}(A^T A) = \sum \sum a_{ij}^2 \end{cases}$$

(4)

Spectral Decomposition Thm: Let  $S \in \mathbb{R}^{P \times P}$  be symmetric, real.

Then  $S = V \Lambda V^T$  for some  $V \in \mathbb{R}^{P \times P}$  :  $V^T V = V V^T = I_P$   
some  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ ,  $\lambda_1 \geq \dots \geq \lambda_p$

Corollary: Let  $V = [v^1 \dots v^p]$ . Then  $v_j$  is an evec of  $S$ , with eval  $\lambda_j$ .

$$\underline{\text{Proof}}: S v_j = V \Lambda V^T v_j = V \Lambda \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \lambda_j \\ 0 \end{bmatrix} = V \begin{bmatrix} 0 \\ \vdots \\ 0 \\ \lambda_j v_j \\ 0 \end{bmatrix} = \lambda_j v_j.$$

Important result: Let  $S = Y^T Y$ . Then  $\lambda_1 \geq \dots \geq \lambda_p \geq 0$ .

(so  $Y^T Y$  is positive semidefinite)

Exercise: Prove result.

(5)

Q: Which evec gives optimal approximation?

A: Minimizer of  $\|Y - Ybb^T\|$ , maximizer of  $b^T Y^T Y b$ .

Let  $v_1, \dots, v_p$  be evecs of  $Y^T Y$  w/ evals  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$

then  $v_j^T Y^T Y v_j = v_j^T [\lambda_j v_j]$  ( $v_j$  is evec w/ eval  $\lambda_j$ )

$$\stackrel{=}{=} \lambda_j$$

In other words,  $\hat{Y} = Y v_j v_j^T$  is the best one dim approx to  $Y$ .

Similarly,  $\hat{Y} = Y B B^T$ ,  $B = [v_1 \dots v_r]$  is best r-dim approx.

- Proof:
- 1) induction
  - 2) matrix differentials
  - 3) brute-force with Lagrange multipliers

In summary:

Goal: Approx  $x_i$  as  $x_i \approx a_{i1} \underline{b}_1 + \dots + a_{ir} \underline{b}_r$ , some common  $v_1, \dots, v_r$   
 $= B \underline{a}_i$

$$\text{Approx } Y = \begin{pmatrix} x_1 & \dots & x_n \end{pmatrix} \approx \begin{pmatrix} a_1^T B^T \\ \vdots \\ a_n^T B^T \end{pmatrix} = \begin{bmatrix} \underline{a}_1^T & \dots & \underline{a}_n^T \end{bmatrix} B^T = A B^T$$

$n \times r \quad r \times p$

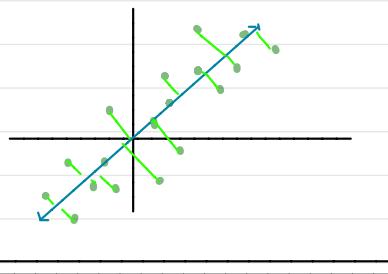
Result: A minimizer of  $\|Y - AB^T\|^2$  is

(princ axes)  $B = [v_1 \dots v_r]$   $v_j = \text{evec}_j(Y^T Y)$

(princ coord)  $A = YB = \text{Coef of projection of rows of } Y \text{ onto evecs of } Y^T Y$

(6)

recall picture:



Singular value decomposition:

Best { $p$ -dim approx to rows of  $Y$  }  
rank- $p$  approx to  $Y$

even of  $Y^T Y$

we can write  $Y = F V^T$ ,  $F = Y V \in \mathbb{R}^{n \times p}$

what are properties of  $F$ ?

$$F^T F = V^T Y^T Y V = V^T [V \Lambda V^T] V = \Lambda \text{ } p \times p \text{ diagonal}$$

Claim:  $F^T F = \Lambda \Rightarrow F = U \Lambda^{1/2}$ , some  $U \in \mathbb{R}^{n \times p}$ :  $U^T U = I_p$ .

Exercise: Prove claim.

$$\begin{aligned} \text{This means } Y \text{ can be written as } Y &= F V^T \\ &= U \Lambda^{1/2} V^T \\ &= U D V^T \end{aligned}$$

where  $U^T U = V^T V = I_p$ ,  $D = \text{diag}(d_1, \dots, d_p)$ ,  $d_1 \geq \dots \geq d_p \geq 0$ .

This representation is called the singular value decomposition (SVD) of  $Y$ .

(7)

Thm (SVD): Let  $Y \in \mathbb{R}^{n \times p}$ ,  $n \geq p$ . Then  $Y$  can be written as

$$Y = UDV^T, \text{ where}$$

$$\textcircled{1} \quad U \in \mathbb{R}^{n \times p}, \quad U^T U = I_p$$

$$\textcircled{2} \quad V \in \mathbb{R}^{p \times p}, \quad V^T V = I_p$$

$$\textcircled{3} \quad D = \text{diag}(d_1, \dots, d_p) \quad d_1 \geq \dots \geq d_p \geq 0.$$

### Comments

\*  $\underline{U}_1, \dots, \underline{U}_p$  are the left sing. vrcs,

\*  $\underline{V}_1, \dots, \underline{V}_p$  are the right sing. vrcs

\*  $d_1, \dots, d_p$  are the singular values

\*  $Y^T Y = V D U^T U D V^T = V D^2 V^T = V \Lambda V^T$  (think: col. cov)
 

- columns of  $V$  are evrcs of  $Y^T Y$
- $D^2$  gives evals of  $Y^T Y$

\*  $Y Y^T = U D V^T V D U^T = U D^2 U^T = U \Lambda U^T$  (think: row cov)
 

- columns of  $U$  are evrcs of  $Y Y^T$
- $D^2$  are evals of  $Y Y^T$
- $Y Y^T$  also has  $n-p$  additional evrcs, with evals = 0.

$$\Rightarrow \text{evrc}(Y Y^T) = [U \ U^\perp], \quad Y Y^T U = U \Lambda$$

$$Y Y^T U^\perp = 0 = U^\perp \cdot 0$$

## Reduced-rank matrix approximation

"Vector" view: data are  $x_1, \dots, x_n \in \mathbb{R}^p$

approximate  $x_i$  with  $\hat{x}_i = \underline{B} \underline{a}_i = b_i a_{i1} + \dots + b_p a_{ip}$

$\Rightarrow B$  is the subspace in which  $x_1, \dots, x_n$  vary

$\Rightarrow$  variation in  $x_1, \dots, x_n \in \mathbb{R}^p$  rep by var. in  $a_1, \dots, a_n \in \mathbb{R}^r$

Matrix view: data are  $Y \in \mathbb{R}^{n \times p}$

approximate  $Y$  with  $\hat{Y} = \underline{A} \underline{B}^T$

$n \times p$  numbers

$r$  numbers  $< n+p + rp = r(n+p)$

## Matrix rank

Let  $X \in \mathbb{R}^{m \times n}$ . Write  $X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ 1 & 1 & \dots & 1 \end{bmatrix} = \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1n} \end{bmatrix}$

$$\text{rank}(X) = \dim(\text{span}(x_1, \dots, x_n))$$

$$= \dim(\text{span}(x_1, \dots, x_n)) \quad \left\} \text{ so } \text{rank}(X) \leq \min(m, n) \right.$$

= smallest  $r$  for which  $\exists A \in \mathbb{R}^{m \times r}, B \in \mathbb{R}^{r \times n}$ ,  
with  $X = AB^T$ .

$$= a_1 b_1^T + \dots + a_r b_r^T$$

Thm (rank matrix): let  $X = \underline{a} \underline{b}^T = \begin{bmatrix} a_1 b_1 & a_1 b_2 \\ a_2 b_1 & \ddots \end{bmatrix}$ . Then  $\text{rank}(X) \leq 1$

Ex: Prove, using 1st 2 definitions of rank.