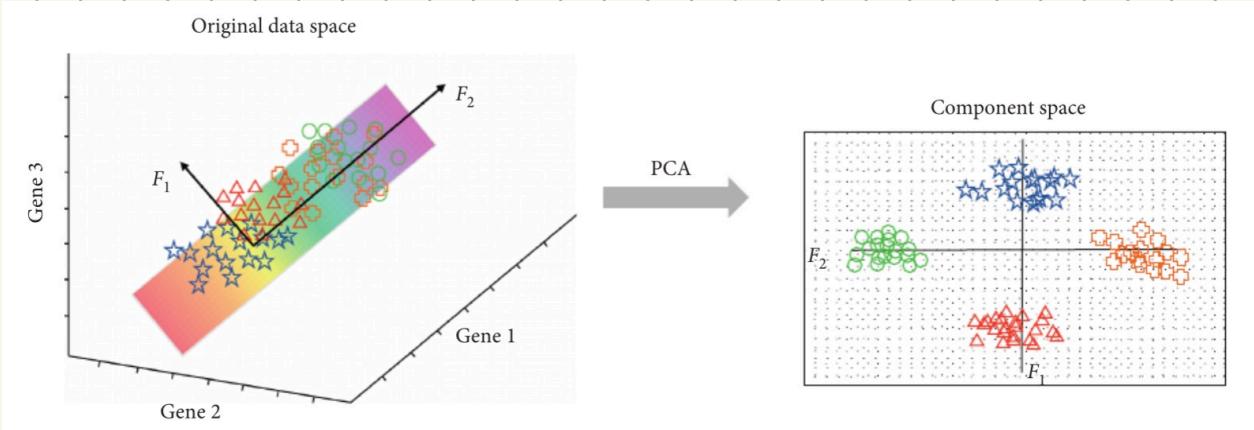


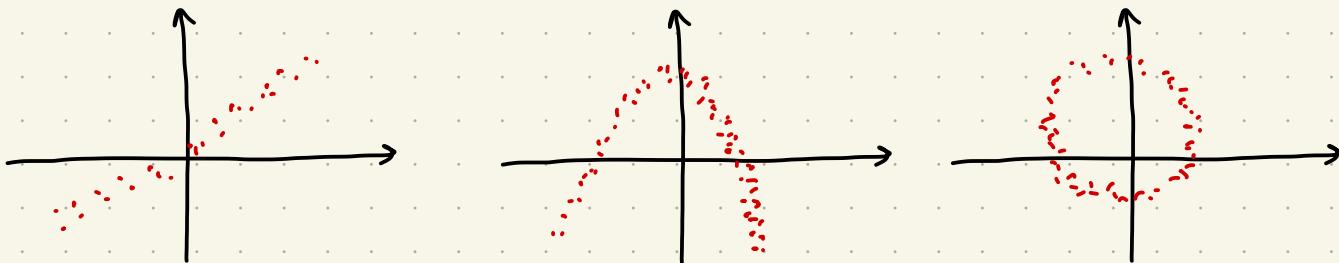
# Principal component analysis (PCA)



## Dimension reduction

Assumption: The data lies along a low-dimensional manifold.

Example: data in 2D, but actually lie on a curve

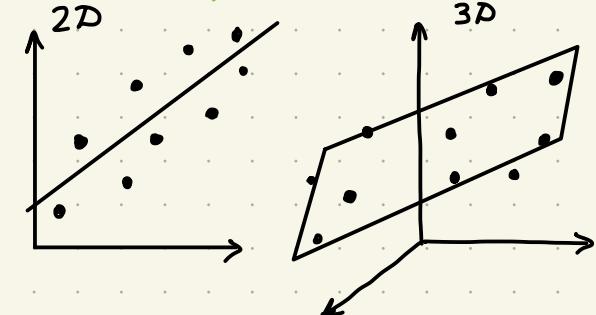


# Principal Component Analysis (PCA)

Data lies on an

$r$ -dimensional affine subspace

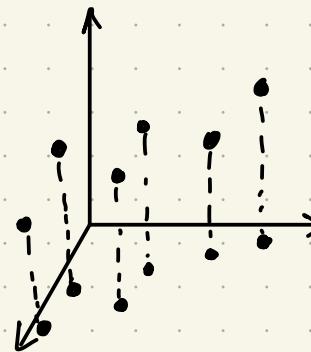
("plane")



Reduce the dimensionality while preserving important information.

How about removing one of the columns?

$$X = \begin{pmatrix} 1 & 1 & 1 \\ f_1 & f_2 & f_3 \\ 1 & 1 & 1 \end{pmatrix}$$



## 1st View of PCA: maximum variance

Suppose we are given  $x_1, \dots, x_n \in \mathbb{R}^p$  points

data matrix  $X = \begin{pmatrix} -x_1^T- \\ \vdots \\ -x_n^T- \end{pmatrix}$

We assume that  $X$  is column-centered, thus,

$$\bar{x} = \frac{X^T \cdot 1}{n} = \frac{\sum_{i=1}^n x_i}{n} = 0$$

$$X \rightarrow X_c = C X = \left( I - \frac{11^T}{n} \right) X, \text{ then}$$

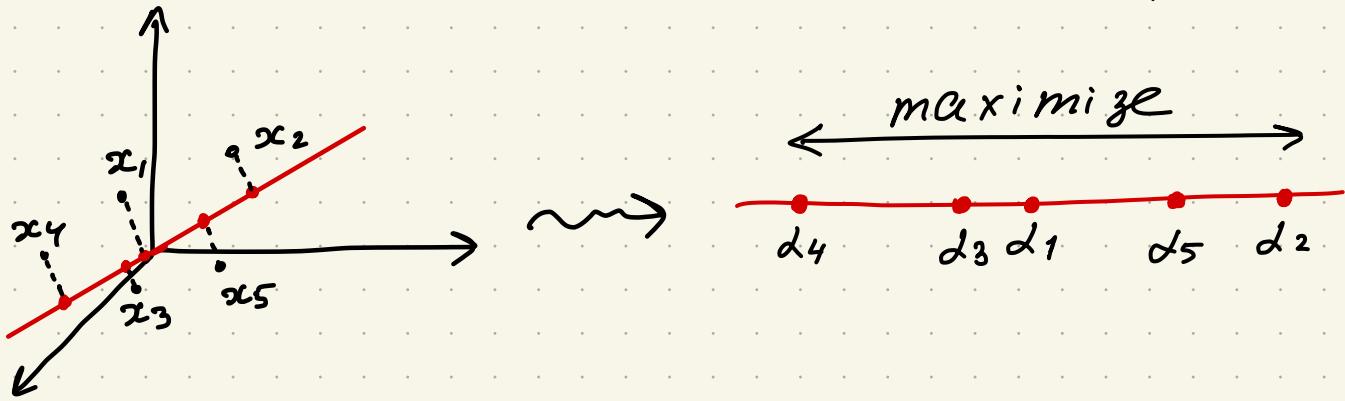
$$\bar{x}_c = (X_c)^T \frac{1}{n} = X^T \cdot \left( I - \frac{11^T}{n} \right) \cdot \frac{1}{n} = X^T \left( \frac{1}{n} - \frac{1}{n} \left( \frac{1^T 1}{n} \right) \right) = 0$$

The covariance matrix is

$$S = \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T = \frac{1}{n-1} X^T X.$$

## PCA: the first principal component ( $PC_1$ )

① Project the data on a line and maximize the variance of the projected points.



The line is defined by  $v \in \mathbb{R}^P$ ,  $\|v\| = 1$

The projection on the line is

$$\hat{x}_i = (x_i^T v) \cdot v = d_i \cdot v$$

|  $P_v = V V^T$ , then  $\hat{x}_i = P_v x_i = V \cdot (V^T x_i)$

Note that  $\bar{z} = \frac{1}{n} \sum_{i=1}^n z_i = 0$

$$\left| \frac{1}{n} \sum_{i=1}^n (v^T x_i) \right| = v^T \left( \frac{1}{n} \sum_{i=1}^n x_i \right) = 0$$

Let's compute the variance of  $z = \begin{pmatrix} z_1 \\ z_n \end{pmatrix} = X \cdot v$

$$\left| \text{var}(z) = \frac{1}{n-1} \sum_{i=1}^n (x_i^T v)^2 = v^T \left( \frac{1}{n-1} \sum_{i=1}^n x_i x_i^T \right) v = v^T S v \right.$$

The goal is to solve

$$\underset{v \in \mathbb{R}^p}{\text{maximize}} \quad v^T S v \quad \text{subject to} \quad \|v\| = 1$$

We will use Lagrange multipliers.

## Détour: method of Lagrange multipliers

Consider  $x \in \mathbb{R}^P$  and two functions

$$f: \mathbb{R}^P \rightarrow \mathbb{R} \text{ and } g: \mathbb{R}^P \rightarrow \mathbb{R}$$

We want to

maximize  $f(x)$  subject to  $g(x) = 0$

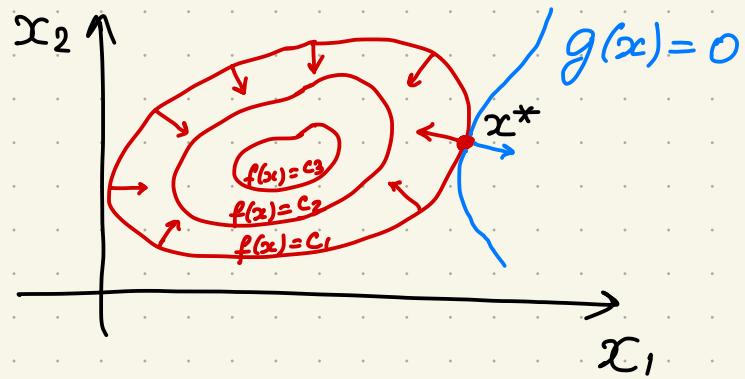
① Form **Lagrangian**  $\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$

$\uparrow$   
**Lagrange multiplier**

② Find critical points of  $\mathcal{L}$ :  $\nabla \mathcal{L}(x, \lambda) = 0$

③ Consider all critical points  $(x^*, \lambda^*)$  of  $\mathcal{L}$   
find  $x^*$  that results in the smallest  $f(x^*)$

- $\nabla \mathcal{L}(x; \lambda) = \begin{pmatrix} \nabla f(x) + \lambda \nabla g(x) \\ g(x) \end{pmatrix} = 0 \Rightarrow \begin{cases} \nabla f(x) = 0 \\ g(x) = 0 \end{cases}$
- contour for  $f(x)$  and  $g(x)$  are tangent at  $x^*$   
 $\nabla f(x)$  and  $\nabla g(x)$  at  $x^*$  are parallel
- there is  $\lambda^*$  such that  $\nabla f(x^*) = -\lambda^* \nabla g(x^*)$



For multiple constraints

minimize  $f(x)$  subject to  $\begin{cases} g_1(x) = 0 \\ \vdots \\ g_k(x) = 0 \end{cases}$

① Form Lagrangian

$$\mathcal{L}(x, \lambda_1, \dots, \lambda_k) = f(x) + \lambda_1 g_1(x) + \dots + \lambda_k g_k(x)$$

② Solve  $\nabla \mathcal{L}(x, \lambda_1, \dots, \lambda_k) = 0$ , find critical points  $x^*, \lambda_1^*, \dots, \lambda_k^*$

③ Use  $x^*$  with the minimum value of  $f(x^*)$ .

maximize  $v^T S v$  subject to  $\|v\|^2 = 1$   
 $v \in \mathbb{R}^{p \times 1}$

- $v$  is an eigenvector of  $S$

$$\begin{cases} \mathcal{L}(v, \lambda) = v^T S v - \lambda(\|v\|^2 - 1) = v^T S v - \lambda(v^T v - 1) \\ \nabla_v = 2Sv - 2\lambda v \Rightarrow Sv = \lambda v \end{cases}$$

- $v$  corresponds to the maximum e. value

$$v^T S v = \lambda \|v\|^2 = \lambda \rightarrow \max$$

We showed that

$$V = v, \text{ and } \text{Var}(z) = v^T S v = \lambda,$$

## Terminology:

- $v_1$  is the **loading vector / direction** of  $PC_1$
- $z_1 = X v_1$  is the **score vector** (or just  $PC_1$ )
- $\lambda_1$  is the **variance explained** by  $PC_1$ .

Comment: the population version

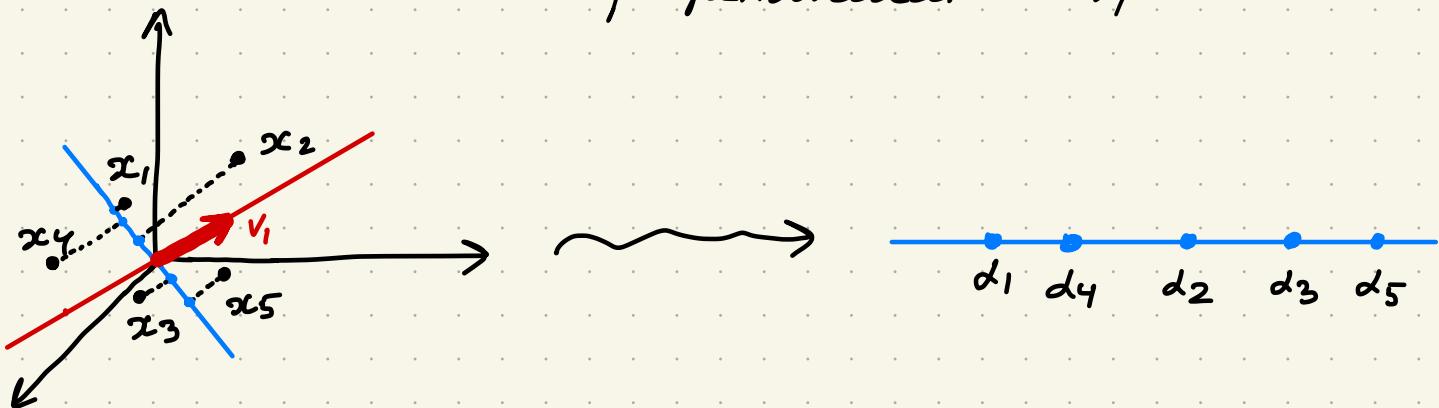
given random  $x = \begin{pmatrix} x_1 \\ x_p \end{pmatrix} \in \mathbb{R}^P$  a linear combination  
 $z = x^T v$  with maximum variance and  $\|v\|^2 = 1$

$$\left| \begin{array}{l} \text{var}(z) = \text{var}(x^T v) = v^T \Sigma v \text{ where} \\ \text{cov}(x) = \Sigma, \text{ then } v \text{ is the first e. vector.} \end{array} \right.$$

## PCA: the second component ( $PC_2$ )

① Project the data on a line and maximize the variance of the projected points.

The line should be perpendicular to  $V_1$



New optimization goal :

$$\text{maximize } V^T S V \text{ subject to } \begin{cases} \|V\| = 1 \\ V^T V = 0 \end{cases}$$

- $v$  is (again) an eigenvector of  $S$

$$Z(v, \lambda, M) = v^T S v - \lambda (v^T v - 1) - M v^T v,$$

$$\nabla_v Z(v, \lambda, M) = 2 S v - 2 \lambda v - M v_1 = 0 \mid \cdot v_1$$

$$2 \underbrace{v_1^T S v}_0 - 2 \lambda \underbrace{v_1^T v}_0 - M \underbrace{v_1^T v_1}_1 = -M = 0$$

So,  $Sv = \lambda v$  and  $v$  and  $\lambda$  are e.vec/e.val

- $v$  corresponds to the second largest e.val.

$$v^T S v = \lambda \|v\|^2 = \lambda \rightarrow \max \text{ and } v^T v_1 = 0$$

Both are satisfied by  $v = v_2$ .

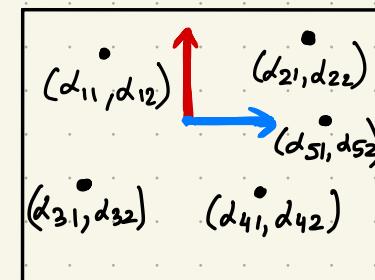
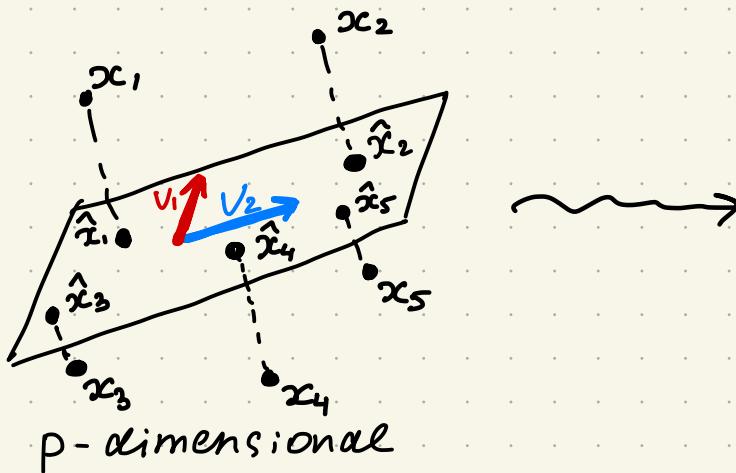
Terminology:  $v_2$  is the direction of  $PC_2$

$z_2 = Xv_2$  is the score vector,  $\lambda_2$  is variance explained.

$$V_{(2)} = \begin{pmatrix} 1 & 1 \\ v_1 & v_2 \\ 1 & 1 \end{pmatrix} \in \mathbb{R}^{p \times 2}$$

$$Z_{(2)} = \begin{pmatrix} 1 & 1 \\ z_1 & z_2 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} d_{11} & d_{12} \\ \dots & \dots \\ d_{n1} & d_{n2} \end{pmatrix} = X V_{(2)} \in \mathbb{R}^{n \times 2}$$

$$\hat{X}_{(2)} = \begin{pmatrix} -\hat{x}_1^T \\ \dots \\ -\hat{x}_n^T \end{pmatrix} = X V_{(2)} V_{(2)}^T \in \mathbb{R}^{n \times p}$$
$$(P_{V_{(2)}} || x^T)^T$$



2-dimensional

## PCA: the r-th component (PC<sub>r</sub>)

$$\text{maximize } V^T S V \text{ subject to } \begin{cases} \|V\| = 1, \\ V^T V = \dots = V^T V_{r-1} = 0 \end{cases}$$

The resulting **directions**  $v_1, \dots, v_r$  are the eigenvectors of  $S$  corresponding to the largest  $r$  eigenvalues. The **scores** are  $z_i = X v_i$  and the **variances explained** are  $d_i = \text{Var}(z_i)$

$$V_{(r)} = \begin{pmatrix} | & | & | \\ v_1 & \dots & v_r \\ | & | & | \end{pmatrix} \in \mathbb{R}^{n \times r} \text{ the basis}$$

$$Z_{(r)} = \begin{pmatrix} | & | & | \\ z_1 & \dots & z_r \\ | & | & | \end{pmatrix} = \begin{pmatrix} d_{11} & \dots & d_{1r} \\ \vdots & \ddots & \vdots \\ d_{nr} & \dots & d_{rr} \end{pmatrix} = X V_{(r)} \in \mathbb{R}^{n \times r} \text{ projections in } \mathbb{R}^n$$

$$\hat{X}_{(r)} = \begin{pmatrix} -x_1^T & - \\ \vdots & \vdots \\ -x_r^T & - \end{pmatrix} = X V_{(r)} V_{(r)}^T \in \mathbb{R}^{n \times p} \text{ projections in } \mathbb{R}^p$$

## PCA via SVD

$$X \text{ is centered}, \quad X = UDV^T = \begin{matrix} U \\ n \times p \end{matrix} \begin{matrix} D \\ p \times p \end{matrix} \begin{matrix} V^T \\ p \times p \end{matrix}$$

$$S = \frac{X^T X}{n-1} = V \frac{D^2}{n-1} V^T$$

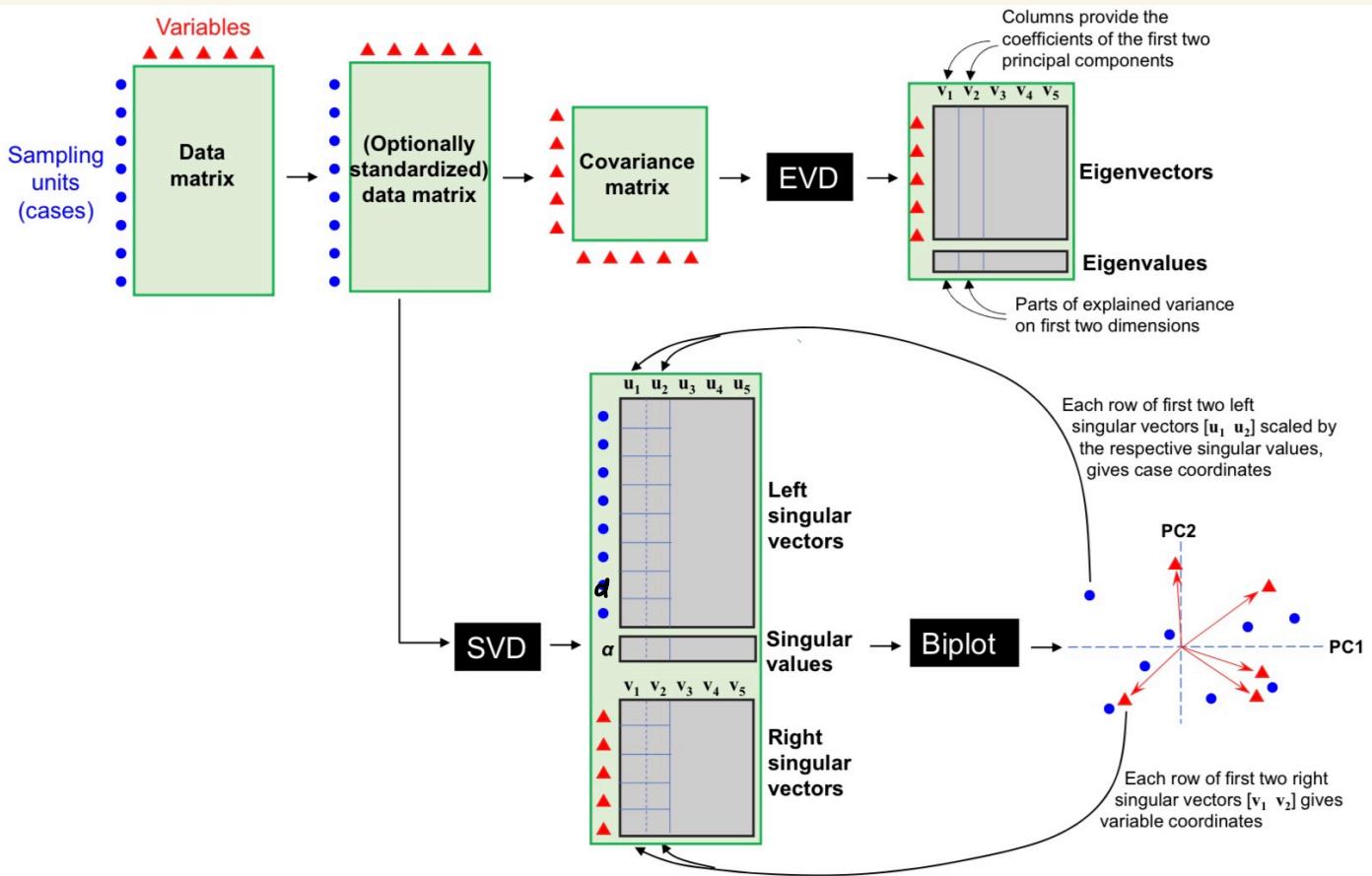
- $V_{(r)}$  are the top- $r$  right singular vectors
- the scores  $Z_{(r)} = X V_{(r)} = U_{(r)} D_{(r)} = \begin{pmatrix} 1 & \dots & 1 \\ u_1 & \dots & u_r \end{pmatrix} \begin{pmatrix} d_1 & \dots & d_r \end{pmatrix}$

$$V^T V_{(r)} = (V_{(r)} (V_{(r)})^\perp)^T V_{(r)} = \begin{pmatrix} V_{(r)}^T V_{(r)} \\ (V_{(r)})^\perp V_{(r)} \end{pmatrix} = \begin{pmatrix} I \\ 0 \end{pmatrix} \in \mathbb{R}^{p \times r}$$

$$\begin{aligned} X V_{(r)} &= U D V^T V_{(r)} = U D \begin{pmatrix} I \\ 0 \end{pmatrix} = U \begin{pmatrix} d_1 & \dots & d_n \end{pmatrix} \begin{pmatrix} I \\ 0 \end{pmatrix} = U \begin{pmatrix} D_{(r)} \\ 0 \end{pmatrix} \\ &= \underbrace{\left( U_{(r)} \underbrace{U_*}_{p-r} \right)}_{r} \begin{pmatrix} D_{(r)} \\ 0 \end{pmatrix} = U_{(r)} D_{(r)} \end{aligned}$$

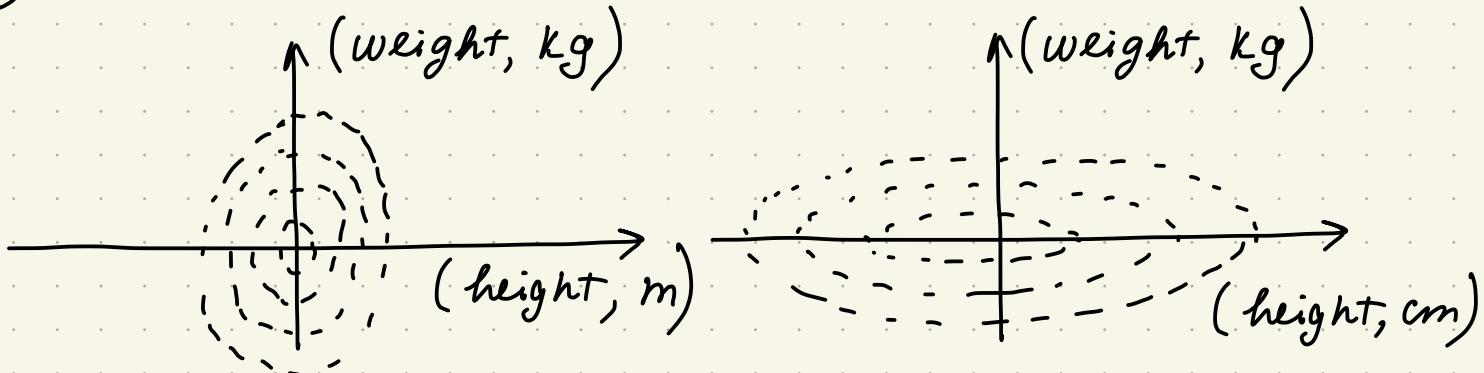
- the projections  $\hat{X}_{(r)} = U_{(r)} D_{(r)} V_{(r)}^T = SVD_r(X)$

$$\hat{X}_{(r)} = d_1 u_1 v_1^T + \dots + d_r u_r v_r^T + d_{r+1} \cancel{u_{r+1} v_{r+1}^T} + \dots + d_p u_p v_p^T \simeq X$$



## Practical aspects

- ① PCA is sensitive to the units



Sometimes you need to **standardize** columns of the data

$$X = \begin{pmatrix} f_1 & \dots & f_p \end{pmatrix} \quad \rightsquigarrow \quad \tilde{X} = \begin{pmatrix} 1 & \dots & 1 \\ \frac{f_1 - \bar{f}_1}{\text{std}(f_1)} & \dots & \frac{f_p - \bar{f}_p}{\text{std}(f_p)} \end{pmatrix}$$

## ② How to pick $r$ ?

- For visualization, use  $r=2$  or  $3$
- For dimension reduction, compute the proportion of variance explained by each PC.

$$\begin{matrix} X \\ n \times p \end{matrix}$$

$$\begin{matrix} S \\ p \times p \end{matrix}$$

$$VE(PC_1) = d_1$$

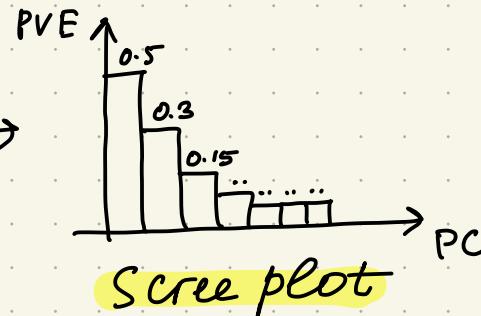
$$VE(PC_p) = d_p$$

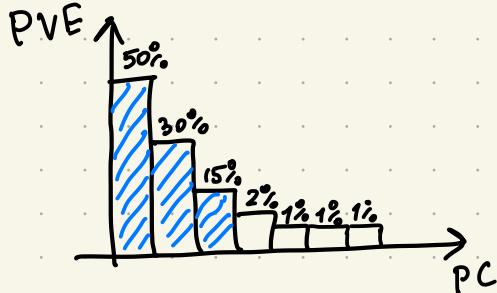
$$\text{Total} = d_1 + \dots + d_p$$

$$PVE(PC_1) = \frac{d_1}{d_1 + \dots + d_p}$$

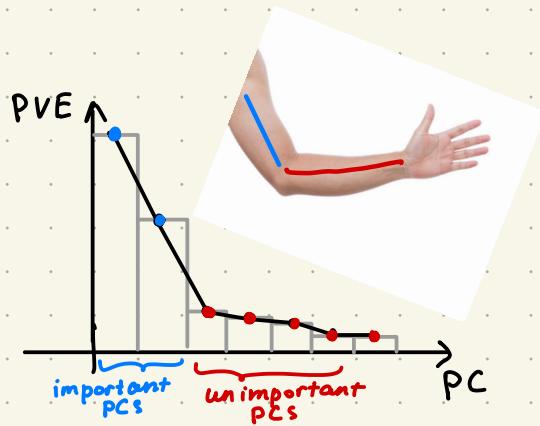
---

$$PVE(PC_p) = \frac{d_p}{d_1 + \dots + d_p}$$



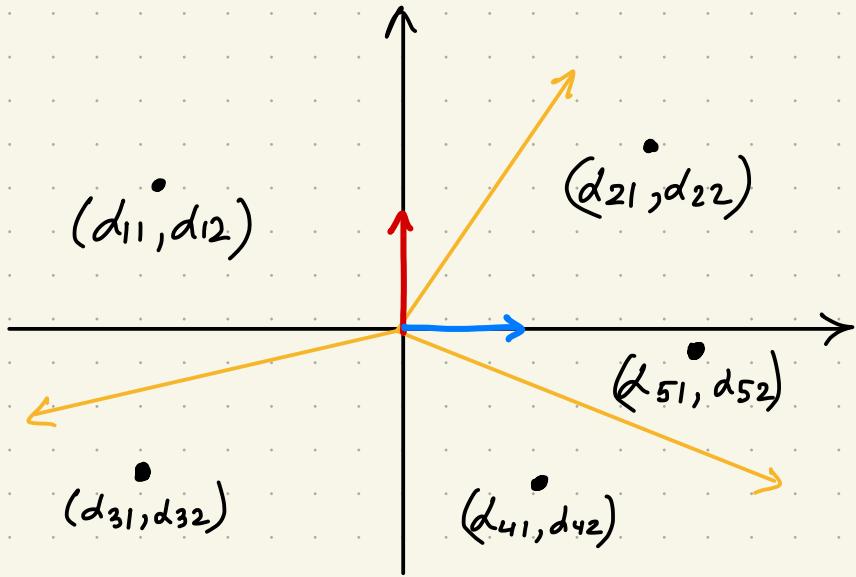


- Select  $r$  such that  $PC_1, \dots, PC_r$  explain, say,  $\geq 90\%$  of the total variance.



- Select  $r$  such that the remaining  $n-k$  PVE values are small (elbow method)

## ② Visualization (2D) : Biplot



In orange : the projections of feature axes on the  $(PC_1, PC_2)$  space

$$\begin{aligned} e_2 &= (0, 1, 0) \\ e_1 &= (1, 0, 0) \\ e_3 &= (0, 0, 1) \end{aligned}$$

$$\begin{pmatrix} -e_1^T & - \\ -e_2^T & - \\ -e_3^T & - \end{pmatrix} = I$$

projections are

$$\begin{aligned} I \cdot V_{(2)} &= \left( V_{(2)} \right) \\ &= \left( \underline{\underline{\underline{\quad}}} \right) \end{aligned}$$

### ③ PCA Complexity

Input:  $X \in \mathbb{R}^{n \times p}$

Step 1: Column-centering: compute  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$O(np)$  and  $x_i \rightarrow x_i - \bar{x}$  for all  $i = 1, \dots, n$

Step 2: Compute  $S = \frac{1}{n-1} X X^T$

$$\begin{matrix} n \\ \text{---} \\ p \end{matrix} \cdot \begin{matrix} p \\ \text{---} \\ n \end{matrix} = \begin{matrix} p \\ \text{---} \\ n \end{matrix}$$

$O(np^2)$

Step 3 Compute eigen decomposition of  $S \in \mathbb{R}^{p \times p}$

$O(p^3)$

Output: directions  $v_1, \dots, v_p$ , variance  $\lambda_1, \dots, \lambda_p$

Total Complexity:  $O(np^2 + p^3)$

Often, we need only  $r$  eigenvectors of  $S \in \mathbb{R}^{p \times p}$

## Power iteration

Find  $\lambda_1$  of  $S \in \mathbb{R}^{P \times P}$  assuming  $\lambda_1 > \lambda_2 \geq \lambda_3 \dots \geq \lambda_P \geq 0$

Input:  $S \in \mathbb{R}^{P \times P}$  and arbitrary  $v \in \mathbb{R}^P$

$$\underline{\text{Step 1}} \quad \tilde{v} = S v$$

$$\underline{\text{Step 2}} \quad v = \frac{\tilde{v}}{\|\tilde{v}\|}$$

} repeat until  $v$  converges

Output:  $v$

Why it works?

Recall  $S^k = U \begin{pmatrix} \lambda_1^k & & \\ & \ddots & \\ & & \lambda_P^k \end{pmatrix} U^T$ . Assume  $V = \sum_{i=1}^P \lambda_i v_i$

$$\text{Then } S^k V = \sum_{i=1}^k \lambda_i^k \lambda_i v_i = \lambda_1^k \sum_{i=1}^k \left(\frac{\lambda_i}{\lambda_1}\right)^k \lambda_i v_i$$

If  $i=1$  then  $\frac{\lambda_i}{\lambda_1} = 1$ , if  $i>1$  then  $\frac{\lambda_i}{\lambda_1} < 1$  and  $\left(\frac{\lambda_i}{\lambda_1}\right)_k \xrightarrow{k \rightarrow \infty} 0$

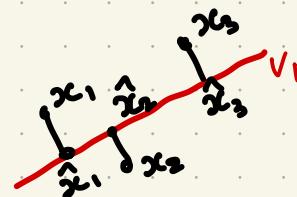
$$S^k V \approx \lambda_1^k \lambda_1 v_1 \quad \text{and} \quad \frac{S^k V}{\|S^k V\|} = \frac{\lambda_1^k \lambda_1 v_1}{\lambda_1^k \lambda_1} = v_1$$

## Application to PCA :

1) Use power iteration to find  $v_1 \Rightarrow d_1 = v_1^T S v_1$

2)  $\tilde{X} = X(I - v_1 v_1^T)$

|  $\hat{X} = X - \tilde{X} = X - X v_1 v_1^T$



3)  $\tilde{S} = \frac{1}{n-1} \hat{X}^T \hat{X} = S - d_1 v_1 v_1^T$

|  $\begin{aligned} \tilde{S} &= (I - v_1 v_1^T) S (I - v_1 v_1^T) = S - v_1 v_1^T S - S v_1 v_1^T + v_1 v_1^T S v_1 v_1^T \\ &= S - \lambda_1 v_1 v_1^T - \lambda_1 v_1 v_1^T + \lambda_1 v_1 v_1^T = S - \lambda_1 v_1 v_1^T \end{aligned}$

4) Use power iteration to find  $\tilde{v}_1$  of  $\tilde{S} \Rightarrow$

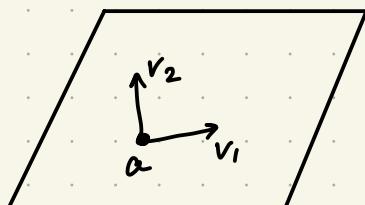
$$v_2 = \tilde{v}_1 \text{ and } d_2 = v_2^T S v_2$$

## 2nd view of PCA: minimum reconstruction error

Given  $n$  points  $x_1, \dots, x_n \in \mathbb{R}^P$  find the "best" approximation by an  $r$ -dimensional plane.

To construct a plane we need:

- a point on the plane  $a \in \mathbb{R}^P$
- a basis  $v_1, \dots, v_r \in \mathbb{R}^P$ , assume  $V^T V = I$

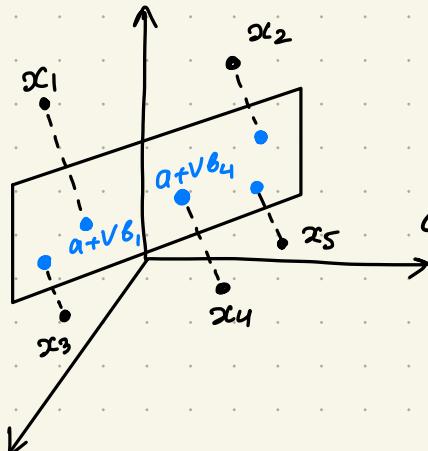


Any point on the plane is:

$$a + \beta_1 v_1 + \dots + \beta_r v_r \text{ for some } \beta_1, \dots, \beta_r$$

Equivalently,  $a + V\beta$ , where  $V = (v_1 \dots v_r) \in \mathbb{R}^{P \times r}$

$$\beta = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_r \end{pmatrix} \in \mathbb{R}^r$$



The PCA goal

$$\text{minimize } \sum_{i=1}^n \|x_i - (a + Vb_i)\|^2 = \ell(a, V, b_i)$$

- If  $a$  and  $V$  are fixed

$$b_i = V^T(x_i - a)$$

$$\text{1 regression } b_i = (V^T V)^{-1} V^T (x_i - a)$$

- $\ell(a, V) = \sum_{i=1}^n \|(\mathbb{I} - VV^T)(x_i - a)\|^2$

minimized at  $a = \bar{x} + VV^T c$  for any  $c \in \mathbb{R}^p$

$$\begin{aligned} \|x_i - (a + VV^T(x_i - a))\|^2 &= \|(\mathbb{I} - VV^T)(x_i - a)\|^2 = \\ &= (x_i - a)^T (\mathbb{I} - VV^T)^2 (x_i - a) \end{aligned}$$

$$\nabla_a \ell(a, V) = 2 \sum_{i=1}^n (\mathbb{I} - VV^T)(x_i - a) = 2n(\mathbb{I} - VV^T)(\bar{x} - a) = 0$$

$$\bar{x} - a = VV^T c \text{ for any } c \in \mathbb{R}^p, \text{ e.g. } c = 0.$$

- $\ell(V) = \sum_{i=1}^n \| (I - VV^T)(x_i - \bar{x}) \|^2$

minimized at  $V = V_{(r)} = (v_1, \dots, v_r)$  top r evecs of S.

$$\begin{aligned}\ell(V) &= \sum_{i=1}^n \text{tr}[(x_i - \bar{x})^T (I - VV^T)(x_i - \bar{x})] = \\ &= \text{tr}[(I - VV^T) \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T] = \\ &= n \text{tr}[(I - VV^T) S] = -n \text{tr}(V^T S V) + \dots\end{aligned}$$

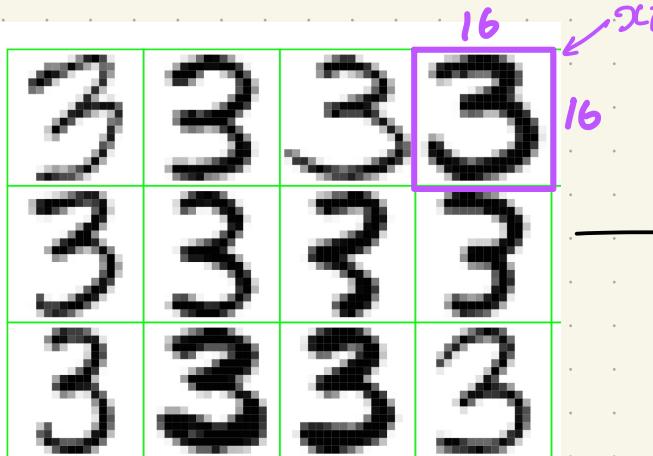
maximize  $\underset{V \in \mathbb{R}^{p \times q}}{\text{tr}(V^T S V)}$  subject to  $V^T V = I$

- $\hat{x}_i = a + VB_i = \bar{x} + V_{(r)} V_{(r)}^T (x_i - \bar{x})$

For centered data,  $\hat{x}_i = V_{(r)} V_{(r)}^T x_i \Leftrightarrow \hat{X} = X \cdot V_{(r)} V_{(r)}^T$

## PCA on images

(from ESLII)



$$X = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} \text{#of images}$$

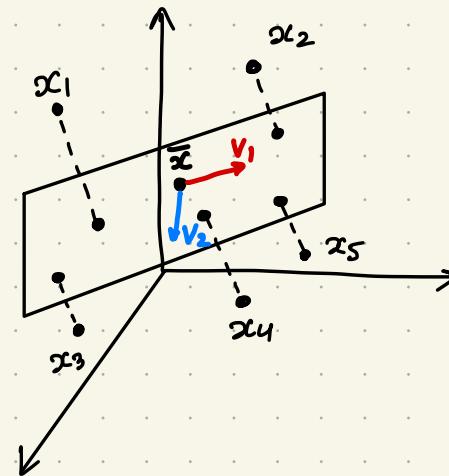
256 pixels

any point in  $(PC_1, PC_2)$  plane is

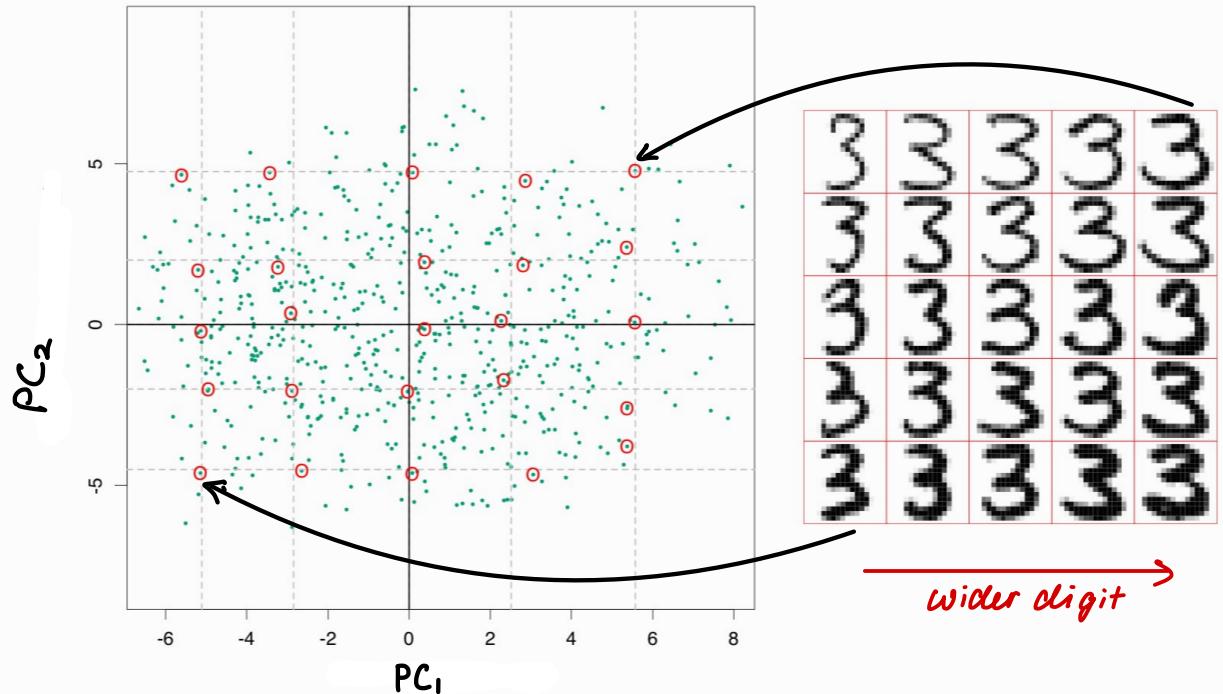
$$\bar{x} + \beta_1 \cdot v_1 + \beta_2 \cdot v_2$$

$$\bar{x} + \beta_1 \cdot v_1 + \beta_2 \cdot v_2$$

"digit width"      "line thickness"

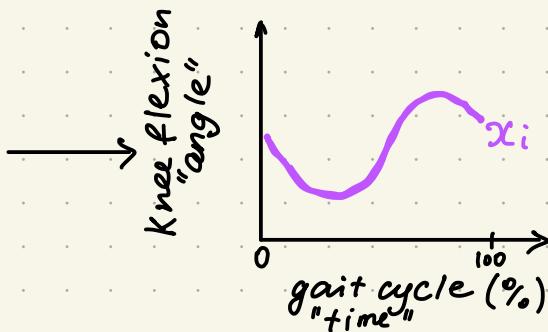
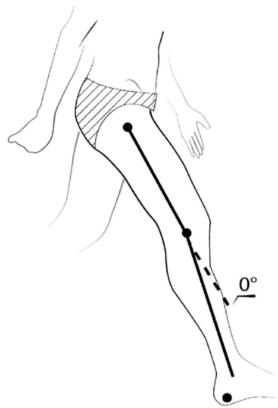


Each point in the biplot corresponds to an image



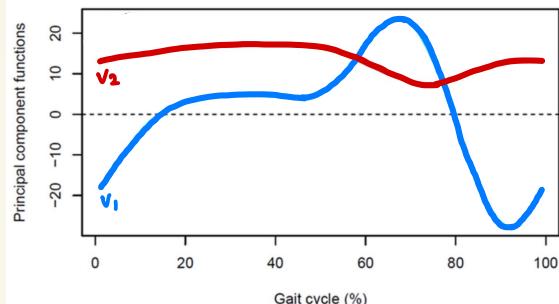
# PCA on functional data (from PCA by Greenacre et al)

Knee angle

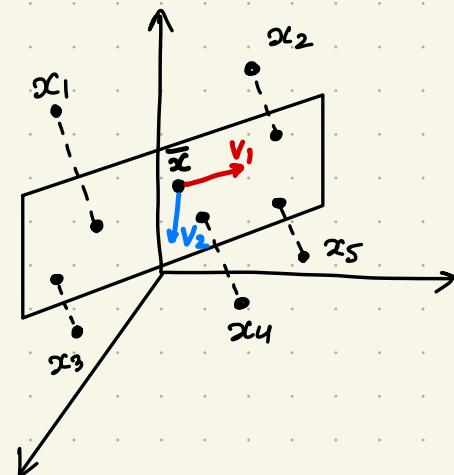


→  $X = \begin{pmatrix} & 100 \\ & \vdots \\ & \text{purple line} \\ & \vdots \\ & 100 \end{pmatrix}_{\text{person}}_{\text{time}}$

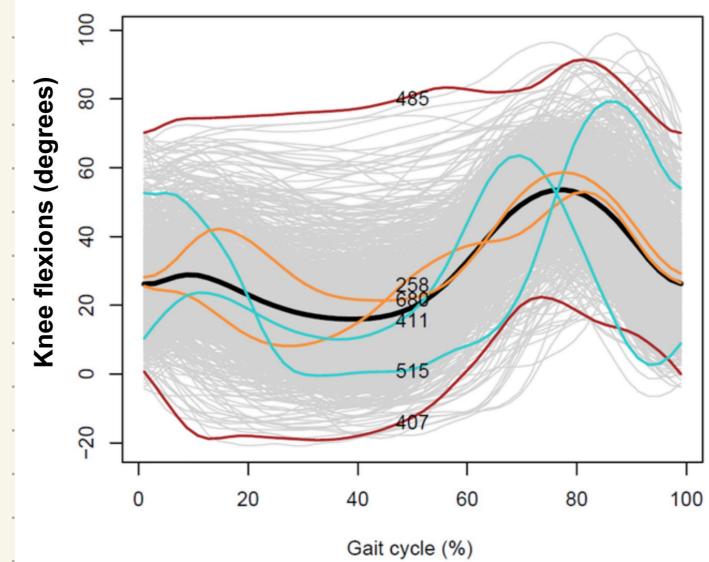
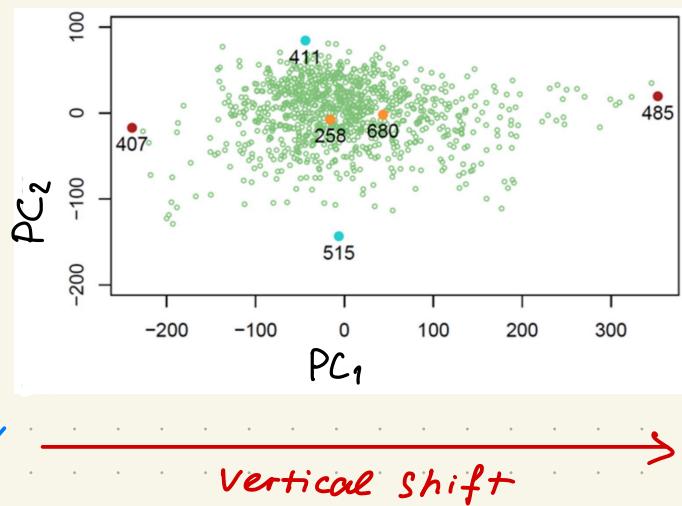
any point in  $(PC_1, PC_2)$  plane is



$$\bar{x} + \beta_1 \cdot v_1 \text{ "Size"} + \beta_2 \cdot v_2 \text{ "Shape"}$$



horiz.  
phase  
shift



Black thick line is  $\bar{x}$ .

## Low-dimensional data

If  $x_1, \dots, x_n \in \mathbb{R}^p$  belong to an  $r$ -dim plane  
then  $X = \begin{pmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{pmatrix} \in \mathbb{R}^{n \times p}$  has rank  $r$  (at most).

There is basis  $v_1, \dots, v_r \in \mathbb{R}^p$  such that  
each  $x_i = \beta_{i1}v_1 + \dots + \beta_{ir}v_r$  for some  $\beta_{i1}, \dots, \beta_{ir}$   
 $= \begin{pmatrix} | & | & | \\ v_1 & \dots & v_r \\ | & | & | \end{pmatrix} \cdot b_i$  where  $b_i \in \mathbb{R}^r$ .

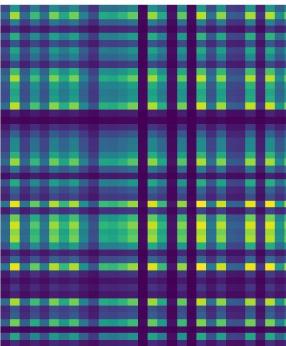
Thus  $\begin{pmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{pmatrix} = \begin{pmatrix} -\beta_1^T - \\ \vdots \\ -\beta_n^T - \end{pmatrix} \begin{pmatrix} -v_1^T - \\ \vdots \\ -v_r^T - \end{pmatrix} = B V^T$

$r=1$

$$X = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_n \\ \hline n \times 1 \end{pmatrix}_{1 \times p} V_i^T = \begin{pmatrix} \beta_{11} & V_i^T \\ \vdots & \vdots \\ \beta_{n1} & V_i^T \end{pmatrix}$$

| . —

$$\begin{aligned} X &= \\ &\text{50} \times 30 \\ \text{rank}(x) &= 1 \end{aligned}$$

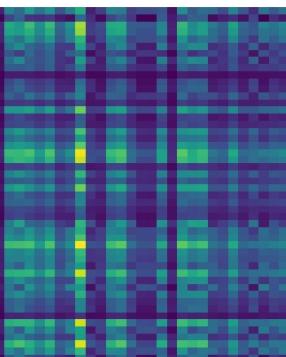


$r=2$

$$X = \begin{pmatrix} -\beta_1^T \\ \vdots \\ -\beta_n^T \\ \hline n \times 2 \end{pmatrix}_{2 \times p} \begin{pmatrix} -V_1^T \\ \vdots \\ -V_2^T \\ \hline 2 \times p \end{pmatrix} = \begin{pmatrix} \beta_{11} V_1^T + \beta_{12} V_2^T \\ \vdots \\ \beta_{n1} V_1^T + \beta_{n2} V_2^T \end{pmatrix}$$

|| . ==

$$\begin{aligned} X &= \\ &\text{50} \times 30 \\ \text{rank}(x) &= 2 \end{aligned}$$

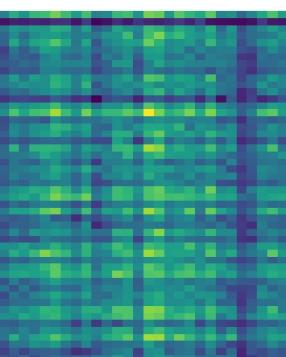


$r=10$

$$X = \begin{pmatrix} -\beta_1^T \\ \vdots \\ -\beta_n^T \\ \hline n \times 10 \end{pmatrix}_{10 \times p} \begin{pmatrix} -V_1^T \\ \vdots \\ -V_{10}^T \\ \hline 10 \times p \end{pmatrix} = \begin{pmatrix} \beta_{11} V_1^T + \dots + \beta_{110} V_{10}^T \\ \vdots \\ \beta_{n1} V_1^T + \dots + \beta_{n10} V_{10}^T \end{pmatrix}$$

|||||        =====

$$\begin{aligned} X &= \\ &\text{50} \times 30 \\ \text{rank}(x) &= 10 \end{aligned}$$



### 3rd view of PCA: low-rank approximation

Given  $X \in \mathbb{R}^{n \times p}$  find the "best" approximation by  $\hat{X} \in \mathbb{R}^{n \times p}$  with rank of  $\hat{X}$  equal to  $r$ .

The approximation error

$$\sum_{i=1}^n \sum_{j=1}^p (x_{ij} - \hat{x}_{ij})^2 = \|X - \hat{X}\|_F^2$$

We aim to solve:

$$\underset{\hat{X}}{\text{minimize}} \|X - \hat{X}\|_F^2 \quad \text{subject to } \text{rank}(\hat{X}) = r$$

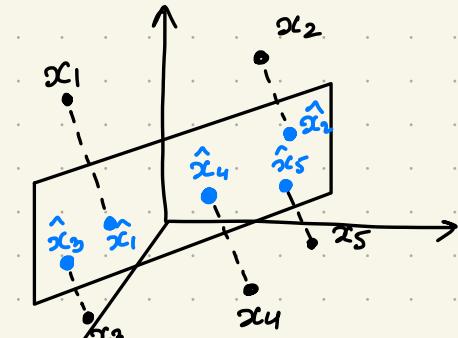
- We can decompose  $\hat{X} = \underset{n \times r}{Q} \underset{r \times p}{A}$  with  $Q^T Q = I$

| rank factorization + QR  $\hat{X} = \underset{n \times r}{Q} \underset{r \times r}{F} \underset{r \times p}{F^T} = \underset{n \times r}{Q} (\underset{r \times r}{R} \underset{r \times p}{F^T}) = \underset{n \times r}{Q} A$

- given  $Q$  find  $A$  as  $A = Q^T \hat{X}$

|  $\|X - QA\|_F^2 = \sum_{j=1}^p \|x_j - Qa_j\|_F^2$  where  $X = \begin{pmatrix} x_1 & \dots & x_p \end{pmatrix}$ ,  $A = \begin{pmatrix} a_1 & \dots & a_p \end{pmatrix}$

| It's a regression problem, so  $A = (Q^T Q)^{-1} Q^T X$



- minimizing  $\|X - QQ^T X\|_F^2$  w.r.t. orthogonal  $Q$   
is equivalent to maximizing  $\text{tr}(\mathcal{D}^2 H)$   
w.r.t. Orthogonal projection  $H \in \mathbb{R}^{n \times n}$  onto  $r$ -dim space.

Here  $\mathcal{D}^2 = \begin{bmatrix} d_1^2 & & 0 \\ & \ddots & \\ 0 & & 0 \end{bmatrix}$  where  $d_1, \dots, d_p$  are s. values of  $X$ .

$$\begin{aligned} \|X - QQ^T X\|_F^2 &= \text{tr}(X^T X) - 2 \text{tr}(X^T Q Q^T X) + \text{tr}(X^T Q Q^T Q Q^T X) = \\ &= -\text{tr}(Q^T X X^T Q) + \dots \end{aligned}$$

$XX^T = \underbrace{U}_{n \times n} \mathcal{D}^2 U^T$  then the goal of minimizing

$$\text{tr}(Q^T X X^T Q) = \text{tr}(\underbrace{Q^T}_{Q_*^T} \underbrace{U}_{\mathcal{D}^2} \underbrace{U^T}_{Q_*} Q) = \text{tr}(Q_*^T \mathcal{D}^2 Q_*) = \text{tr}(\mathcal{D}^2 \underbrace{Q_*^T Q_*}_{H})$$

Here  $Q_* \in \mathbb{R}^{n \times r}$  and  $Q_*^T Q_* = I$  and  $H$  is projection onto  $Q_*$ .

- Optimal  $Q$  is  $Q = U_{(r)}$ , where  $U_{(r)} = [u_1 | \dots | u_r]$  are left s. vectors.

$$\text{tr}(\mathcal{D}^2 H) = \sum_{i=1}^p d_i^2 h_{ii} + \sum_{i=p+1}^n 0 \cdot h_{ii} \quad \text{where } d_1^2 \geq \dots \geq d_p^2$$

$$\text{tr}(H) = \sum_{i=1}^n h_{ii} = r [= \text{tr}(Q_*^T Q_*)] \text{ and } 0 \leq h_{ii} \leq 1$$

Optimal  $H$  would have  $h_{11} = \dots = h_{rr} = 1$

What  $H$  is projection with  $h_{11} = \dots = h_{rr} = 1$ ?

Take  $Q = U_{(r)}$ , then  $Q^* = U^T Q = (U_{(r)} (U_{(r)})^\perp)^T U_{(r)} = \begin{pmatrix} I_r \\ 0 \end{pmatrix}$

$H = Q_*^T Q_* = \begin{pmatrix} I_r \\ 0 \end{pmatrix} (I_r \ 0) = \begin{pmatrix} I_r & 0 \\ 0 & 0 \end{pmatrix}$  follows the condition.

- $A = D_{(r)} V_{(r)}^T = \begin{pmatrix} d_1 & \dots & d_r \end{pmatrix} \begin{pmatrix} -v_1^T \\ \vdots \\ -v_r^T \end{pmatrix}$

$$A = Q^T X = U_{(r)}^T U D V^T = \begin{pmatrix} I_r & 0 \end{pmatrix} D V^T = (D_{(r)}, 0) V^T = D_{(r)} V_{(r)}^T$$

- The solution is  $\hat{x} = U_{(r)} D_{(r)} V_{(r)}^T$ .

Thus  $\hat{x}$  has rank  $r$ .