

## Principle Components

$x_1, \dots, x_n$   $\begin{cases} n \text{ obs} \\ p \text{ variables} \end{cases}$

$$X = \begin{pmatrix} x_{1T} \\ \vdots \\ x_{nT} \end{pmatrix} \xrightarrow[\text{scale}]{\text{centre}} \tilde{X}$$

$n \times p$

**PCA:** Transform original data to obtain  $p$  uncorrelated variables

SVD:  $\tilde{X} = UDV^T$

$\begin{matrix} \downarrow & \downarrow \\ n \times p & p \times p \text{ diagonal} \end{matrix}$

$$(U^T U = I = V^T V)$$

$$\tilde{X}V = UD = \begin{pmatrix} \underline{u}_1 & \dots & \underline{u}_p \end{pmatrix} \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_p \end{pmatrix} = \begin{pmatrix} d_1 \underline{u}_1 & \dots & d_p \underline{u}_p \end{pmatrix}$$

$d_1 \geq d_2 \geq \dots \geq d_p > 0$

$\uparrow \quad \uparrow$   
principal component scores

Covariance matrix

$$S = \frac{1}{n-1} \sum_{i=1}^n (\underline{x}_i - \bar{\underline{x}})(\underline{x}_i - \bar{\underline{x}})^T$$

$$= \frac{1}{n-1} \tilde{X}^T \tilde{X}$$

$$S = V \Lambda V^T \quad \left( \Lambda = \frac{D^2}{n-1} \right)$$

- 1st PC maximizes  $\underline{a}^T S \underline{a}$  s.t.  $\underline{a}^T \underline{a} = 1$   
 $\Rightarrow \underline{a} = \underline{v}_1, \quad V = (\underline{v}_1, \dots, \underline{v}_p)$
- $k$ th PC maximizes  $\underline{a}^T S \underline{a}$  s.t.  $\underline{a}^T \underline{a} = 1$  and  $\underline{a}^T \underline{v}_j = 0$  for  $j=1, \dots, k-1$
- Principal component scores:  
columns of  $\tilde{X}V$
- Principal components loadings:  
columns of  $V = (\underline{v}_1, \dots, \underline{v}_p)$

Hopefully, PC  $\begin{cases} \text{scores} \\ \text{loadings} \end{cases}$  are easily interpretable

Example: Athletics records (Blackboard)

8 records

35 countries

PCA on correlation matrix ( $\tilde{X}$  centred & scaled)

$$\left. \begin{array}{l} \lambda_1 = 0.79 \\ \lambda_2 = 0.11 \\ \lambda_3 = 0.04 \end{array} \right\} \begin{array}{l} \lambda_1 + \lambda_2 = 0.90 \\ \lambda_1 + \lambda_2 + \lambda_3 = 0.94 \end{array}$$

Loadings	PC # 1	PC # 2
100m	-0.33	-0.54
200m	-0.34	-0.47
400m	-0.36	-0.25
800m	-0.38	0
1500m	-0.39	0.13
5000m	-0.37	0.31
10000m	-0.33	0.35
marathon	-0.35	0.47

PC #1 is a measure of overall strength  
PC #2 is a measure of contrast b/w distance & shorter events.

### Biplot

- Very useful plot. Plot 1st PC scores versus 2nd PC scores
- Biplot: adds variable information.

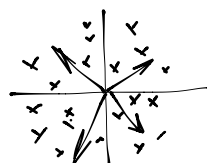
Rank 2 approx to  $\tilde{X}$  ← largest two singular values

$$\tilde{X}_2 = (\underbrace{\underline{u}_1, \underline{u}_2}_{\text{scores}})^{n \times 2} \begin{pmatrix} d_1 & 0 \\ 0 & d_2 \end{pmatrix} \begin{pmatrix} \underline{v}_1^T \\ \underline{v}_2^T \end{pmatrix}^{2 \times p}$$

$$(\tilde{X}\underline{v}_1, \tilde{X}\underline{v}_2) = (d_1 \underline{u}_1, d_2 \underline{u}_2) \quad \text{1st two PC scores}$$

- Plot  $\underline{u}_1$  vs  $\underline{u}_2$  (plot of 1st 2 PC scores)
- $(d_1 \underline{v}_1, d_2 \underline{v}_2)$  represent rows variables of this matrix by vectors

- $p \times 2$
- length of vector represent variance
- angles between vectors represent correlation.



$$\tilde{X}\underline{v} = \begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} \begin{pmatrix} v_1 \\ \vdots \\ v_p \end{pmatrix} \leftarrow \text{loadings}$$

### Multidimensional scaling

Problem:  $\underline{x}_1, \dots, \underline{x}_n$  not necessarily observed  
But ... given an  $n \times n$

← symmetric matrix of distances  $d(\underline{x}_i, \underline{x}_j)$

Find  $\underline{y}_1, \dots, \underline{y}_n$  (low dimensional) such that  $d(\underline{y}_i, \underline{y}_j) \approx d(\underline{x}_i, \underline{x}_j)$

Define  $D = (d_{ij} = d(\underline{x}_i, \underline{x}_j))$

→ centred & scaled

Special case: Assume  $\underline{x}_1, \dots, \underline{x}_n$  known and  $d(\underline{x}_i, \underline{x}_j) = \text{Euclidean distance}$   
 $d(\underline{u}, \underline{v}) = [\sum_{j=1}^p (u_j - v_j)^2]^{\frac{1}{2}}$

Define  $X = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix}^{n \times p}$  and  $B = \underbrace{XX^T}_{n \times n}$

Claim: Can recover  $d_{ij} = d(\underline{x}_i, \underline{x}_j)$  from elements of  $B$ .

$$B = \begin{pmatrix} b_{11} & \dots & b_{1n} \\ \vdots & \ddots & \vdots \\ b_{n1} & \dots & b_{nn} \end{pmatrix}$$

$$d_{ij}^2 = b_{ii} + b_{jj} - 2b_{ij}$$

Moreover, if we rotate each row of  $X$  then  $B$  &  $D$  remain unchanged:  $X \underbrace{OO^T}_{\text{orthogonal}} X^T = B$

$$\begin{aligned} \text{SVD: } X &= U \Lambda^{\frac{1}{2}} V^T \\ X X^T &= U \Lambda^{\frac{1}{2}} \underbrace{V^T V}_{=I} \Lambda^{\frac{1}{2}} U^T \\ &= U \Lambda U^T \end{aligned}$$

$\nearrow$   $\uparrow$   $\nwarrow$   
 $n \times p$   $p \times p$   $p \times n$

$\sim$  result is  $n \times n$

- suggests that best  $r$  dimensional distance preserving transformation is first  $r$  columns of  $U$