

## Principal component analysis (PCA)

Multivariate observations  $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_n$

-  $n$  observations -  $p$  variables

Goals: Try to represent/approximate data by  $r < p$  ( $r \ll p$ ) dimensions,  $r=1, 2, 3$  optimal.

Earlier: Try to find "interesting" projections of the data  $\rightarrow$  projection pursuit.  
i.e. find  $\underline{q}$  s.t.  $\underline{q}^T \underline{x}_1, \dots, \underline{q}^T \underline{x}_n$  look look non-normal (e.g. maximize Kurtosis)

Find structure (outliers, clusters etc)

Trying to find  $\underline{y}_1, \dots, \underline{y}_n$  based on  $\underline{x}_1, \dots, \underline{x}_n$  having similar structure to  $\underline{x}_1, \dots, \underline{x}_n$

Example: Distance measure  $d(\underline{x}_i, \underline{x}_j)$   $\left\{ \begin{array}{l} \text{Euclidean distances} \\ \text{Manhattan} \end{array} \right.$

Try to find  $\underline{y}_1, \dots, \underline{y}_n$  s.t. distance structure is approximately maintained  
 $d(\underline{y}_i, \underline{y}_j) \approx d(\underline{x}_i, \underline{x}_j)$

PCA: Two approaches

- ① Low rank matrix approximation
- ② Projection pursuit  $\rightarrow$  maximize variance

### Low rank matrix approximation

- assume  $p < n$  and define

$$\underset{n \times p}{X} = \begin{pmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_n^T \end{pmatrix} \quad \text{rank}(X) = p$$

- can we approximate  $X$  by a matrix with rank  $r \leq p$  (and how)?

- find  $X^*$  with rank  $r$  to minimize  $\|X - M\|_F^2 = \sum_{i=1}^n \sum_{j=1}^p |X_{ij} - M_{ij}|^2$

$\downarrow$  Frobenius norm

### Singular value decomposition (SVD)

Idea: Write  $\underset{n \times p}{X} = \underset{n \times r}{U} \underset{r \times p}{D} \underset{r \times p}{V}^T$  where  $U$  &  $V$  have orthonormal columns

$U = (\underline{u}_1, \dots, \underline{u}_r)$ ,  $V = (\underline{v}_1, \dots, \underline{v}_r)$  with  $\underline{u}_i^T \underline{u}_j = \underline{v}_i^T \underline{v}_j = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{if } i \neq j \end{cases}$

$D = \begin{pmatrix} d_1 & & 0 \\ 0 & \ddots & d_p \end{pmatrix}$   $\underbrace{d_1 \geq d_2 \geq \dots \geq d_p > 0}_{\text{singular values of } X}$

Now Diagonal  $D^* = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_r & 0 & \dots & 0 \end{pmatrix}$  and define  $X^* = U D^* V^T$

$$\begin{aligned} \text{rank}(X^*) &= \min(\text{rank}(U), \text{rank}(D^*), \text{rank}(V)) \\ &= \text{rank}(D^*) \\ &= r \end{aligned}$$

$X^*$  minimizes  $\|X - M\|_F^2$  over  $M$  with  $\text{rank}(M) = r$

- approximation is good if  $d_{r+1}, d_p$  close to 0.

$$- X^* V = U D^* V^T V = U D^* = \begin{pmatrix} \underline{u}_1 & \dots & \underline{u}_r \end{pmatrix} \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_r & 0 & \dots & 0 \end{pmatrix} = (d_1 \underline{u}_1 \dots d_r \underline{u}_r \ 0 \dots 0)$$

- R function: svd

Example: Exam mark data

- look at centred & scaled data

$$d_1 = 16.64, d_3 = 6.22, d_5 = 4.63 \\ d_2 = 8.02, d_4 = 5.81$$

$$V = \begin{pmatrix} -0.40 & -0.65 \\ -0.43 & -0.44 \\ -0.50 & 0.13 \\ -0.46 & 0.39 \\ -0.44 & 0.47 \end{pmatrix}$$

overall strength

comparison of  
(MECH, VECT)  
to (ALG, ANA, STAT)

From SVD to PCA

Define  $\tilde{X}$  to be the  $n \times p$  matrix of centred (and possibly scaled) data.

$$\text{SVD: } \tilde{X} = U D V^T \\ S = \frac{1}{n-1} \tilde{X}^T \tilde{X} = \frac{1}{n-1} V D U^T U D V^T \\ \text{sample covariance matrix (correlation matrix if vars scaled)}$$

- Columns of  $V$  are the eigenvectors  $v_1, \dots, v_p$  of  $S$  with eigenvalues the diagonal elements of  $\frac{D^2}{n-1} \Rightarrow \frac{d_1^2}{n-1}, \dots, \frac{d_p^2}{n-1}$

Principal components:

Define  $n \times p$  matrix

$$y = \tilde{X} V = \tilde{X} (v_1, \dots, v_p) = (\tilde{X} v_1, \dots, \tilde{X} v_p) \xrightarrow[p \text{ new variables}]{} \text{PC scores}$$

Projection pursuit approach

Data  $x_1, \dots, x_n$  with sample covariance  $S$  or correlation matrix  $R$

- find vector  $a$  with  $\|a\|^2 = a^T a = 1$  maximizing sample variance of  $\{a^T x_i\}$

$$\frac{1}{n-1} \sum_{i=1}^n (a^T x_i - \bar{a}^T \bar{x})^2 = a^T S a$$

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

$$a^T S a = a^T V \Lambda V^T a = \sum_{k=1}^p \lambda_k (v_k^T a)^2$$

$$\Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix}$$

$\Rightarrow$  maximized at  $a = v_1$

Now maximize  $a^T S a$  s.t.  $\|a\|^2 = 1$  and  $a^T v_1 = 0 \Rightarrow$  maximized at  $a = v_2$

Maximize  $a^T S a$  s.t.  $\|a\|^2 = 1$  and  $a^T v_1 = \dots = a^T v_{p-1} = 0$

$\Rightarrow$  maximized at  $a = v_p$

Also note that  $\frac{1}{n-1} \sum_{i=1}^n (v_k^T x_i - \bar{v}_k^T \bar{x})^2 = v_k^T S v_k = \lambda_k$

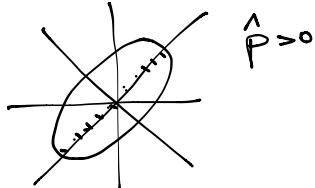
Example:  $p=2$ ,  $R = \begin{pmatrix} 1 & \hat{p} \\ \hat{p} & 1 \end{pmatrix}$ :

eigenvalues

$$\begin{aligned} \text{solve } \det(R - \lambda I) &= 0 \\ \Rightarrow \lambda^2 - 2\lambda + 1 - |\hat{p}|^2 &= 0 \\ \Rightarrow \lambda = 1 \pm |\hat{p}| & \\ \lambda_1 &= 1 + |\hat{p}| \\ \lambda_2 &= 1 - |\hat{p}| \end{aligned}$$

eigenvectors  $\underline{v}_1 = \left( \begin{array}{c} \frac{1}{\sqrt{2}} \\ \text{sgn}(\hat{p}) \frac{1}{\sqrt{2}} \end{array} \right)$   $\underline{v}_2 = \left( \begin{array}{c} \frac{1}{\sqrt{2}} \\ -\text{sgn}(\hat{p}) \frac{1}{\sqrt{2}} \end{array} \right)$   
(loadings)

PCs with  $p=2$  using  $R = \text{sums and differences between 2 variables}$



What happens if  $\hat{p} = 0$ ?  $R = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I$

- any two orthogonal vectors  $\underline{v}_1, \underline{v}_2$  will define PC loadings.

How to use PCs?

①  $R$  vs  $S$  ← centred

↑  
centred, scaled

Generally:  $R$  is the better choice.

- PCs are not invariant to scaling.

- If  $S$  is used in PCA, variables with high variance will tend to dominate.

- If all variables are measures of size (i.e. length, volume, mass) then taking logs and using  $S$  is often OK.

② What PCs to look at? → corresponding to largest eigenvalues or singular values  
- Typically, first few PCs contain most information  
- often reveal interesting structure ⇒ projection pursuit

③ Plot of 1st 2 PC scores is particularly useful.

- biplot (R function: `biplot`): scatterplot of 1st 2 PC scores + information on original variables.

④ PCs with smallest variance can be useful also.

- reveals relationship between variables or a subset of variables.

Examples: Athlete records (Blackboard)

55 countries: USA, CAN, etc, Cook Islands, Western Samoa

8 men's national records: 100m, 200m, 400m, 800m, 1500m, 5000m, 10000m, marathon

- from 1980s

- do PCA on correlation matrix

`princomp(..., cor = T)`

- correlations highest for similar events

$\lambda_1, \dots, \lambda_8$  = variances of PCs

$$\lambda_1 + \dots + \lambda_8 = 8$$

$$\frac{\lambda_1}{8} = 0.79, \frac{\lambda_1 + \lambda_2}{8} = 0.90, \frac{\lambda_1 + \lambda_2 + \lambda_3}{8} = 0.94,$$

Loadings: PC #1 PC #2

100m	-0.33	-0.54	Interpretation?
200m	-0.34	-0.47	
400m	-0.36	-0.25	
800m	-0.38	0.00	
1500m	-0.39	0.13	
5000m	-0.37	0.31	
10000m	-0.33	0.35	
marathon	-0.33	0.47	

- rank countries based on PC scores