

**PS5841**

Data Science in Finance & Insurance

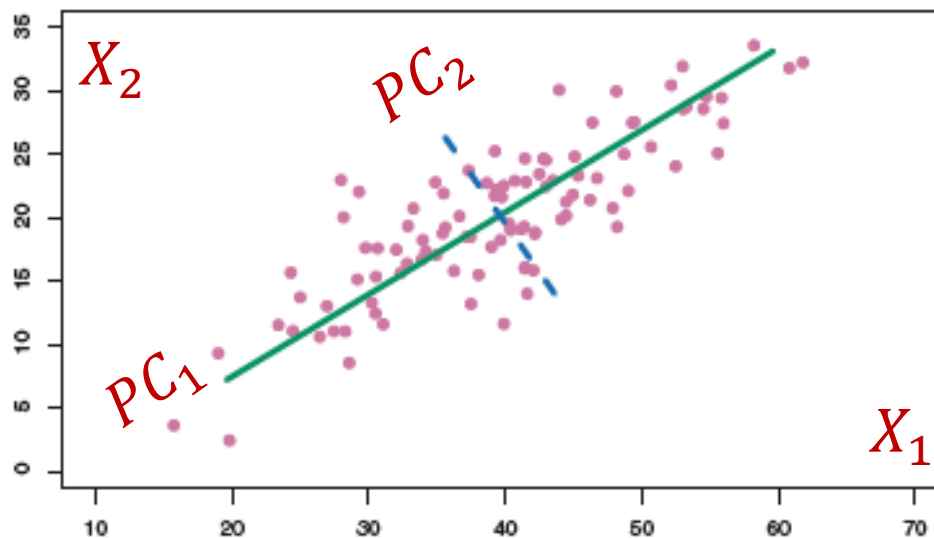
# Dimension Reduction

Yubo Wang

Autumn 2022

# Principal Component Analysis

At most  $\min(n - 1, p)$  principal components



# notations

scores

loadings

$$\mathbf{Z} = \mathbf{X} \mathbf{\Phi}$$

$n \times q \quad n \times p \quad p \times q$

$$(\mathbf{Z}_1, \dots, \mathbf{Z}_q) = (\mathbf{X}_1, \dots, \mathbf{X}_p)(\mathbf{\Phi}_1, \dots, \mathbf{\Phi}_q)$$

$$\begin{bmatrix} z_{11} & \cdots & z_{1q} \\ \vdots & \ddots & \vdots \\ z_{n1} & \cdots & z_{nq} \end{bmatrix} = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \phi_{11} & \cdots & \phi_{1q} \\ \vdots & \ddots & \vdots \\ \phi_{p1} & \cdots & \phi_{pq} \end{bmatrix}$$

$$\mathbf{X}_j = (x_{1j}, \dots, x_{nj})^T, \quad \mathbf{\Sigma} = \text{Var}(\mathbf{X}), \quad \mathbf{\Phi}_k = (\phi_{1k}, \dots, \phi_{pk})^T$$

$$\mathbf{Z}_k = (z_{1k}, \dots, z_{nk})^T = \sum_{j=1}^p \phi_{jk} \mathbf{X}_j = \phi_{1k} \mathbf{X}_1 + \cdots + \phi_{pk} \mathbf{X}_p$$

$$z_{ik} = \sum_{j=1}^p x_{ij} \phi_{jk} = \phi_{1k} x_{i1} + \cdots + \phi_{pk} x_{ip}$$

$$\mathbf{\Phi} \mathbf{\Phi}^T = \mathbf{I}$$

# Principal Components

---

$$\max_{\Phi_k} \text{Var}(\mathbf{Z}_k) = \max_{\Phi_k} \Phi_k^T \Sigma \Phi_k$$

such that

$$\Phi_k^T \Phi_k = 1$$

$$\text{Cov}(\mathbf{Z}_{k'}, \mathbf{Z}_k) = \Phi_{k'}^T \Sigma \Phi_k = 0, \quad k' = k - 1, \dots, 1$$

- Solution:

$\Phi_k = \mathbf{e}_k$ , the  $k$ -th eigenvector of  $\Sigma$

$\text{Var}(\mathbf{Z}_k) = \lambda_k$ , the  $k$ -th eigenvalue of  $\Sigma$

- In practice, use sample variance

$$\hat{\Sigma} = \frac{1}{n-1} \mathbf{X}_{centered}^T \mathbf{X}_{centered}$$

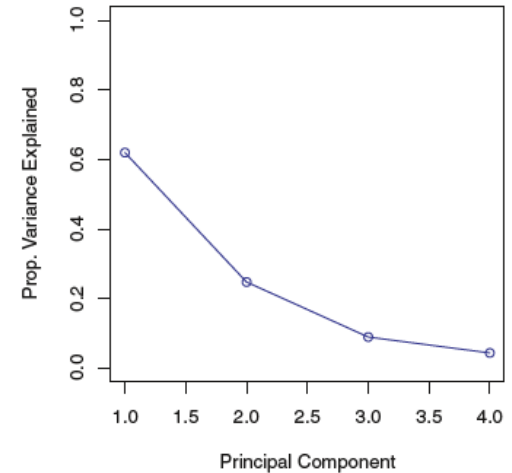
# Proportion of Variance Explained

- Total Variance of  $\mathbf{X}$

$$\text{trace}(\mathbf{\Sigma}) = \sum_{j=1}^p \lambda_j$$

- PVE

$$PVE = \frac{\text{Var}(\mathbf{Z}_k)}{\text{trace}(\mathbf{\Sigma})} = \frac{\lambda_k}{\sum_{j=1}^p \lambda_j}$$



# Uniqueness

---

- Loading vectors are unique up to a sign flip
  - sign flip does not alter the coordinate system
- Score vectors are unique up to a sign flip
  - Variance of  $\mathbf{Z}_k$  and  $-\mathbf{Z}_k$  are the same
- But the right sign may improve interpretability

# Coding (PCA)

---

- R

`base::prcomp()`

- Python

`sklearn.decomposition.PCA()`

# Dimension Reduction

- From  $p$  predictors to  $M < p$  predictors
  - Standardizing the predictors necessary to have predictors on the same scale

$$\begin{aligned}y_i &= \theta_0 + \sum_{m=1}^M \theta_m z_{im} + \varepsilon_i \\&= \theta_0 + \sum_{m=1}^M \theta_m \sum_{j=1}^p x_{ij} \phi_{jm} + \varepsilon_i \\&= \theta_0 + \sum_{m=1}^M \sum_{j=1}^p x_{ij} \phi_{jm} \theta_m + \varepsilon_i \\&= \beta_0 + \sum_{j=1}^p x_{ij} \beta_j + \varepsilon_i\end{aligned}$$

$$\beta_j = \sum_{m=1}^M \phi_{jm} \theta_m$$

$$\beta_0 = \theta_0$$

~ a representation of the original regression  
recall  $\mathbf{X} = \mathbf{Z}\Phi^T$



# Principal Components Regression (PCR)

---

- A small number of PCs may be able to explain most of the variability in data, as well as the relationship with the response (no guarantee)
- Assumption: the directions in which the predictors show the most variation are the direction that are associated with the response
- Not a variable selection method
  - Each PC is a linear combination of all original predictors

# Coding (PCR)

---

- R

`pls::pcr()`

- Python

`sklearn.decomposition.PCA()`

`sklearn.linear_model.LinearRegression()`

# Partial Least Squares Regression (PLS)

---

- $\mathbf{Z}_1 = (z_{11}, \dots, z_{n1})^T = \sum_{j=1}^p \phi_{j1} \mathbf{X}_j$

where  $\phi_{j1}$  is the SLR slope of  $Y$  on  $\mathbf{X}_j$

- $\mathbf{Z}_2 = (z_{12}, \dots, z_{n2})^T = \sum_{j=1}^p \phi_{j2} \mathbf{X}_j$

where  $\phi_{j1}$  is the SLR slope of the residual  
(of  $Y$  on  $\mathbf{Z}_1$ ) on  $\mathbf{X}_j$

- ...

# Coding (PLS)

---

- R

`pls::plsr()`

- Python

`sklearn.decomposition.PLSRegression()`

# That was

