

PS5841

# Data Science in Finance & Insurance

PCA

Yubo Wang

Spring 2021

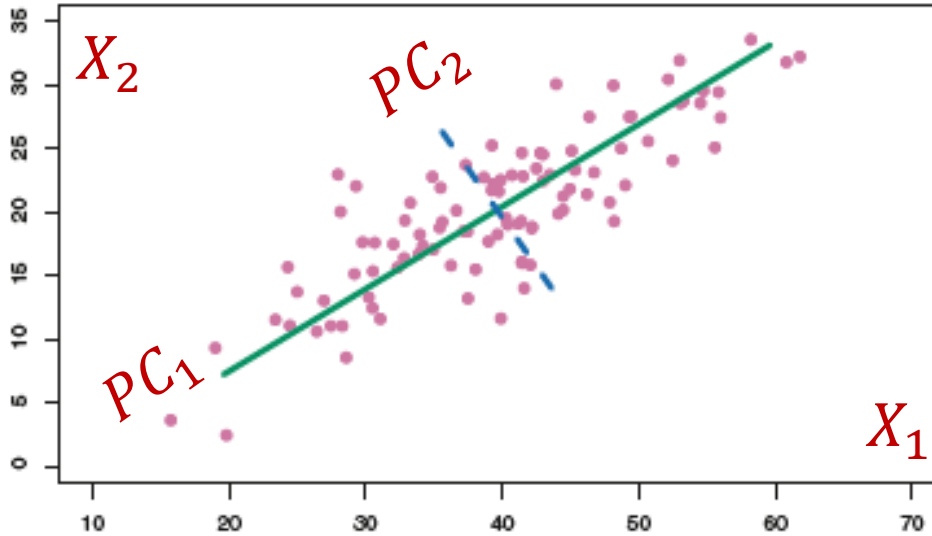
# Unsupervised Learning

- No responses to supervise learning
  - Difficult to access the quality of unsupervised learning
- Feature visualization
- Data pre-processing
- Discover unknown subgroups in data

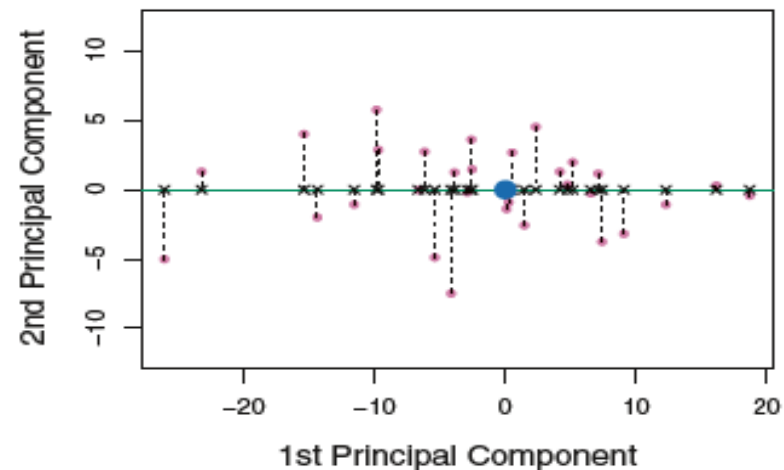
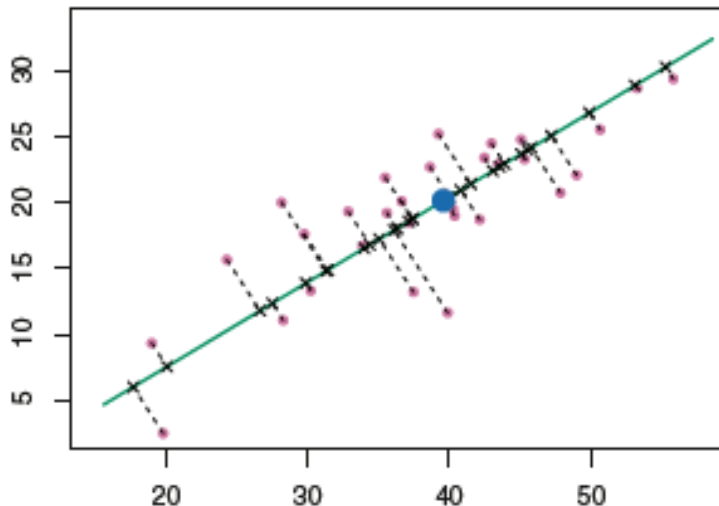
# PCA

- Low(er)-dimensional representation of feature space capturing variation as much as possible
  - Loading vectors: orthogonal unit vectors in feature space with the most variations
  - Score vectors: projections along loading vectors
- Q-dimensional hyperplane that is closest (in terms of Euclidean distance) to the observations.
- At most  $\min(n - 1, p)$  principal components.

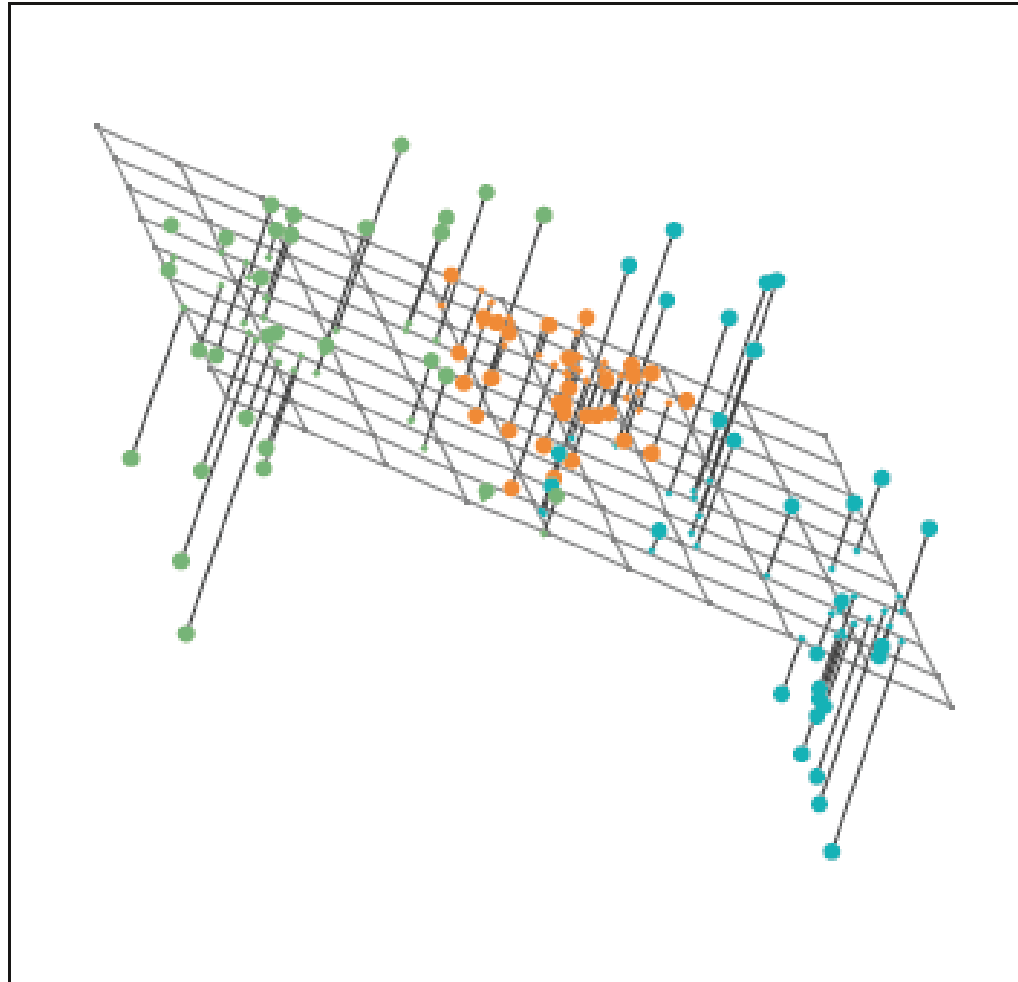
# Example: PC in 2D Data



- Centered “rotation”
- $PC_1$  minimizes SS distance from data to projection onto  $PC_1 \Leftrightarrow$
- $PC_1$  maximizes SS distance from projection onto  $PC_1$  to the center



# Example: PC in 3D Data



# Standardizing Data

- De-mean (zero mean) to make variance calculation more tractable
- Unit variance unless measured in the same units

# PC notation

$$\begin{aligned} \mathbf{Z} &= \mathbf{X} \mathbf{\Phi} \\ n \times q & \quad n \times p \quad p \times q \\ (\mathbf{Z}_1, \dots, \mathbf{Z}_q)^T &= \mathbf{X} (\mathbf{\Phi}_1, \dots, \mathbf{\Phi}_q)^T \\ \begin{bmatrix} z_{11} & \cdots & z_{1q} \\ \vdots & \ddots & \vdots \\ z_{n1} & \cdots & z_{nq} \end{bmatrix} &= \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \phi_{11} & \cdots & \phi_{1q} \\ \vdots & \ddots & \vdots \\ \phi_{p1} & \cdots & \phi_{pq} \end{bmatrix} \\ \mathbf{Z}_k &= \sum_{j=1}^p \phi_{jk} \mathbf{X}_j = \phi_{1k} \mathbf{X}_1 + \cdots + \phi_{pk} \mathbf{X}_p \\ z_{ik} &= \sum_{j=1}^p \phi_{jk} x_{ij} = \phi_{1k} x_{i1} + \cdots + \phi_{pk} x_{ip} \end{aligned}$$

# $PC_1$

- Perform PCA on standardized data (zero mean, plus unit variance unless measured in the same units)
- $\Phi_1$  maximizes

$$\frac{1}{n} \sum_{i=1}^n z_{i1}^2 = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{j1} x_{ij} \right)^2$$

subject to  $\sum_{j=1}^p \phi_{j1}^2 = 1$

where

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



# $PC_k$

- Perform PCA on standardized data (zero mean, plus unit variance unless measured in the same units)
- $\Phi_k$  maximizes

$$\frac{1}{n} \sum_{i=1}^n z_{ik}^2 = \frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{jk} x_{ij} \right)^2$$

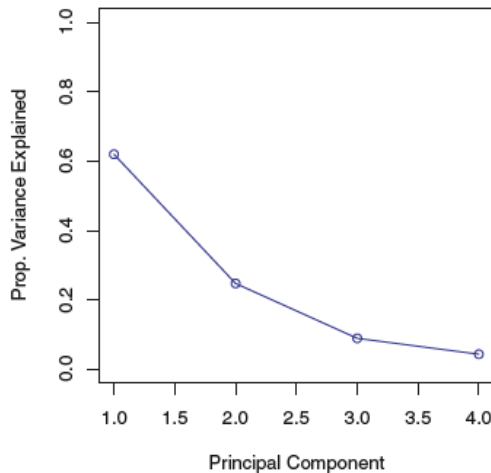
subject to

$$\sum_{j=1}^p \phi_{jk}^2 = 1$$

$\Phi_2$  orthogonal to  $\Phi_1, \dots, \Phi_{k-1}$

# PVE

- Proportion of variance (total variation) explained by  $PC_k$



$$\frac{\frac{1}{n} \sum_{i=1}^n z_{ik}^2}{\sum_{j=1}^p \text{Var}(\mathbf{X}_j)} = \frac{\frac{1}{n} \sum_{i=1}^n \left( \sum_{j=1}^p \phi_{jk} x_{ij} \right)^2}{\sum_{j=1}^p \sum_{i=1}^n x_{ij}^2}$$

That was

