# STAT 3690 Lecture 22

zhiyanggeezhou.github.io

Zhiyang Zhou (zhiyang.zhou@umanitoba.ca)

Mar 23, 2022

## Sample PCA

- Data $\mathbf{X} = [\mathbf{X}_1, \ldots, \mathbf{X}_n]_{n \times p}^\top$

  - Each row $\mathbf{X}_i \overset{\text{iid}}{\sim} (\boldsymbol{\mu}, \boldsymbol{\Sigma})$

- Estimate the loadings $\boldsymbol{w}_j$ through the eigenvectors of sample covariance matrix $\mathbf{S}$ or sample correlation matrix $\hat{\mathbf{R}}$

$$
\hat{\mathbf{R}} = \left[ \begin{array}{ccc} \{\widehat{\text{var}}(X_1)\}^{-1/2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \{\widehat{\text{var}}(X_p)\}^{-1/2} \end{array} \right] \mathbf{S} \left[ \begin{array}{ccc} \{\widehat{\text{var}}(X_1)\}^{-1/2} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \{\widehat{\text{var}}(X_p)\}^{-1/2} \end{array} \right]
$$

- Matrix of scores of the first $s$ principal components

$$
\mathbf{Z} = [Z_{ij}]_{n \times s} = \tilde{\mathbf{X}}\widehat{\mathbf{W}}
$$

  - $\tilde{\mathbf{X}} = [\mathbf{X}_1 - \bar{\mathbf{X}}, \ldots, \mathbf{X}_n - \bar{\mathbf{X}}]_{n \times p}^\top$: row-centered $\mathbf{X}$ (i.e. the sample mean has been subtracted from each row of $\mathbf{X}$)
    * $\bar{\mathbf{X}} = n^{-1} \sum_{i=1}^n \mathbf{X}_i$
  - $\widehat{\mathbf{W}} = [\hat{\boldsymbol{w}}_1, \ldots, \hat{\boldsymbol{w}}_s]_{p \times s}$: $\hat{\boldsymbol{w}}_j$ is the estimate of $\boldsymbol{w}_j$
  - $Z_{ij} = (\mathbf{X}_i - \bar{\mathbf{X}})^\top \hat{\boldsymbol{w}}_j$: the $j$th PC score for the $i$th observation

---

## Geometric interpretation of (sample) PCA

- The definition of PCA as a linear combination that maximises variance is due to H. Hotelling (1933, Journal of Educational Psychology, 24, 417–441).

- PCA was introduced earlier by K. Pearson (1901, Philosophical Magazine, Series 6, 2(11), 559–572) to minimize the overall error in reconstructing data points

$$
(\bar{\mathbf{X}}, \widehat{\mathbf{W}}, \mathbf{Z}_i) = \arg \min_{\boldsymbol{\theta}, \mathbf{A}, \mathbf{B}_i} \sum_{i=1}^n \|\mathbf{X}_i - \boldsymbol{\theta} - \mathbf{A}\mathbf{B}_i\|^2
$$

  - $\mathbf{Z}_i$: the $i$th row of score matrix $\mathbf{Z}$