

High-dimensional data analysis

Lecture 3

Singular Value Decomposition and Principal Component Analysis

Innopolis University

Fall 2018

Similar matrices

- Matrices $A, B \in \mathbb{R}^{d \times d}$ are called similar if $\exists C \in \mathbb{R}^{d \times d}$ such that
 - C is invertible
 - $B = C^{-1}AC$
- Similar matrices have same:
 - Rank
 - Determinant
 - Trace
 - Eigenvalues
 - Frobenius norm ($\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2\right)^{1/2} = \sqrt{\text{tr } A^T A}$)

Spectral Decomposition

- Let $A \in \mathbb{R}^{d \times d}$, symmetric and positive-definite, then
$$A = \Gamma \Lambda \Gamma^T,$$

where

- $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_d)$
 - $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$
 - $\Gamma = [\boldsymbol{\eta}_1 \quad \boldsymbol{\eta}_2 \quad \dots \quad \boldsymbol{\eta}_d]$
 - $\forall i = 1 \dots d: A\boldsymbol{\eta}_i = \lambda_i\boldsymbol{\eta}_i, \|\boldsymbol{\eta}_i\| = 1$
-
- Property:
 - Γ is orthogonal ($\Gamma^T \Gamma = \Gamma \Gamma^T = \mathbf{I}$)

Spectral Decomposition

- Let $A \in \mathbb{R}^{d \times d}$, $r \stackrel{\text{def}}{=} \text{rank } A < d$, then
$$A = \Gamma_r \Lambda_r \Gamma_r^T,$$

where

- $\Lambda_r = \text{diag}(\lambda_1, \dots, \lambda_r) \in \mathbb{R}^{r \times r}$
 - $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > 0, \lambda_{r+1} = \dots = \lambda_d = 0$
 - $\Gamma_r = [\boldsymbol{\eta}_1 \quad \boldsymbol{\eta}_2 \quad \dots \quad \boldsymbol{\eta}_r] \in \mathbb{R}^{d \times r}$
 - $\forall i = 1 \dots d: A\boldsymbol{\eta}_i = \lambda_i \boldsymbol{\eta}_i, \|\boldsymbol{\eta}_i\| = 1$
- Property:
 - Γ is r -orthogonal ($\Gamma_r^T \Gamma_r = \mathbf{I}, \Gamma_r \Gamma_r^T \neq \mathbf{I}$)

Spectral Decomposition

- Let $A \in \mathbb{R}^{d \times d}$, $\text{rank } A = d$, $A = \Gamma \Lambda \Gamma^T$, then
 - $A = \sum_{i=1}^d \lambda_i \boldsymbol{\eta}_i \boldsymbol{\eta}_i^T$
 - $\exists \Theta: A = \Theta^T \Theta$ (e.g., $\Theta = \Lambda^{1/2} \Gamma$)
 - $\forall q \in \mathbb{Q}: A^q = \Gamma \Lambda^q \Gamma^T$

Singular Value Decomposition

- Let $A \in \mathbb{R}^{d \times n}$, $\text{rank } A = r \leq d$,
$$A = UDV^T$$

where

- $D = \text{diag}(d_1, d_2, \dots, d_r) \in \mathbb{R}^{r \times r}$
 - $d_1 \geq d_2 \geq \dots \geq d_r > 0$
- $U = [\mathbf{u}_1 \quad \mathbf{u}_2 \quad \dots \quad \mathbf{u}_r] \in \mathbb{R}^{d \times r}$
- $V = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_r] \in \mathbb{R}^{n \times r}$
 - $\forall i = 1 \dots r: A^T \mathbf{u}_i = d_i \mathbf{v}_i, A \mathbf{v}_i = d_i \mathbf{u}_i$
- Property:
 - $V^T V = U^T U = I$

Decompositions for the Sample Case

- Let $\mathbb{X} \in \mathbb{R}^{d \times n}$, $\mathbb{X} \sim \text{Sam}(\bar{\mathbf{X}}, S)$, $\mathbb{X}_0 \stackrel{\text{def}}{=} \mathbb{X} - \bar{\mathbf{X}}$,
 - $\mathbb{X}_0 = UDV^T$,
 - $S = \Gamma\Lambda\Gamma^T$,
- then:
 - $U = \Gamma$
 - $D^2 = (n - 1)\Lambda$

Principal Component Analysis

- Goals of Multivariate data analysis:
 - understand the structure in the data;**
 - summarize data in simpler ways;**
 - find the relationship between parts of the data;
 - make decisions based on the data.

Principal Component Analysis

- Instead of looking at individual features, let's look at their combinations
- How to choose these combinations?
 - Linear combinations,
 - that best explain the variability of data
- How many combinations to choose?
 - More combinations: higher accuracy
 - Less combinations: faster computations, better interpretability

PCA: Population case

- $\mathbf{X} \sim (\boldsymbol{\mu}, \Sigma)$, $\Sigma = \Gamma \Lambda \Gamma^T$, $r = \text{rank } \Sigma$, then for $k = 1 \dots r$:
 - $W_k = \boldsymbol{\eta}_k^T (\mathbf{X} - \boldsymbol{\mu})$ – k -th principal component score
 - $\mathbf{W}^{(k)} = [W_1, \dots, W_k]^T = \Gamma_k^T (\mathbf{X} - \boldsymbol{\mu})$ – k -th principal component vector
 - $\mathbf{P}_k = \boldsymbol{\eta}_k \boldsymbol{\eta}_k^T (\mathbf{X} - \boldsymbol{\mu}) = W_k \boldsymbol{\eta}_k$ – k -th principal component projection
- First PC – largest eigenvalue – direction in which data has largest variance

PCA: Population case

- Example:

- 2D Multivariate normal distribution

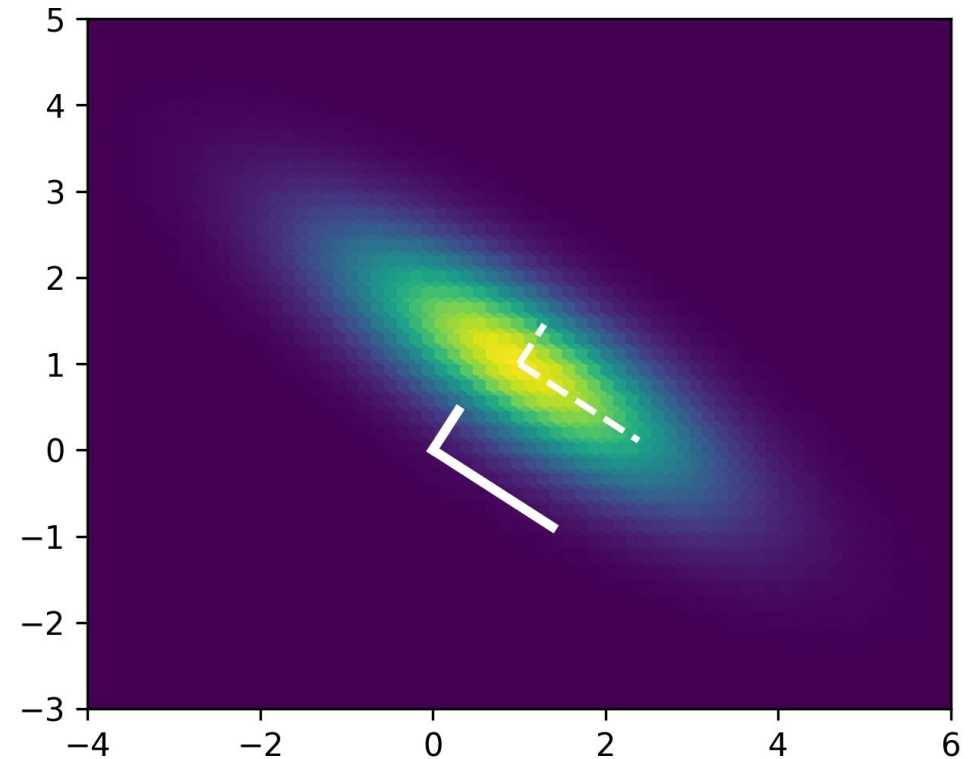
- $\mu = [1 \ 1]^T$

- $\Sigma = \begin{bmatrix} 2 & -1.1 \\ -1.1 & 1 \end{bmatrix}$

- $\Lambda \approx \text{diag}(2.7, 0.3)$

- $\Gamma \approx \begin{bmatrix} 0.84 & 0.54 \\ -0.54 & 0.84 \end{bmatrix}$

- Solid lines – principal components
- Dashed lines – PCs, transposed by μ



PCA: Sample case

- $\mathbb{X} \in \mathbb{R}^{d \times n}$, $\mathbb{X} \sim \text{Sam}(\bar{\mathbf{X}}, S)$, $S = \Gamma \Lambda \Gamma^T$, $r = \text{rank } S$, then for $k = 1 \dots r$:
 - $\mathbf{W}_{\blacksquare k} = \boldsymbol{\eta}_k^T (\mathbb{X} - \bar{\mathbf{X}}) \in \mathbb{R}^{1 \times n}$ – k -th principal component score
 - $\mathbb{W}^{(k)} = [\mathbf{W}_{\blacksquare 1}, \dots, \mathbf{W}_{\blacksquare k}]^T = \Gamma_k^T (\mathbb{X} - \bar{\mathbf{X}}) \in \mathbb{R}^{k \times n}$ – k -th principal component data
 - $\mathbb{P}_{\blacksquare k} = \boldsymbol{\eta}_k \boldsymbol{\eta}_k^T (\mathbb{X} - \bar{\mathbf{X}}) = \boldsymbol{\eta}_k \mathbf{W}_{\blacksquare k} \in \mathbb{R}^{d \times n}$ – k -th principal component projection
- First PC – largest eigenvalue – direction in which data has largest variance
 - $\text{Var}(\boldsymbol{\eta}_k^T \mathbb{X}) = \boldsymbol{\eta}_k^T S \boldsymbol{\eta}_k = \boldsymbol{\eta}_k^T \Gamma \Lambda \Gamma^T \boldsymbol{\eta}_k = \lambda_k$

PCA: Sample case

- Example:

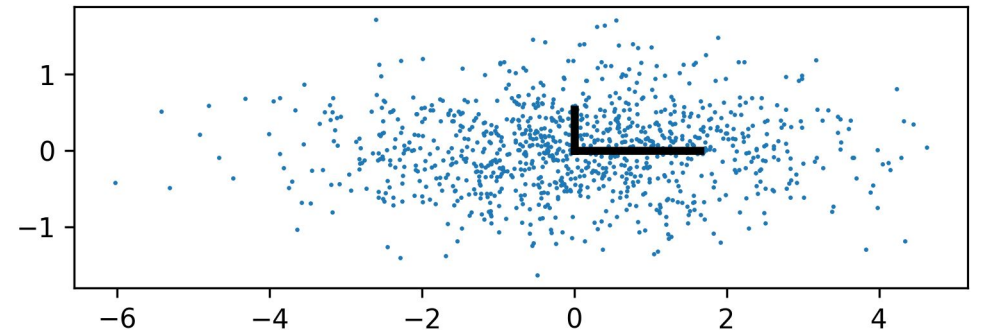
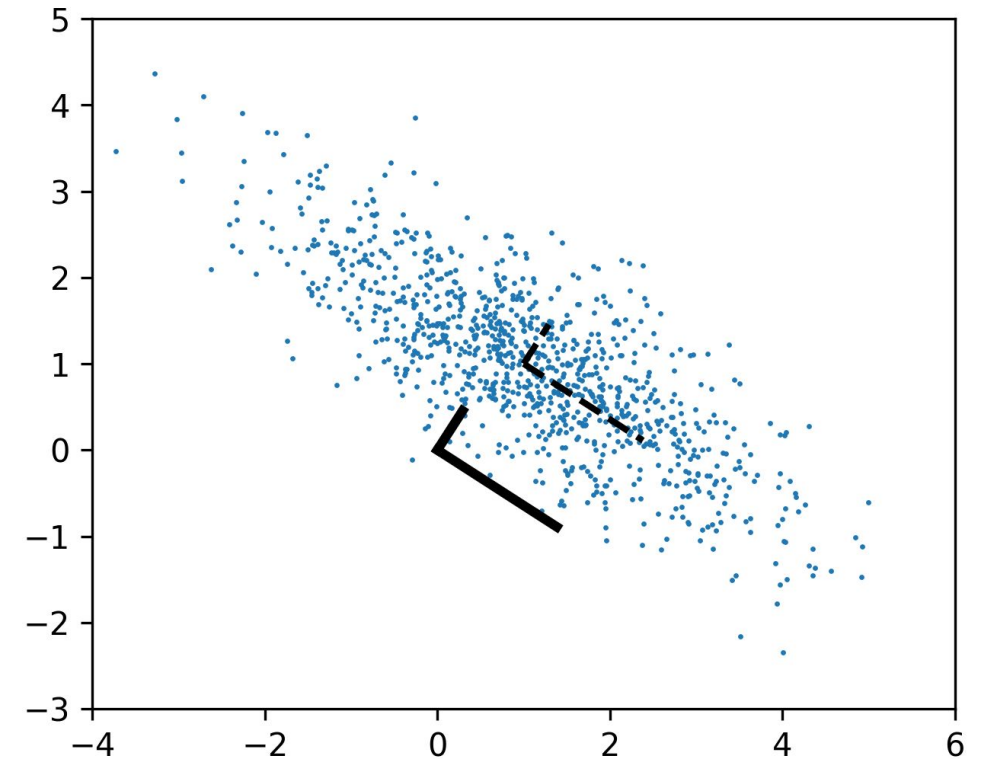
- 2D data

- $\mathbb{X} \sim \text{Sam}\left(\begin{bmatrix} 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 2 & -1.1 \\ -1.1 & 1 \end{bmatrix}\right)$

- Solid lines – principal components

- Dashed lines – PCs, transposed by μ

- Lower figure: $\mathbb{W}^{(2)}$ (2nd PC data)

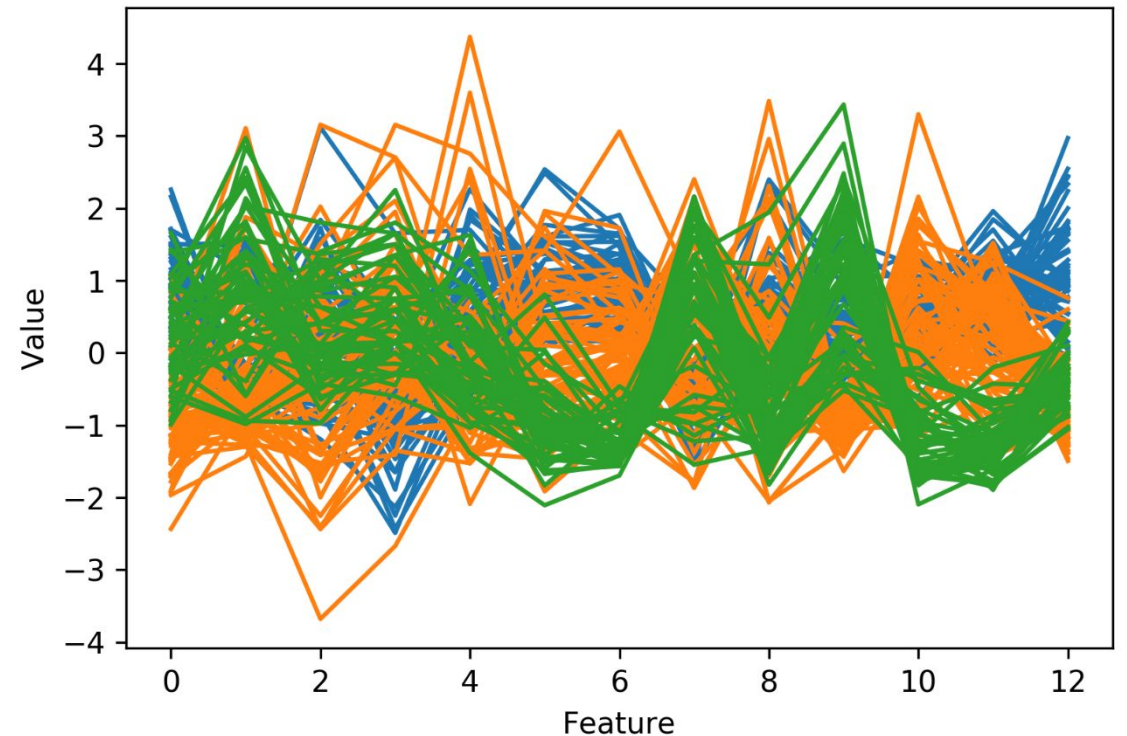


PCA: Dimensionality reduction

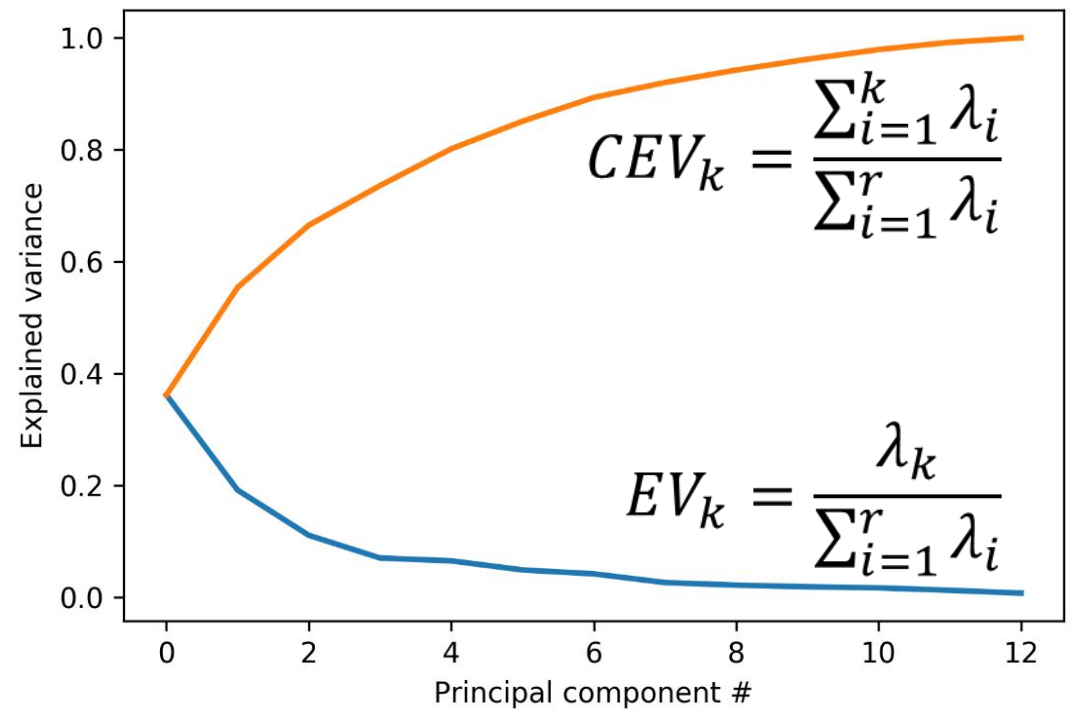
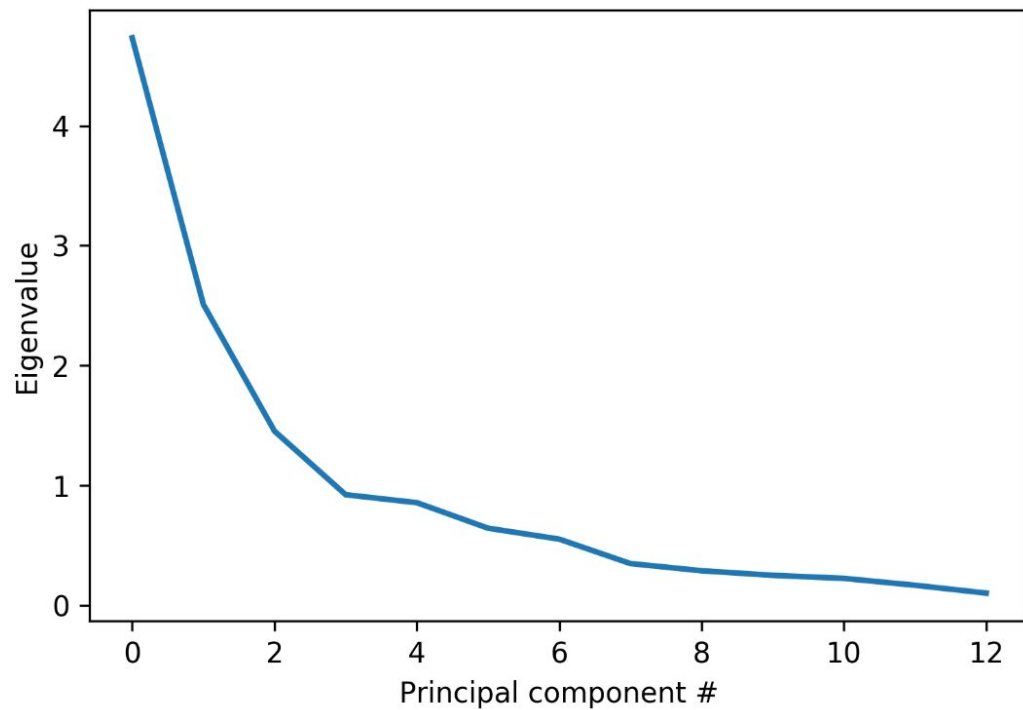
- We look for the subspace of dimension k onto which we can project our d -dimensional data
 - And we want to keep as much information as possible
- PCA “preserves” covariance: $\text{Var}(\boldsymbol{\eta}_k^T \mathbb{X}) = \lambda_k$

PCA: Dimensionality reduction

- UCI ML Wine Data Set:
 - 13 features (x-axis)
 - 178 samples (lines)
 - Values scaled to $\text{Sam}(0, 1)$ (y-axis)
 - 3 classes (color)
- Principal components often correspond to underlying structure of the data
- As with most linear methods, data must be scaled for practical applications



PCA: Dimensionality reduction



PCA: Dimensionality reduction

