

Principal Component Analysis

Data Mining

Prof. Sujee Lee

Department of Systems Management Engineering

Sungkyunkwan University

Dimensionality Reduction

■ Dimensionality Reduction

- Find a new way to represent this data that summarizes the essential characteristics with **fewer features**.
- **Motivation:** Algorithms like k-means are more computationally intensive in higher dimension and/or might take longer to converge
- **Benefits:**
 - Computational Savings: compress data -> saving in **time/space efficiency**
 - Statistical Benefits: fewer dimension -> **better generalization** from fewer observations
 - **Visualization:** look at data in 1, 2, or 3 dimension (for identifying outliers, etc.)
- **Methods :**
 - (Projection) Principal Component Analysis (PCA)
 - (Manifold Learning) t-distributed Stochastic Neighbor Embedding (t-SNE)

Principal Component Analysis

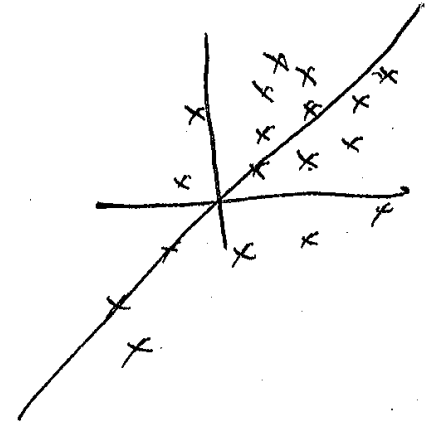
■ Principal Component Analysis

- Given a (training) dataset $D = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ such that $\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,p})$ is the i -th input vector of d features
- **Goal** : Find “best” linear transformation $f: \mathbb{R}^p \rightarrow \mathbb{R}^q$ where $q < p$, which maps data into a q -dimensional space.
 - “best” means that the transformations maximally captures variation in the data.

Principal Component Analysis: simple case

- Consider a case $p = 2$, $q = 1$ (1-dimensional projection)
 - Just like LDA, we want to find a direction $u \in \mathbb{R}^d$ (unit vector) that maximizes variance of data, so we consider the projected points $\{u^T x_i\}_{i=1}^n$, and we want **maximize sample variance**:

$$\frac{1}{n} \sum_i^n (u^T x_i - u^T \bar{x})^2, \text{ where } \bar{x} = \frac{1}{n} \sum_i^n x_i$$



- We can rewrite the **objective function** as

$$\begin{aligned} \frac{1}{n} \sum_i^n (u^T (x_i - \bar{x}))^2 &= \frac{1}{n} \sum_i^n (u^T (x_i - \bar{x})(x_i - \bar{x})^T u) \\ &= u^T \left(\frac{1}{n} \sum_i^n (x_i - \bar{x})(x_i - \bar{x})^T \right) u = u^T S u \end{aligned}$$

for S defined from data

Principal Component Analysis: simple case

- So PCA amounts to solving:

$$\max_{\|u\|_2=1} u^T S u, \quad \text{where } S = \frac{1}{n} \sum_i^n (x_i - \bar{x})(x_i - \bar{x})^T$$

- The solution is the top eigenvector of S , i.e., the eigenvector corresponding to the max eigenvalues

- **PCA algorithm for $q = 1$:**

1. Construct the matrix $S = \frac{1}{n} \sum_i^n (x_i - \bar{x})(x_i - \bar{x})^T$
2. Compute the top eigenvector u
3. Project all the data to $\{u^T x_i\}_{i=1}^n$

Principal Component Analysis: simple case

- It is convenient to write \mathbf{S} in terms of data matrix \mathbf{X} :

- Recall: If $\mathbf{X} = \begin{pmatrix} -\mathbf{x}_1^T & - \\ \cdots & \\ -\mathbf{x}_n^T & - \end{pmatrix}$, then $\frac{\mathbf{X}^T \mathbf{X}}{n} = \frac{1}{n} \sum_i^n \mathbf{x}_i \mathbf{x}_i^T$

- To construct \mathbf{S} , we can just construct a matrix $\tilde{\mathbf{X}}$ by subtracting column means from row of \mathbf{X} :

$$\tilde{\mathbf{X}} = \begin{pmatrix} -(\mathbf{x}_1 - \bar{\mathbf{x}})^T & - \\ \cdots & \\ -(\mathbf{x}_n - \bar{\mathbf{x}})^T & - \end{pmatrix}, \text{ then } \mathbf{S} = \frac{\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}}{n}$$

Singular Value Decomposition

■ Singular Value Decomposition (SVD)

- Why is $\max_{\|u\|_2=1} u^T S u$ solved with max eigenvector of S ?
- Fact1: Any real, symmetric matrix has an orthonormal basis of eigenvectors
- Fact2: If $A \in \mathbb{R}^{m \times n}$ (real, not necessarily symmetric), we can always construct a singular value decomposition (SVD) of A :

$$A = U \begin{pmatrix} D & 0 \\ 0 & 0 \end{pmatrix} V^T$$

where $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$, the columns of U are orthonormal vectors, columns of V are orthonormal, and D is diagonal with positive entries.

- Columns of U are left singular vectors of A
- Columns of V are right singular vectors of A
- Values in D are singular values of A
- In particular, if A is symmetric, then $U = V$ and columns correspond to eigenvectors of A .

Principal Component Analysis

- For any vector $\mathbf{u} \in \mathbb{R}^p$, we can write it as a linear combination of orthonormal eigenvectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p\}$ so that

$$\mathbf{u} = c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_p \mathbf{u}_p$$

where $c_1^2 + c_2^2 + \dots + c_p^2 = 1$ if \mathbf{u} is a unit vector

- Now, we can write

$$\begin{aligned} \mathbf{u}^T \mathbf{S} \mathbf{u} &= (c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_p \mathbf{u}_p)^T \mathbf{S} (c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_p \mathbf{u}_p) \\ &= (c_1 \mathbf{u}_1 + c_2 \mathbf{u}_2 + \dots + c_p \mathbf{u}_p)^T (c_1 \lambda_1 \mathbf{u}_1 + c_2 \lambda_2 \mathbf{u}_2 + \dots + c_p \lambda_p \mathbf{u}_p) \\ &= c_1^2 \lambda_1 + c_2^2 \lambda_2 + \dots + c_p^2 \lambda_p \end{aligned}$$

this is a weighted sum of $\lambda_1, \dots, \lambda_p$, and maximized if $c_1 = 1$, and $c_2 = c_3 = \dots = c_p = 0$.

i.e., $\mathbf{u} = \mathbf{u}_1$ is maximizer.

Principal Component Analysis

- Returning to PCA with $q > 1$, we want to project, the data $\{\mathbf{x}_i\}_{i=1}^n$ onto a subspace spanned by orthogonal unit vectors $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q\}$
- The coordinates with respect to the new basis are

$$\tilde{\mathbf{x}}_i \mapsto (\mathbf{u}_1^T \tilde{\mathbf{x}}_i, \mathbf{u}_2^T \tilde{\mathbf{x}}_i, \dots, \mathbf{u}_q^T \tilde{\mathbf{x}}_i) \in \mathbb{R}^q$$

We want to maximize the sum of squared lengths of projected vectors:

$$\max_{\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q\}} \frac{1}{n} \sum_i^n \|(\mathbf{u}_1^T \tilde{\mathbf{x}}_i, \mathbf{u}_2^T \tilde{\mathbf{x}}_i, \dots, \mathbf{u}_q^T \tilde{\mathbf{x}}_i)\|_2^2$$

Principal Component Analysis

- Using the matrix notation $\mathbf{U} = (\mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_q) \in \mathbb{R}^{p \times q}$,

PCA is :

$$\max_{\mathbf{U} \in \mathbb{R}^{p \times q} : \mathbf{U}^T \mathbf{U} = \mathbf{I}} \frac{1}{n} \sum_i^n \|\mathbf{U}^T \tilde{\mathbf{x}}_i\|_2^2$$

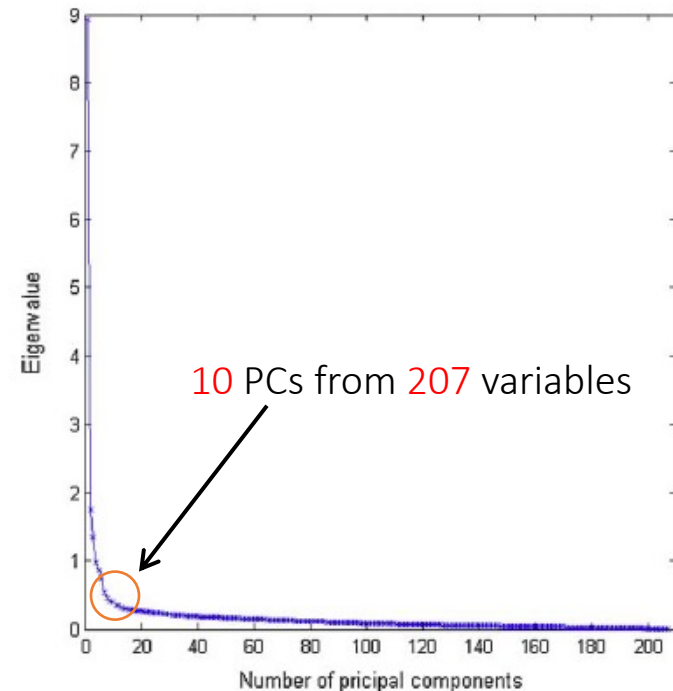
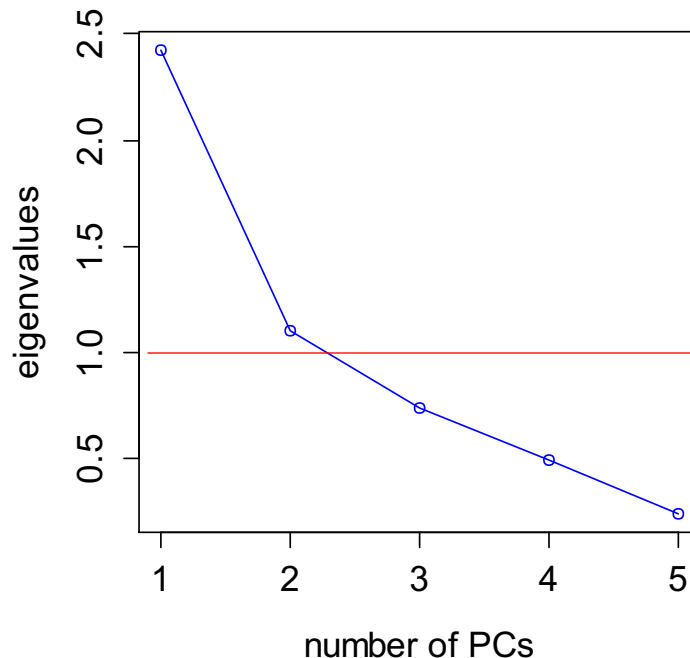
Maximizer is when \mathbf{U} consists of top q eigenvectors of $\mathbf{S} = \frac{\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}}{n}$

- PCA output is to map each \mathbf{x}_i to $\mathbf{U}^T \tilde{\mathbf{x}}_i$
- Note: another interpretation of the objective function is to minimize sum of squared distances to space spanned by $\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_q\}$

Number of Principal Components

- How many PCs? – PCA as Dimension Reduction Tech.

- 1) visualization: $q=1,2$, or 3
- 2) computational consideration: q is at most something
- 3) Elbow plot: plot eigenvalues of S , find elbow



Number of Principal Components

- 4) compute % of variation explained by principal components:

Consider eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$,

Compute ratio $\frac{\lambda_1 + \lambda_2 + \dots + \lambda_q}{\sum_i^p \lambda_i}$, which is a fraction of explained variance.

Determine a proper q considering it

Principal Component Analysis

■ PCA recap.

- PCA maximizes variation in p -dimensional data when projected onto a q -dimensional subspace $q < p$

1. Recenter the data: $\mathbf{x}_i \rightarrow \tilde{\mathbf{x}}_i$, where $\tilde{\mathbf{x}}_i = \mathbf{x}_i - \bar{\mathbf{x}}$

this gives a data matrix instead of $\tilde{\mathbf{X}}$ instead of \mathbf{X}

2. Construct the matrix $\mathbf{S} = \frac{1}{n} \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} = \frac{1}{n} \sum_i \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \in \mathbb{R}^{p \times p}$

3. Compute the q eigenvectors corresponding to the largest q eigenvalues of \mathbf{S} . Store them in matrix $\mathbf{U} = (\mathbf{u}_1 \mathbf{u}_2 \dots \mathbf{u}_q) \in \mathbb{R}^{p \times q}$

4. Map data to $\{\mathbf{U}^T \tilde{\mathbf{x}}_i\}_{i=1}^n$