

Dimensionality Reduction

Dr. Supaporn Erjongmanee

Department of Computer Engineering
Kasetsart University
fengspe@ku.ac.th

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 1



Department of Computer Engineering
Kasetsart University

1

Outline

- Introduction
- Principal Component Analysis

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 2



Department of Computer Engineering
Kasetsart University

2

Dimensionality Reduction

- Transformation process of data with many dimensions (many variables) to lower dimensions
- Two main approaches
 1. Feature selection
 - Selection subset of variables
 2. Feature extraction
 - Reduce higher dimensionality space to lower dimensionality subspace

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 3

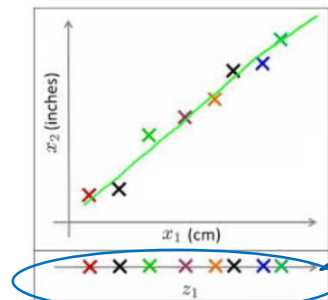


Department of Computer Engineering
Kasetsart University

3

Feature Extraction Techniques

- Example:



- Three commonly-used techniques
 1. Principal Component Analysis (PCA)
 2. Linear Discriminant Analysis (LDA)
 3. Generalized Discriminant Analysis (GDA)

Image source: <https://www.analyticsvidhya.com/blog/2015/07/dimension-reduction-methods/>

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 4



Department of Computer Engineering
Kasetsart University

4

Feature Extraction Techniques (cont.)

- Advantages
 - Lower data storage
 - Lower computation time
- Disadvantages
 - Information loss
 - Meaning interpretation of lower dimension

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 5



Department of Computer Engineering
Kasetsart University

5

Review: Covariance Matrix

- Given random variables X_1, X_2, \dots, X_p
- Each entry in covariance matrix = $Cov(X_i, X_j)$
 - $Cov(X_i, X_j) = E[(X_i - E[X_i])(X_j - E[X_j])]$
- Note that $Cov(X_i, X_i) = \sigma_{X_i}^2$

$$\Sigma = \begin{bmatrix} \sigma_{X_1}^2 & \cdots & Cov(X_1, X_p) \\ \vdots & \ddots & \vdots \\ Cov(X_p, X_1) & \cdots & \sigma_{X_p}^2 \end{bmatrix}$$

Square (p x p)
and Symmetric
matrix

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 6

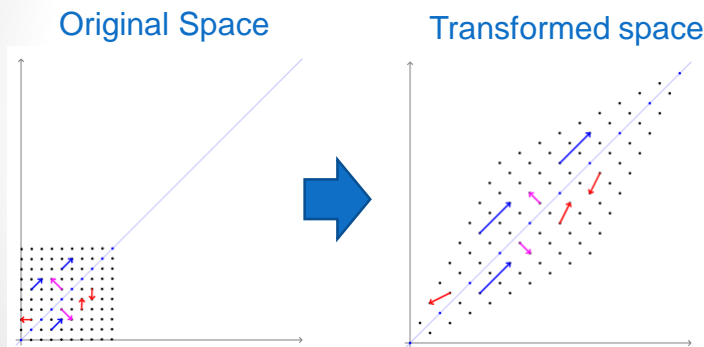


Department of Computer Engineering
Kasetsart University

6

Review: Eigenvectors and Eigenvalues

- Transformation of vector x in another space by changing its length (with the changing amount $= \lambda$), but not its direction



Notice that red and blue vectors maintain the same direction, but not their lengths

Let A = Transformation matrix

$$Ax = \lambda x$$

Image source: <https://pathmind.com/wiki/eigenvector>

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 7



Department of Computer Engineering
Kasetsart University

7

Review: Eigenvectors and Eigenvalues

To find eigenvalues, we use $Ax = \lambda x$
 $(A - \lambda I)x = 0$

A transforms vector x from one space to the other space

x is positive. For the above equation to have solution, $\det(A - \lambda I) = 0$

We solve for eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$

After deriving λ_i , use $Ae_i = \lambda_i e_i$ to solve for eigenvector e_i

Let $\Phi = [e_1 \ e_2 \ \dots \ e_p]$, and $\Lambda = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_p \end{bmatrix}$

We can rewrite $A = \Phi \Lambda \Phi^T$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 8



Department of Computer Engineering
Kasetsart University


8

Outline

- Introduction
- Principal Component Analysis



Principal Component Analysis

- Process of mapping data from higher dimension space to lower dimension subspace
- From p dimension  q dimension

$$p > q$$

Each dimension can be called component



Principal Component Analysis

- It is known that variable with large variance tends to explain output better than one with small variance
- Idea:
 1. Select the first component with largest variance.
 2. Then, select the next component that is orthogonal to the first one and also has large variance.
 3. Continue for all components
 4. Choose components that **most** explained original data (maybe with some information loss)

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 11

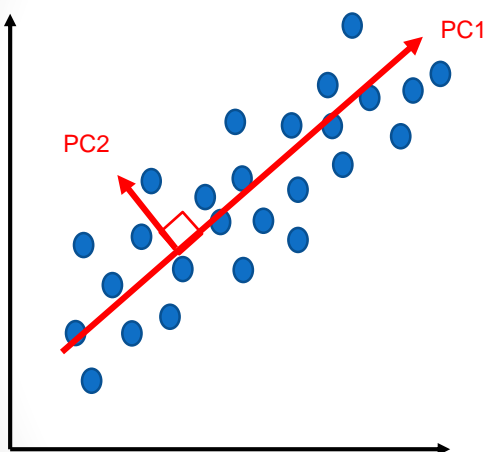


Department of Computer Engineering
Kasetsart University

11

PCA Concept

• Example



- Choose principle components 1 (PC1) with largest variance
- Choose principle components 2 (PC2) orthogonal to PC1 and with large variance
- Do we need both PC1 and PC2?
- Is PC1 sufficient to explain data?

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 12



Department of Computer Engineering
Kasetsart University

12

PCA Definition

- Let X be n samples, each sample has p variables

$$X = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \cdots & \cdots & \ddots & \cdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix}$$

- Let $\mu = [\mu_1, \mu_2, \dots, \mu_p] = \text{mean vector}$
- Let $\hat{X} = \text{centered data matrix}$

$$\hat{X} = \begin{bmatrix} x_{11} - \mu_1 & x_{12} - \mu_2 & \cdots & x_{1p} - \mu_p \\ x_{21} - \mu_1 & x_{22} - \mu_2 & \cdots & x_{2p} - \mu_p \\ \cdots & \cdots & \ddots & \cdots \\ x_{n1} - \mu_1 & x_{n2} - \mu_2 & \cdots & x_{np} - \mu_p \end{bmatrix}$$

Note that

$$\begin{aligned} E[\hat{X}] &= E[X - \mu] \\ &= E[X] - \mu = 0 \end{aligned}$$



PCA Definition (cont.)

- Let $Q = \hat{X}^T \hat{X}$

$$Q = \hat{X}^T \hat{X}$$

$$= \begin{bmatrix} x_{11} - \mu_1 & x_{12} - \mu_2 & \cdots & x_{1p} - \mu_p \\ x_{21} - \mu_1 & x_{22} - \mu_2 & \cdots & x_{2p} - \mu_p \\ \cdots & \cdots & \ddots & \cdots \\ x_{n1} - \mu_1 & x_{n2} - \mu_2 & \cdots & x_{np} - \mu_p \end{bmatrix}^T \begin{bmatrix} x_{11} - \mu_1 & x_{12} - \mu_2 & \cdots & x_{1p} - \mu_p \\ x_{21} - \mu_1 & x_{22} - \mu_2 & \cdots & x_{2p} - \mu_p \\ \cdots & \cdots & \ddots & \cdots \\ x_{n1} - \mu_1 & x_{n2} - \mu_2 & \cdots & x_{np} - \mu_p \end{bmatrix}$$

$Q = \Sigma = \text{Covariance Matrix}$

Square (p x p) and Symmetric



Projection Residuals

- Note that centered data have mean $E[\hat{X}] = 0$
- Let one centered data sample be \vec{x}_i
- Let \vec{w} be *unit vector* of first component
- Then, the projected vector of \vec{x}_i on first component will be $(\vec{w} \cdot \vec{x}_i) \vec{w}$
- Projection residual = difference between the centered data and the projected vector = $\|\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}\|^2$



PCA Setup

- From the p-dimensional space, find q-dimensional subspace
- Use data (p-dimensional vector) to project to different component Which component to choose?
- Goal: Each component will maximize variance
- Instead of finding component with maximum variance, we will show that it is equivalent to find component with smallest projection residual
 - To find component with maximum variance = To minimize projection residual



Minimize Projection Residuals (cont.)

$$\begin{aligned}
 \|\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}\|^2 &= ((\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}) \cdot (\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w})) \\
 &= \vec{x}_i \cdot \vec{x}_i - \vec{x}_i \cdot (\vec{w} \cdot \vec{x}_i) \vec{w} - (\vec{w} \cdot \vec{x}_i) \vec{w} \cdot \vec{x}_i + (\vec{w} \cdot \vec{x}_i) \vec{w} \cdot (\vec{w} \cdot \vec{x}_i) \vec{w} \\
 &= \|\vec{x}_i\|^2 - 2(\vec{w} \cdot \vec{x}_i)^2 + (\vec{w} \cdot \vec{x}_i)^2 \vec{w} \cdot \vec{w} \\
 &= \underbrace{\|\vec{x}_i\|^2 - (\vec{w} \cdot \vec{x}_i)^2}_{\text{Residual of one data vector}} \quad \text{since } \vec{w} \cdot \vec{w} = \|\vec{w}\|^2 = 1
 \end{aligned}$$

Residual of one data vector

Average residuals of all n data vectors = MSE

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 17



Department of Computer Engineering
Kasetsart University

17

Minimize Projection Residuals (cont.)

$$\begin{aligned}
 MSE(\vec{w}) &= \frac{1}{n} \sum_{i=1}^n (\|\vec{x}_i\|^2 - (\vec{w} \cdot \vec{x}_i)^2) \quad \left. \vphantom{\sum_{i=1}^n} \right\} \text{Average residuals of all n sample data vector} \\
 &= \frac{1}{n} \left(\underbrace{\sum_{i=1}^n \|\vec{x}_i\|^2}_{\text{Do not depend on } \vec{w}} - \underbrace{\sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2}_{\text{To minimize MSE, make this term very large}} \right)
 \end{aligned}$$

Do not depend on \vec{w} To minimize MSE, make this term very large

Hence, to minimize MSE, we want to maximize $\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 18



Department of Computer Engineering
Kasetsart University

18

Minimize Projection Residuals (cont.)

$$Var[Y] = E[Y^2] - E[Y]^2 = \frac{1}{n} \left(\sum_{i=1}^n y_i^2 \right) - \left(\frac{1}{n} \sum_{i=1}^n y_i \right)^2$$

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2}_{\text{To minimize MSE, we maximize this}} = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \vec{w} \cdot \vec{x}_i \right)^2}_{\text{Equal to zero (see below)}} + \underbrace{Var[\vec{w} \cdot \vec{x}_i]}_{\text{Hence, to minimize MSE, we maximize this}}$$

To minimize MSE,
we maximize this

Equal to zero
(see below)

Hence, to minimize MSE,
we maximize this

$$\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i) = \left(\frac{1}{n} \vec{w} \right) \left(\sum_{i=1}^n \hat{x}_i \right)$$

Note that \hat{X} = centered data, $E[\hat{X}] = 0$

$$\|\vec{x}_i - (\vec{w} \cdot \vec{x}_i) \vec{w}\|^2$$

To minimize projection residual
= To find projection with
maximum variance

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 19



Department of Computer Engineering
Kasetsart University

19

Minimize Projection Residuals (cont.)

$$Var[Y] = E[Y^2] - E[Y]^2 = \frac{1}{n} \left(\sum_{i=1}^n y_i^2 \right) - \left(\frac{1}{n} \sum_{i=1}^n y_i \right)^2$$

$$\underbrace{\frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2}_{\text{To minimize MSE, we maximize this}} = \underbrace{\left(\frac{1}{n} \sum_{i=1}^n \vec{w} \cdot \vec{x}_i \right)^2}_{\text{Equal to zero (see below)}} + \underbrace{Var[\vec{w} \cdot \vec{x}_i]}_{\text{Hence, to minimize MSE, we maximize this}}$$

To minimize MSE,
we maximize this

Equal to zero
(see below)

Hence, to minimize MSE,
we maximize this

This is variance of projection to first principal component

We need to project all \vec{x}_i' 's to other principal components.

Hence, to minimize MSE = to maximize sum of variances to all components

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 20



Department of Computer Engineering
Kasetsart University

20

Maximize Projection Variance

- Note that $\hat{\mathbf{x}}$ be n centered data vectors with p dimension
 - $\hat{\mathbf{x}} = n \times p$ matrix
- Given \vec{w} be *unit vector* of each projection, then projection of all $\hat{\mathbf{x}}$ onto all projections $= \hat{\mathbf{x}}\mathbf{w}$
 - $\hat{\mathbf{x}}\mathbf{w} = n \times 1$

$$\sigma_w^2 = \frac{1}{n} \sum_{i=1}^n (\vec{w} \cdot \vec{x}_i)^2 = \frac{1}{n} (\hat{\mathbf{x}}\mathbf{w})^T (\hat{\mathbf{x}}\mathbf{w}) = \frac{1}{n} \mathbf{w}^T \hat{\mathbf{x}}^T \hat{\mathbf{x}} \mathbf{w}$$

Variance of one projection Using n data

$$= \mathbf{w}^T \frac{\hat{\mathbf{x}}^T \hat{\mathbf{x}}}{n} \mathbf{w} = \mathbf{w}^T \Sigma \mathbf{w}$$

$\Sigma = \hat{\mathbf{x}}^T \hat{\mathbf{x}}$ ($p \times p$ matrix)

Find \mathbf{w} that maximizes σ_w^2

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 21



Department of Computer Engineering
Kasetsart University

21

Maximize Projection Variance (cont.)

- Maximize σ_w^2 with constraint $\mathbf{w}^T \mathbf{w} = 1$
- Solve by using Lagrange Multiplier

$$\mathcal{L}(\mathbf{w}, \lambda) = \sigma_w^2 - \lambda(\mathbf{w}^T \mathbf{w} - 1) = \mathbf{w}^T \Sigma \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{w} - 1)$$

$$\frac{\partial \mathcal{L}}{\partial \lambda} = \mathbf{w}^T \mathbf{w} - 1 \xrightarrow{\text{Set to zero}} \mathbf{w}^T \mathbf{w} = 1$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 2\Sigma\mathbf{w} - 2\lambda\mathbf{w} \xrightarrow{\text{Set to zero}} \Sigma\mathbf{w} = \lambda\mathbf{w}$$

\mathbf{w} that maximizes σ_w^2
= eigenvector of Σ with
eigenvalue = λ

Eigenvector of covariance matrix

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 22



Department of Computer Engineering
Kasetsart University

22

Principal Components

$$\Sigma \mathbf{w} = \lambda \mathbf{w}$$

- Since Σ is $p \times p$ matrix, there will be p eigenvectors
- Σ is covariance matrix $\rightarrow \Sigma$ is symmetric
 - Eigenvectors are orthogonal to each other
- Σ is covariance matrix \rightarrow values in Σ are positive
 - Eigenvalues are positive
- **p eigenvectors of $\Sigma = p$ principal components of $\hat{\mathbf{x}}$**

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 23



Department of Computer Engineering
Kasetsart University

23

Principal Components (cont.)

- **p Eigenvectors of $\Sigma = p$ principal components of $\hat{\mathbf{x}}$**
 - Eigenvector with largest eigenvalue = \mathbf{e}_1 = first principal component
 - Eigenvector with second largest eigenvalue = \mathbf{e}_2 = second principal component
 - ...
 - Eigenvector with the p^{th} largest eigenvalue = \mathbf{e}_p = p^{th} principal component
- Note that each principal component is orthogonal to each other
- Each eigenvalue = variance described by each component
- Sum of eigenvalues = total variances described by all components

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 24



Department of Computer Engineering
Kasetsart University

24

Principal Components (cont.)

- There are p eigenvectors $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$ that correspond to eigenvalues $\lambda_1, \lambda_2, \dots, \lambda_p$ and $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$

$$\Sigma \mathbf{e}_1 = \lambda_1 \mathbf{e}_1, \Sigma \mathbf{e}_2 = \lambda_2 \mathbf{e}_2, \dots, \Sigma \mathbf{e}_p = \lambda_p \mathbf{e}_p$$

$$\underbrace{\Sigma [\mathbf{e}_1 \quad \mathbf{e}_2 \quad \dots \quad \mathbf{e}_p]}_{\Phi} = [\mathbf{e}_1 \quad \mathbf{e}_2 \quad \dots \quad \mathbf{e}_p] \underbrace{\begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_p \end{bmatrix}}_{\Lambda}$$

$$\Sigma \Phi = \Phi \Lambda$$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 25



Department of Computer Engineering
Kasetsart University

25

Principal Components (cont.)

- If we normalize eigenvector to unit vector

$$\Phi \Phi^T = \Phi^T \Phi = I$$

$$\Sigma \Phi = \Phi \Lambda$$

$$\Rightarrow \Phi^T \Sigma \Phi = \Lambda$$

$$\Sigma = \Phi \Lambda \Phi^T$$



Eigenvector matrix Eigenvector matrix

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 26



Department of Computer Engineering
Kasetsart University

26

PCA Procedure

1. Compute \hat{X} Dimension of $X = [n \times p]$
2. Compute Σ = covariance matrix of \hat{X}
3. Compute single value decomposition (SVD) of $\Sigma = \Phi \Lambda \Phi^T$
 - $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p$
 - $\Phi = [e_1, e_2, \dots, e_p]$
4. Choose $q < p$ and obtain $\Phi_q = [e_1, e_2, \dots, e_q]$
5. Obtain $y = \Phi_q^T x$ Dimension of $Y = [n \times q]$

Question: How to choose q ?

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 27



Department of Computer Engineering
Kasetsart University

27

How to choose q ?

- $\sum_{i=1}^p \lambda_i$ = total variance described by all components
- $\sum_{i=1}^q \lambda_i$ = variance explained by PCA
- Define fraction of original variance and input vectors
 - $R^2 = \frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i} \geq 1 - \alpha$
 - where α = accepted error and $0 < \alpha < 1$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 28



Department of Computer Engineering
Kasetsart University

28

Example

- Beside variances, **eigenvalues also specify percent of transformation along each component**
- If there are 3 eigenvectors with eigenvalues
 - $\{\lambda_1 = 3, \lambda_2 = 2, \lambda_3 = 1\}$
 - Percent of transformation occurs in 1st component = 50%
 - Percent of transformation occurs in 2nd component = 33.33%
 - Percent of transformation occurs in 3rd component = 16.67%
- If we choose only 2 principal components ($q=2$),
 - Percent of data explained by PCA = 83.33%
 - Some information will be lost.

Question: Is it acceptable?



Projection Residuals vs. q

- Let $\Phi_q = q$ eigenvectors
 - Projection residual = $x - \Phi_q x$
- If data are really q -dimensional, there are q eigenvectors with q positive eigenvalues.
 - The remaining $p-q$ eigenvalues ≈ 0
 - Projection residual = 0
- If data approximately have q -dimensions, projection residual is small
- If data is larger than q -dimensions, projection residual is large



Example: Data Projection

- Obtain $y_i = \Phi_q^T x_i$

- Let $p = 3, q = 3$

No dimensionality
reduction

$$x_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix}$$

$$\Phi = [e_1 \quad e_2 \quad e_3]$$

$$\Phi^T x_i = \begin{bmatrix} e_1 \\ e_2 \\ e_3 \end{bmatrix} \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \\ e_{31} & e_{32} & e_{33} \end{bmatrix} \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = \begin{bmatrix} y_{i,1} \\ y_{i,2} \\ y_{i,3} \end{bmatrix} = y_i$$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 31



Department of Computer Engineering
Kasetsart University

31

Example 2: Data Projection

- Obtain $y_i = \Phi_q^T x_i$

- Let $p = 3, q = 2$

Dimensionality
reduction

$$x_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix}$$

$$\Phi = [e_1 \quad e_2 \quad e_3]$$

$$\Phi_q = [e_1 \quad e_2]$$

$$\Phi_q^T x_i = \begin{bmatrix} e_1 \\ e_2 \end{bmatrix} \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = \begin{bmatrix} e_{11} & e_{12} & e_{13} \\ e_{21} & e_{22} & e_{23} \end{bmatrix} \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = \begin{bmatrix} y_{i,1} \\ y_{i,2} \end{bmatrix} = y_i$$

Supaporn Erjongmanee
fengspe@ku.ac.th

Statistics in Computer Engineering
Slide 32



Department of Computer Engineering
Kasetsart University

32

Example: Data Reconstruction

- Reconstruct $x_i = \Phi_q y_i$

- Let $p = 3, q = 3$

$$x_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix}$$

$$y_i = \begin{bmatrix} y_{i,1} \\ y_{i,2} \\ y_{i,3} \end{bmatrix}$$

$$\Phi = [e_1 \quad e_2 \quad e_3]$$

$$\Phi_q y_i = [e_1 \quad e_2 \quad e_3] \begin{bmatrix} y_{i,1} \\ y_{i,2} \\ y_{i,3} \end{bmatrix} = \begin{bmatrix} e_{11} & e_{21} & e_{31} \\ e_{12} & e_{22} & e_{32} \\ e_{13} & e_{23} & e_{33} \end{bmatrix} \begin{bmatrix} y_{i,1} \\ y_{i,2} \\ y_{i,3} \end{bmatrix} = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = x_i$$

No information loss



Example 2: Data Reconstruction

- Reconstruct $x_i = \Phi_q y_i$

- Let $p = 3, q = 2$

$$x_i = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix}$$

$$y_i = \begin{bmatrix} y_{i,1} \\ y_{i,2} \end{bmatrix}$$

$$\Phi = [e_1 \quad e_2 \quad e_3]$$

$$\Phi_q = [e_1 \quad e_2]$$

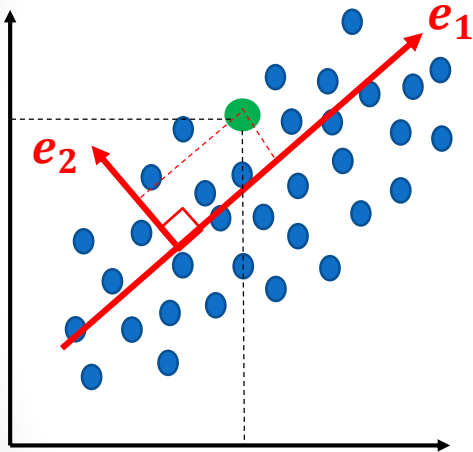
$$\Phi_q y_i = [e_1 \quad e_2] \begin{bmatrix} y_{i,1} \\ y_{i,2} \end{bmatrix} = \begin{bmatrix} e_{11} & e_{21} \\ e_{12} & e_{22} \\ e_{13} & e_{23} \end{bmatrix} \begin{bmatrix} y_{i,1} \\ y_{i,2} \end{bmatrix} = \begin{bmatrix} x_{i,1} \\ x_{i,2} \\ x_{i,3} \end{bmatrix} = x_i$$

Some information loss



PCA Transformation

• Visualization



$$\Sigma = \Phi \Lambda \Phi^T$$

Calculate eigenvectors of covariance matrix Σ
 $\Phi = [e_1, e_2]$

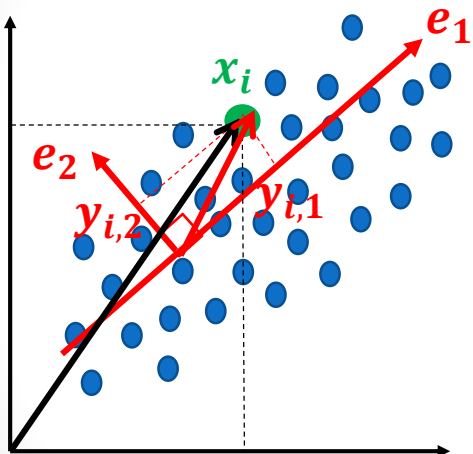


We derive PC1 = e_1 ,
 PC2 = e_2



Transformation

• Visualization



We derive
 PC1 = e_1
 PC2 = e_2



We can project each data point
 to principal component space

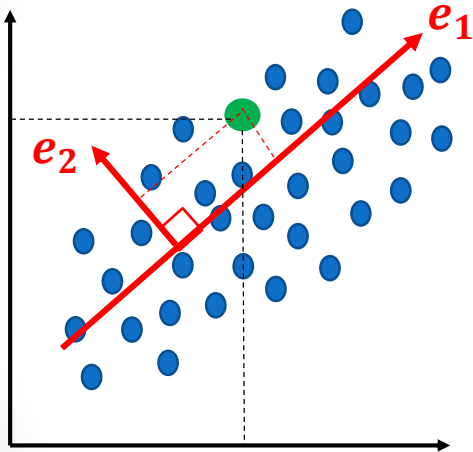
$$y_i = [e_1, e_2]^T \vec{x}_i$$

$$= \Phi^T \vec{x}_i$$



Transformation

- Visualization



With $PC1 = e_1$, $PC2 = e_2$



We can reconstruct each data point from principal components

$$\tilde{x}_i = [e_1, e_2]y_i$$

The reconstructed \tilde{x}_i may not have the same information as \vec{x}_i if $q < p$



Conclusion

- Dimensionality reduction is to reduce higher dimensional space to lower dimensional subspace
 - Some original information may be lost
- Principal component analysis is to find new components of data where each component span over maximum variances and orthogonal to other component
 - Some component with small eigenvalues (variances) can be ignored
 - This allows us to reduce data dimensionality



References

1. C. Shalizi, Principal Component Analysis, available at <https://www.stat.cmu.edu/~cshalizi/uADA/12/lectures/ch18.pdf>
2. L. Wiskott, Lecture Notes on Principal Component Analysis, available at <http://cs233.stanford.edu/ReferencedPapers/LectureNotes-PCA.pdf>
3. P. Rigollet, Principal Component Analysis, Statistics for Applications, available at: https://ocw.mit.edu/courses/mathematics/18-650-statistics-for-applications-fall-2016/lecture-slides/MIT18_650F16_PCA.pdf

