

Principal Component Analysis

COMP4211



THE DEPARTMENT OF
COMPUTER SCIENCE & ENGINEERING
計算機科學及工程學系

Pattern preprocessing may be necessary because

- some features are **irrelevant** to the classification task
- strong **correlations** exist between sets of features (i.e., the same information is repeated in several features)

feature extraction

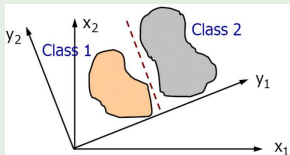
- the original feature space is **transformed** linearly or nonlinearly to a new space, usually of **lower** dimensionality

Unsupervised feature extraction

- feature extraction techniques rely **entirely** on the **input** data itself without reference to the corresponding target output data

Feature Extraction...

Example



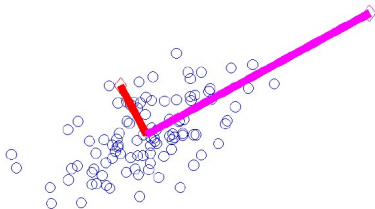
- Either x_1 or x_2 alone is not sufficient for classifying the two classes completely
- After transforming into the new y_1 – y_2 space, one feature (y_1) is sufficient for classifying the two classes

Goal

To preserve as much of the **relevant** information as possible during **dimensionality reduction**

Principal Component Analysis (PCA)

- Given: n d -dimensional points x_1, \dots, x_n



- In many practical applications, the data in \mathbb{R}^d is usually have a true/**intrinsic dimensionality** much lower than d
- PCA** is a powerful technique for extracting (lower-dimensional) structure (**feature extraction**) from possibly high-dimensional data sets
- aka **Karhunen-Loève (K-L)** transformation, **Hotelling** transformation

Zero-D Representation

How to find x_0 that represents x_1, \dots, x_n ?

Criterion: find x_0 such that the **sum of the squared distances** between x_0 and the various x_k is as small as possible

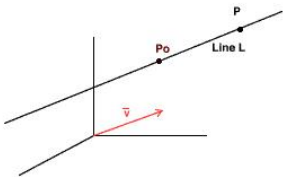
$$J_0(x_0) = \sum_{k=1}^n \|x_0 - x_k\|^2$$

- minimize $J_0(x_0) \rightarrow x_0 = m = \sum_{k=1}^n x_k / n$

The “best” zero-dimensional representation of the data set is the **sample mean**

One-D Representation

How to represent the set of points by a **line** through m ?



- w : **unit vector** along the line
- $x = m + aw$

$$J_1(a_1, \dots, a_n, w) = \sum_{k=1}^n \|(m + a_k w) - x_k\|^2$$

first consider the case where w is known

- recall that $\|w\| = 1$

$$\frac{\partial}{\partial a_k} J_1(a_1, \dots, a_n, w) = 0 \Rightarrow a_k = w^t (x_k - m)$$

Project x_k onto the line in the direction of w that passes through the sample mean

What is the Best Direction?

$$\begin{aligned}J_1(w) &= \sum_{k=1}^n a_k^2 \|w\|^2 - 2 \sum_{k=1}^n a_k w^t (x_k - m) + \sum_{k=1}^n \|x_k - m\|^2 \\&= \sum_{k=1}^n a_k^2 - 2 \sum_{k=1}^n a_k^2 + \sum_{k=1}^n \|x_k - m\|^2 \\&= - \sum_{k=1}^n (w^t (x_k - m))^2 + \sum_{k=1}^n \|x_k - m\|^2 \\&= - \sum_{k=1}^n w^t (x_k - m) (x_k - m)^t w + \sum_{k=1}^n \|x_k - m\|^2 \\&= - w^t \underbrace{\sum_{k=1}^n (x_k - m) (x_k - m)^t}_={S(\text{scatter matrix})} w + \sum_{k=1}^n \|x_k - m\|^2\end{aligned}$$

What is the Best Direction?...

$$\max w^t S w \quad \text{subject to } \|w\| = 1$$

- constrained optimization
- method of Lagrange multipliers

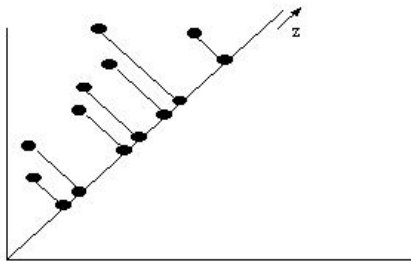
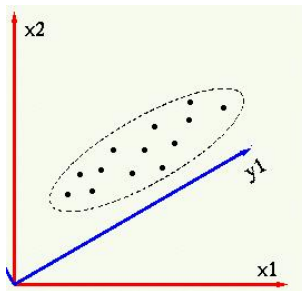
$$S w = \lambda w \quad (\text{i.e., } w \text{ is an eigenvector of } S)$$

Which eigenvector?

- find the one that is best in sum-of-squared-error
 - maximize $w^t S w = \lambda w^t w = \lambda$
 - select the eigenvector w corresponding to the largest eigenvalue of S

Dimensionality Reduction

Can be used to simplify a dataset by choosing a **new coordinate system**



- if we **only** keep y_1 but ignore y_2 , a 50% compression rate can be achieved without losing much information in the signal

Second Best Direction?

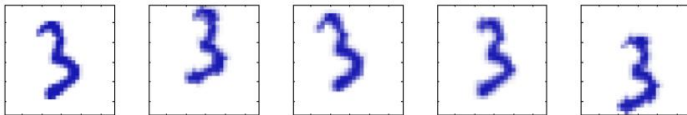
What is the **second best** direction?

- the **second best** direction should be **orthogonal** to the first best direction

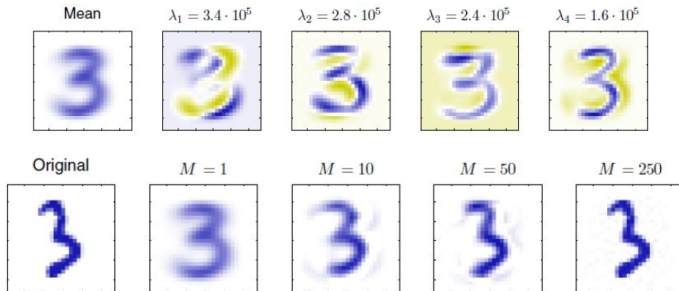
$$\max_w w^T S w \quad \text{s.t. } w^T w = 1 \text{ and } w^T w_1 = 0$$

- select eigenvector w_2 corresponding to **2nd largest** eigenvalue of S
- Similarly, the n th best direction is the eigenvector w_n corresponding to the n th largest eigenvalue of S

Example



- a collection of 100×100 images created from one image by introducing random displacement and rotation eigenvectors



Another Derivation

Find the projection w s.t. $\text{var}(w^t x)$ is maximized

$$\begin{aligned}\text{var}(w^t x) &= E[(w^t x - w^t m)^2] \\ &= E[(w^t x - w^t m)(w^t x - w^t m)] \\ &= E[w^t (x - m)(x - m)^t w] \\ &= w^t E[(x - m)(x - m)^t] w \\ &= w^t \Sigma w\end{aligned}$$

- $E[(x - m)(x - m)^t] = \Sigma$

maximize $\text{var}(w^t x)$ subject to $\|w\| = 1$

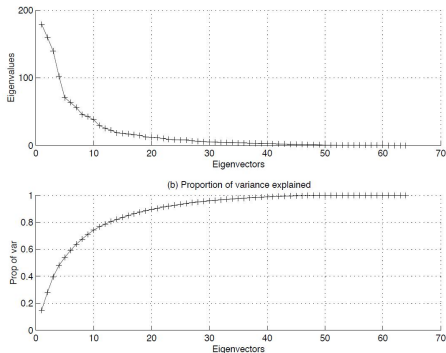
- choose the eigenvector with the **largest eigenvalue** for the variance to be maximum

How to Choose k ?

Proportion of variance explained:

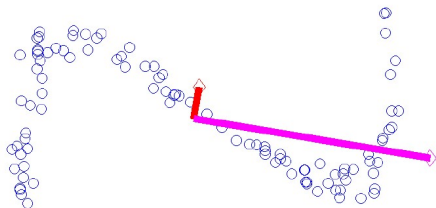
$$\frac{\lambda_1 + \lambda_2 + \cdots + \lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_d}$$

- λ_i are sorted in descending order



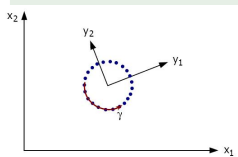
- e.g., stop at proportion of variance > 0.9

Limitations of PCA (1)



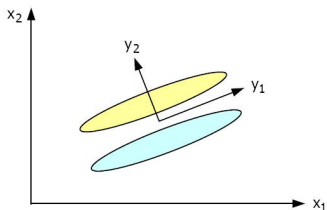
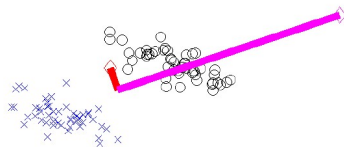
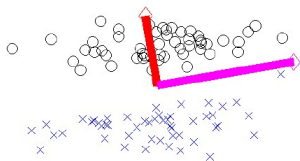
PCA is **linear**

Example



- dimensionality of the feature space is 2
- **intrinsic** dimensionality of the data distribution is 1
- Each point x in the data set can be specified (parametrically) by a single parameter (γ), instead of two variables x_1 and x_2

Limitations of PCA (2)



- data variance is largest along the y_1 direction
 - PCA transforms to one dimension will remove all the ability to **discriminate** the two classes
-
- For PCA to be effective in extracting useful features for classification, large variance in the data should correspond to large variance **between** classes rather than large variance **within** each class