

# Principal Component Analysis

Hee-il Hahn

Professor

Department of Information and Communications Engineering

Hankuk University of Foreign Studies

hihahn@hufs.ac.kr

- The most popular algorithm for dimensionality reduction
  - Finds the direction vectors along which the data has maximum variance.
  - Most useful in the case when the data lies on or close to a linear subspace of the data set.
  - Suppose that  $\mathbf{X} \in \mathbf{R}^{n \times p}$  is a matrix whose rows are  $p$  –dimensional data points.
  - We are looking for the  $d$  –dimensional linear subspace of  $\mathbf{R}^p$  along which the data has maximum variance.
  - The objective function it optimizes is

$$\max_V \text{var}(\mathbf{XV}) \quad \text{where } \mathbf{V} \text{ is an orthogonal } p \times d \text{ matrix.}$$

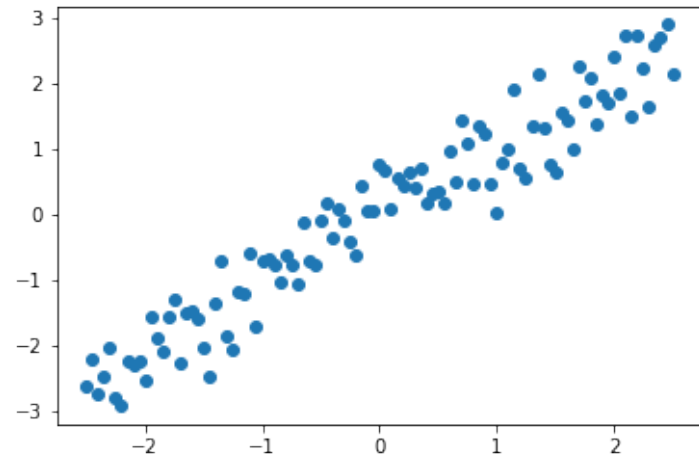
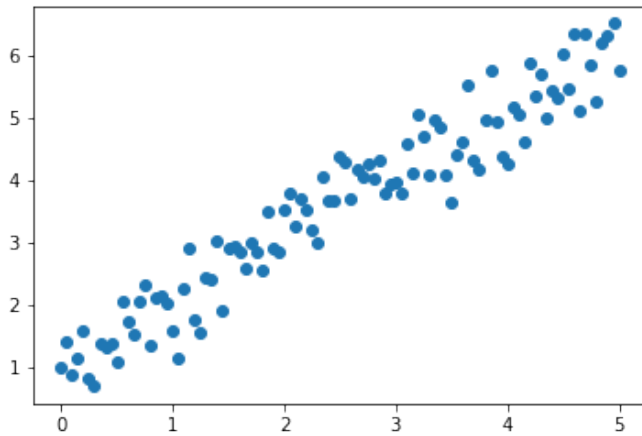
- If  $d = 1$ , then  $\mathbf{V}$  is simply a unit-length vector which gives the direction of maximum variance.

# Deriving the vector of maximum variance

$$\begin{aligned}\max_{\|\mathbf{v}\|=1} \text{var}(\mathbf{X}\mathbf{v}) &= \max_{\|\mathbf{v}\|=1} E(\mathbf{X}\mathbf{v})^2 - (E(\mathbf{X}\mathbf{v}))^2 \\ &= \max_{\|\mathbf{v}\|=1} E(\mathbf{X}\mathbf{v})^2 \quad \text{if mean-centered} \\ &= \max_{\|\mathbf{v}\|=1} \sum_{i=1}^n (\mathbf{x}_i \mathbf{v})^2 \\ &= \max_{\|\mathbf{v}\|=1} \sum_{i=1}^n (\mathbf{x}_i \mathbf{v})^T \mathbf{x}_i \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T \left( \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i \right) \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} \\ &= \max_{\|\mathbf{v}\|=1} \mathbf{v}^T \boldsymbol{\lambda} \mathbf{v} = \lambda_{\max}\end{aligned}$$

# PCA

## ■ Mean-centered data



# Principal component analysis (another ver.)

- Suppose a collection of  $m$  points  $\{x^{(1)} \dots x^{(m)}\}$  in  $\mathbf{R}^n$ 
  - lossy compression assumed.
  - for each point  $x^{(k)}$ , corresponding code vector  $c^{(k)} \in \mathbf{R}^l, l \leq n$
  - Find some encoding function  $f(x) = c$  and decoding function  $x \approx g(f(x))$
  - Let  $g(c) = Dc$  where  $D \in \mathbf{R}^{n \times l}$   
**constraints:** columns of  $D$  are orthonormal to each other.
  - For optimal code  $\Rightarrow$  minimize the distance between  $x$  and  $g(c)$
  - i.e.

$$\begin{aligned} c^* &= \operatorname{argmin}_c \|x - g(c)\|_2 \\ &= \operatorname{argmin}_c (x - g(c))^T (x - g(c)) \\ &= \operatorname{argmin}_c (x^T x - x^T g(c) - g(c)^T x + g(c)^T g(c)) \\ &= \operatorname{argmin}_c (-2x^T g(c) + g(c)^T g(c)) \end{aligned}$$

# Principal component analysis

$$\begin{aligned}\mathbf{c}^* &= \underset{\mathbf{c}}{\operatorname{argmin}} \|\mathbf{x} - g(\mathbf{c})\|_2 \\&= \underset{\mathbf{c}}{\operatorname{argmin}} (\mathbf{x} - g(\mathbf{c}))^T (\mathbf{x} - g(\mathbf{c})) \\&= \underset{\mathbf{c}}{\operatorname{argmin}} (\mathbf{x}^T \mathbf{x} - \mathbf{x}^T g(\mathbf{c}) - g(\mathbf{c})^T \mathbf{x} + g(\mathbf{c})^T g(\mathbf{c})) \\&= \underset{\mathbf{c}}{\operatorname{argmin}} (-2\mathbf{x}^T g(\mathbf{c}) + g(\mathbf{c})^T g(\mathbf{c})) \\&= \underset{\mathbf{c}}{\operatorname{argmin}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{D}^T \mathbf{D}\mathbf{c}) \\&= \underset{\mathbf{c}}{\operatorname{argmin}} (-2\mathbf{x}^T \mathbf{D}\mathbf{c} + \mathbf{c}^T \mathbf{c}) \quad (\mathbf{D}^T \mathbf{D} = \mathbf{I}) \\ \therefore \mathbf{c} &= \mathbf{D}^T \mathbf{x} \quad \Leftrightarrow \quad f(\mathbf{x}) = \mathbf{D}^T \mathbf{x} \text{ and} \\ &\quad r(\mathbf{x}) = g(f(\mathbf{x})) = \mathbf{D}\mathbf{D}^T \mathbf{x}\end{aligned}$$

# Principal component analysis

- How to choose the encoding matrix  $\mathbf{D} \in \mathbb{R}^{n \times l}$

$$\mathbf{D}^* = \underset{\mathbf{D}}{\operatorname{argmin}} \sqrt{\sum_{i,j} \left( x_j^{(i)} - r(\mathbf{x}^{(i)})_j \right)^2} \quad \text{subject to } \mathbf{D}^T \mathbf{D} = \mathbf{I}$$

Consider the case of  $l = 1$ .

$$\begin{aligned} \mathbf{d}^* &= \underset{\mathbf{d}}{\operatorname{argmin}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{d} \mathbf{d}^T \mathbf{x}^{(i)} \right\|_2^2 \quad \text{subject to } \|\mathbf{d}\|_2 = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{d}^T \mathbf{x}^{(i)} \mathbf{d} \right\|_2^2 \quad \text{subject to } \|\mathbf{d}\|_2 = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} \sum_i \left\| \mathbf{x}^{(i)} - \mathbf{x}^{(i)T} \mathbf{d} \mathbf{d} \right\|_2^2 \quad \text{subject to } \|\mathbf{d}\|_2 = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^T\|_F^2 \quad \text{subject to } \mathbf{d}^T \mathbf{d} = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} \operatorname{Tr}((\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^T)^T (\mathbf{X} - \mathbf{X} \mathbf{d} \mathbf{d}^T)) \quad \text{subject to } \|\mathbf{d}\|_2 = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} (\operatorname{Tr}(\mathbf{X}^T \mathbf{X} - \mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T - \mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X} + \mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T)) \end{aligned}$$

# Principal component analysis

- How to choose the encoding matrix  $\mathbf{D} \in \mathbf{R}^{n \times l}$

$$\begin{aligned}\mathbf{d}^* &= \underset{\mathbf{d}}{\operatorname{argmin}} \sum_i \|\mathbf{x}^{(i)} - \mathbf{d}\mathbf{d}^T \mathbf{x}^{(i)}\|_2^2 \quad \text{subject to } \|\mathbf{d}\|_2 = 1 \quad (= \mathbf{d}^T \mathbf{d} = 1) \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} (-2\operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \operatorname{Tr}(\mathbf{d} \mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T)) \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} (-2\operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T \mathbf{d} \mathbf{d}^T)) \quad \text{subject to } \mathbf{d}^T \mathbf{d} = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmin}} (-2\operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T) + \operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T)) \quad \text{subject to } \mathbf{d}^T \mathbf{d} = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmax}} (\operatorname{Tr}(\mathbf{X}^T \mathbf{X} \mathbf{d} \mathbf{d}^T)) \quad \text{subject to } \mathbf{d}^T \mathbf{d} = 1 \\ &= \underset{\mathbf{d}}{\operatorname{argmax}} (\operatorname{Tr}(\mathbf{d}^T \mathbf{X}^T \mathbf{X} \mathbf{d})) \quad \text{subject to } \mathbf{d}^T \mathbf{d} = 1\end{aligned}$$

$\Leftrightarrow$  optimal  $\mathbf{d} = \mathbf{d}^*$  = eigenvector of  $\mathbf{X}^T \mathbf{X}$  corresponding to the largest eigenvalue