

Machine Learning

5강

# 데이터 표현: 특징추출

컴퓨터과학과 이관용 교수

# 학습목차

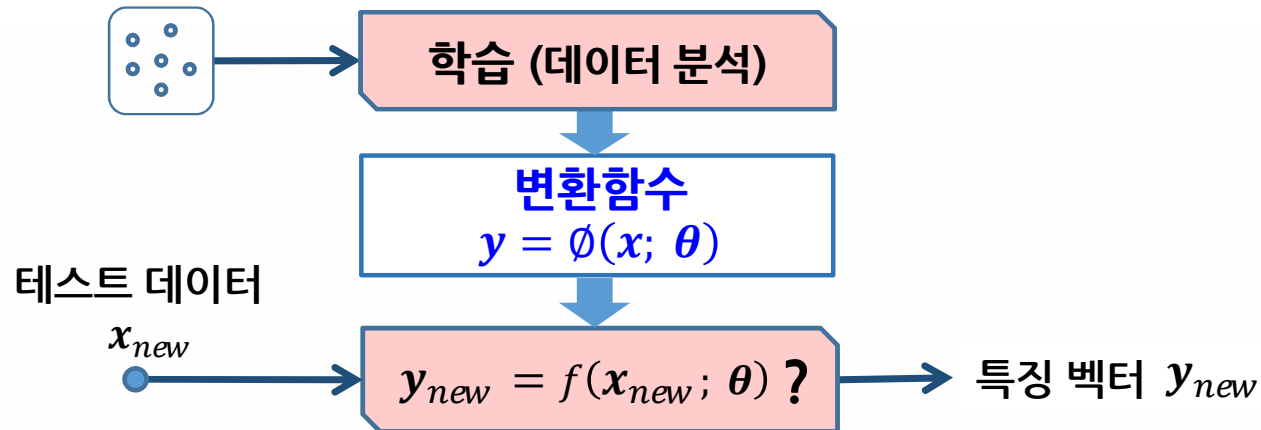
- 01 선형변환에 의한 특징추출
- 02 주성분분석법
- 03 선형판별분석법
- 04 거리 기반 차원 축소 방법

1

## 선형변환에 의한 특징추출

# 특징추출?

학습 데이터 집합 ( $D = \{x_i\}_{i=1, \dots, N}$ ,  $D = \{(x_i, y_i)\}_{i=1, \dots, N}$ )



## ○ 학습의 목적

- ☐ 분석에 불필요한 정보를 제거하고 핵심 정보만 추출
- ☐ 차원 축소를 통한 분석 시스템의 효율 향상

## 변환함수

### ○ 변환함수 embedding function, transformation function의 종류

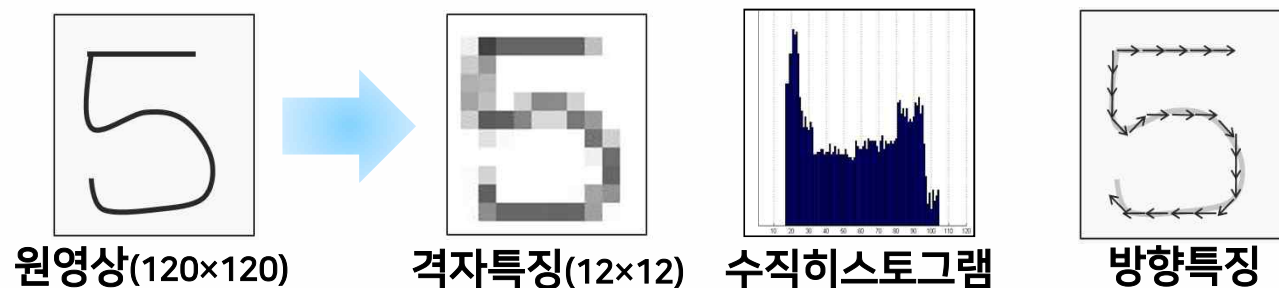
- 선형변환 linear transformation  $y = \phi(x) = W^T x$ 
  - ✓  $n$ 차원 열벡터  $x$ 에 변환행렬  $W$  ( $n \times m$ )을 곱하여  $m$ 차원 특징을 획득
  - ✓ 통계적 방법으로 특징벡터  $y$ 가 원하는 분포가 되도록 하는  $W$ 를 찾음
- 비선형변환
  - ✓ 복잡한 비선형함수  $\phi(x)$ 를 이용하여  $n$ 차원 벡터를  $m$ 차원 벡터로 매핑
  - ✓ 수작업 handcrafted에 의한 특징추출
  - ✓ 표현학습 representation learning

## 특징추출을 위한 접근법

### ○ 수작업에 의한 특징추출

□ 입력 데이터의 특성과 분석 목적에 맞는 특징을 개발자가 설계함

✓ 숫자인식을 위한 특징



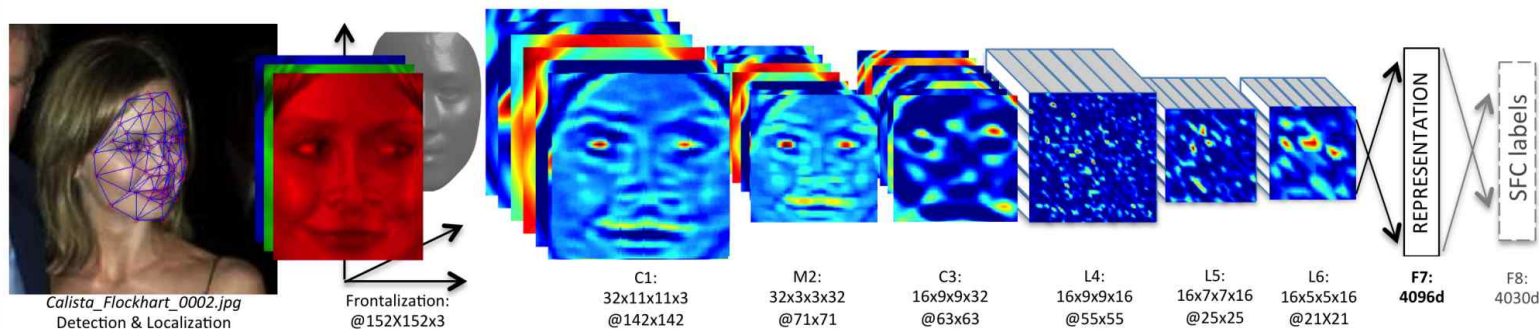
✓ 영상 분석을 위한 특징 → 에지, 가로/세로 방향 성분 등

✓ 문서 분석을 위한 특징 → 단어의 발생 빈도 등

## 특징추출을 위한 접근법

### 표현학습

- 특징추출을 위한 비선형 변환함수를 신경망 등의 머신러닝 모델로 표현
- 학습을 통해 분석이 잘 되도록 변환함수의 최적화가 가능
- 예: 딥러닝 모델을 이용한 얼굴인식용 특징추출



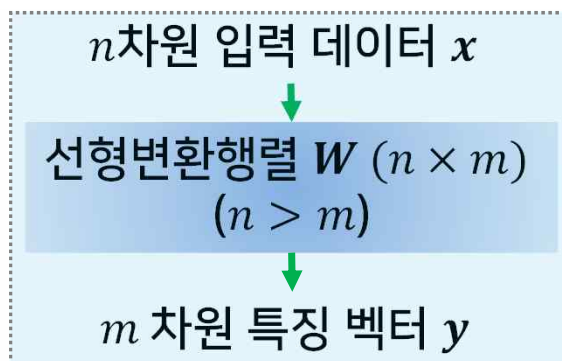
Y. Taigman 등, "DeepFace: Closing the gap to Human-Level Performance in Face Verification", CVPR2014

# 선형변환에 의한 특징추출

## ○ 차원 축소 관점에서의 특징추출



차원축소?



$$y = W^T x \rightarrow y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} w_1^T x \\ w_2^T x \\ \vdots \\ w_m^T x \end{bmatrix} = \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_m^T \end{bmatrix} x = [w_1, w_2, \dots, w_m]^T x = W^T x$$

□ 특징값  $y_i$

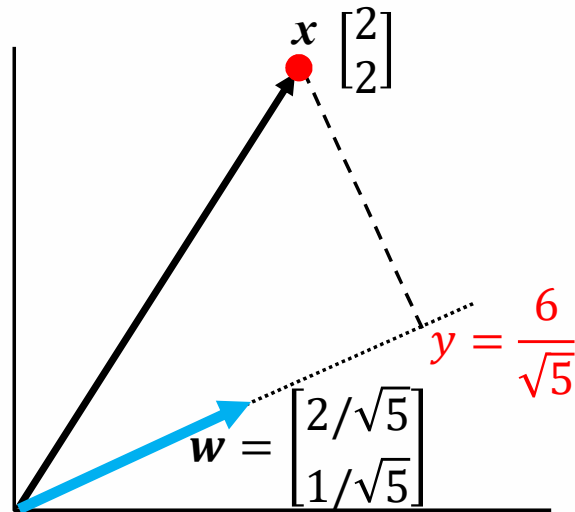
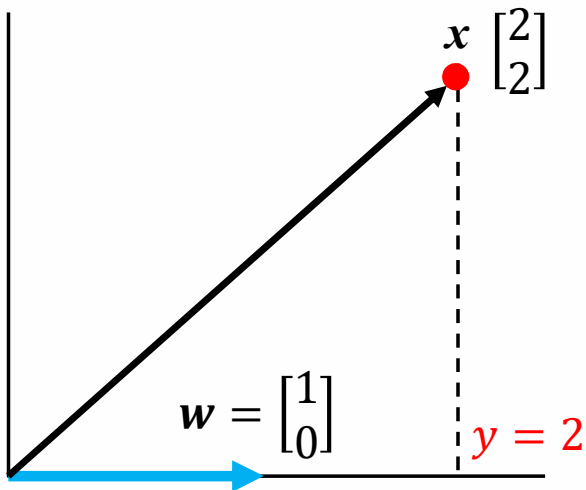
✓  $x$ 를  $W$ 의 열벡터  $w_i$  위로 사영한 값 (단,  $w_i$ 는 단위벡터)



## 선형변환에 의한 특징추출

○ 2차원 데이터  $x$ 를 1차원 특징  $y$ 로 변환

□ 벡터  $x$ 를  $w$  위로 사영  $\rightarrow y = w^T x$  (단,  $w$ 는 단위벡터)

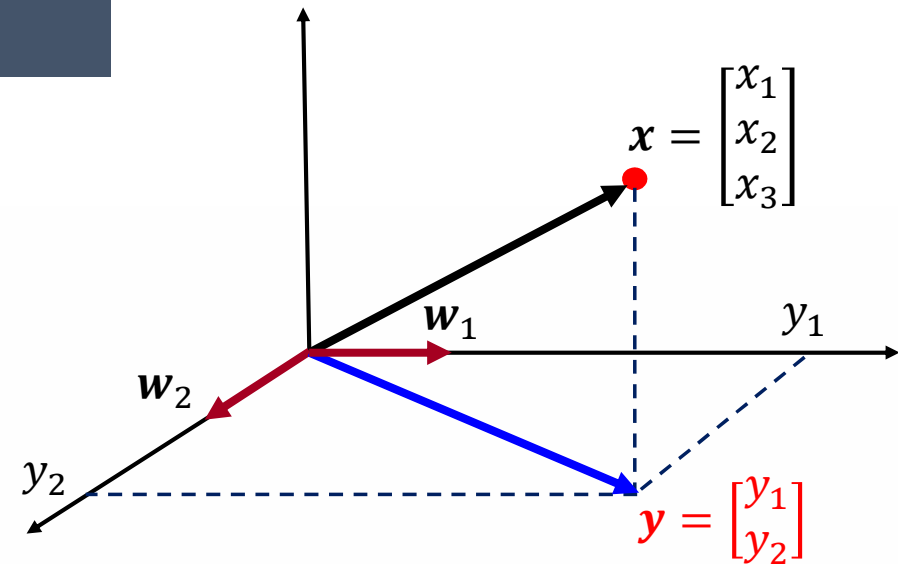


## 선형변환에 의한 특징추출

### ○ 3차원 데이터 → 2차원 특징추출

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{W} = [\mathbf{w}_1, \mathbf{w}_2] = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$

$$\mathbf{y} = \mathbf{W}^T \mathbf{x} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \end{bmatrix} = \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}$$



$y$ 는  $x$ 를 열벡터  $w_1, w_2$ 가 이루는 2차원 평면 위로의 사영으로 얻어지는 2차원 벡터

### ○ $n$ 차원 데이터 → $m$ 차원 특징추출

$$\mathbf{y} = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m]^T \mathbf{x} = \begin{bmatrix} \mathbf{w}_1^T \\ \mathbf{w}_2^T \\ \vdots \\ \mathbf{w}_m^T \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T \mathbf{x} \\ \mathbf{w}_2^T \mathbf{x} \\ \vdots \\ \mathbf{w}_m^T \mathbf{x} \end{bmatrix} = \mathbf{W}^T \mathbf{x}$$

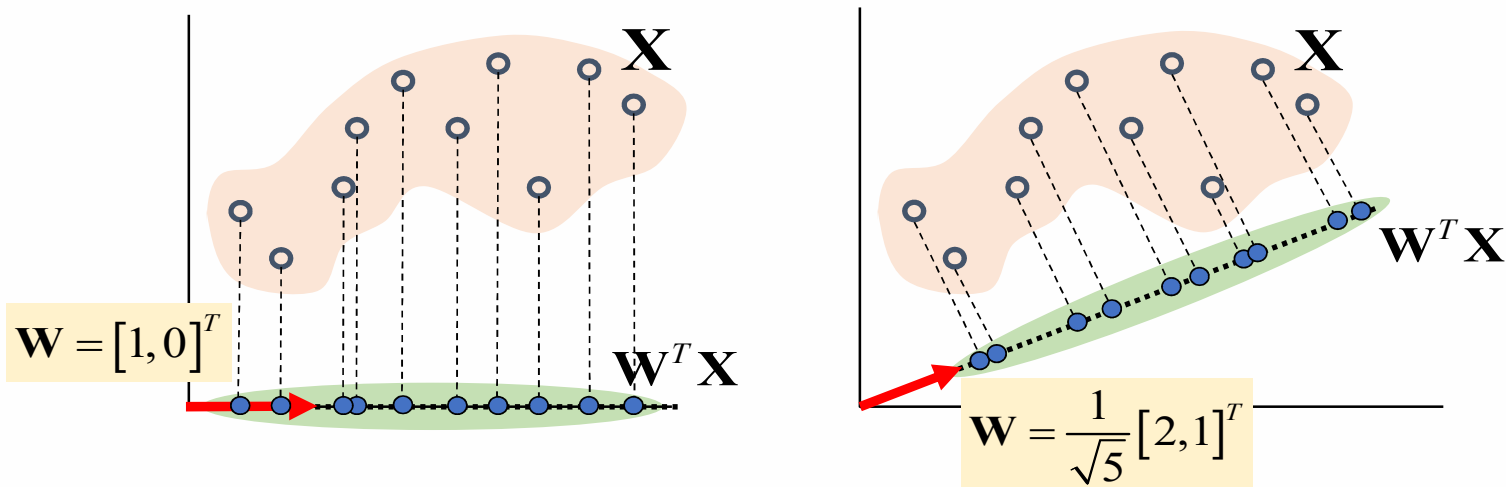
# 선형변환에 의한 특징추출

## ○ 전체 데이터 집합 $X$ 에 대한 특징추출

□  $X = [x_1, x_2, \dots, x_N]$  ( $n \times N$  행렬)

□  $Y = [y_1, y_2, \dots, y_N]$  ( $m \times N$  행렬)

→  $Y = W^T X = [W^T x_1, W^T x_2, \dots, W^T x_N]$

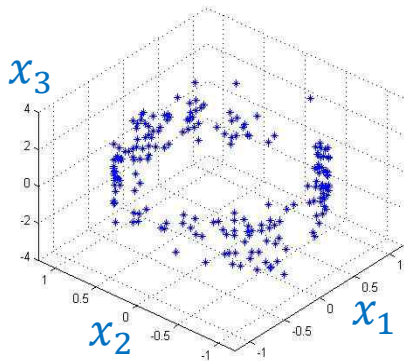


선형변환에 의한 특징추출? 주어진 데이터를 변환행렬  $W$ 에 의해  
정해지는 방향으로 사영함으로 저차원의 특징값을 얻는 것

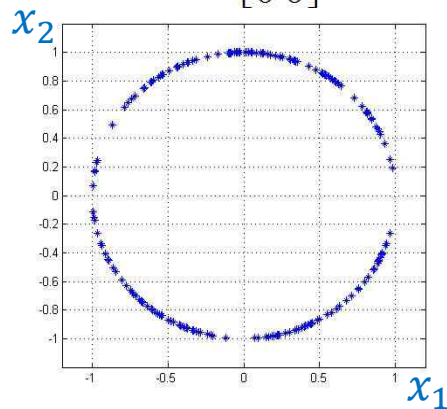
# 선형변환에 의한 특징추출

## ○ 변환행렬에 따른 특징의 분포

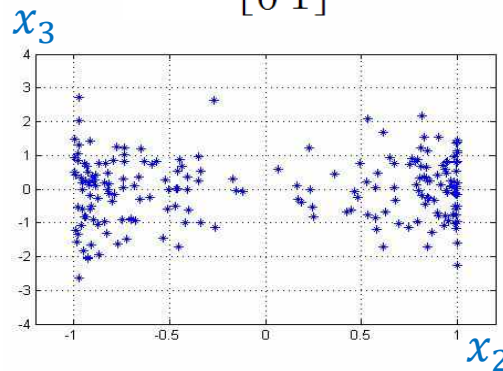
원래 데이터



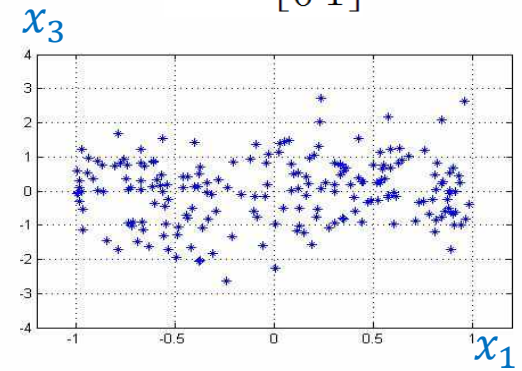
$$W_a = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix}$$



$$W_b = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$W_c = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \end{bmatrix}$$



### □ 좋은 특징추출이란?

✓ 변환행렬  $W$ 를 적절히 조절해서 분석 목적에 맞는 특징 분포를 만드는 것

→ 통계적 특징추출

## 통계적 특징추출

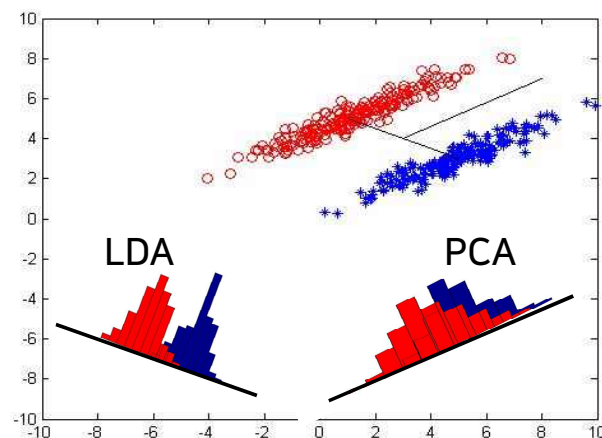
### ○ 선형변환을 사용하는 대표적인 통계적 특징추출 방법

□ 주성분분석법 Principal Component Analysis: PCA

✓ 클래스 정보 미사용 → 비지도 학습

□ 선형판별분석법 Linear Discriminant Analysis: LDA

✓ 클래스 정보 사용 → 지도 학습

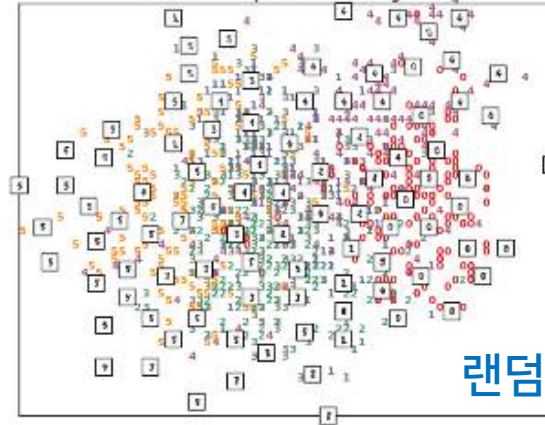


# 통계적 특징추출의 예

A selection from the 64-dimensional digits dataset

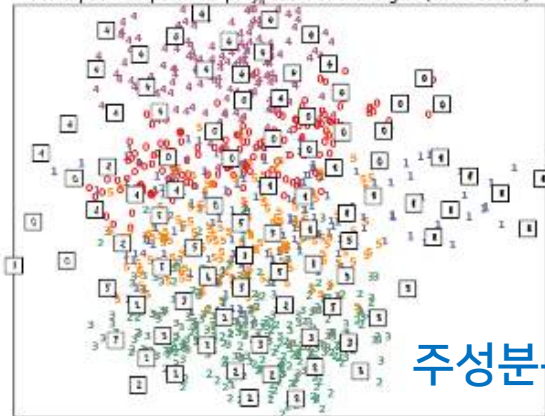


Random Projection of the digits



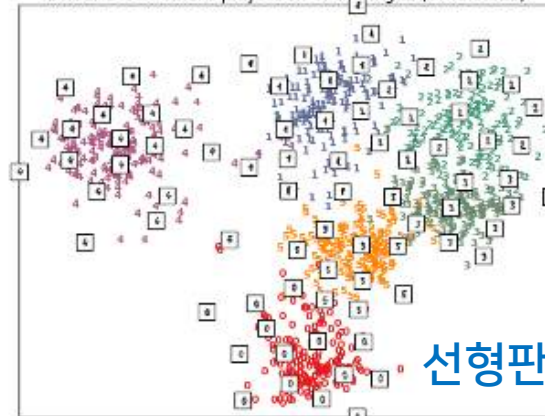
랜덤 사영

Principal Components projection of the digits (time 0.02s)



주성분분석법

Linear Discriminant projection of the digits (time 0.02s)



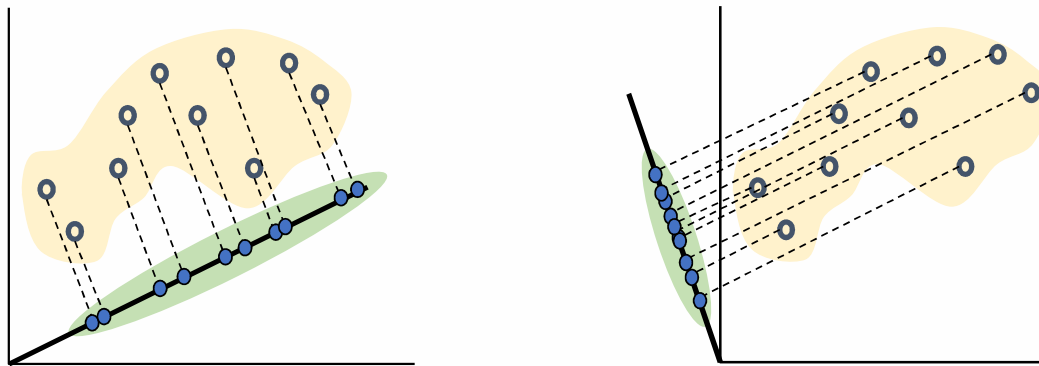
선형판별분석법

2

## 주성분분석법

# 주성분분석법

- 목적 → 변환 전의 데이터  $X$ 가 가지고 있는 정보를 차원 축소 후에도 최대한 유지



- 데이터 집합이 가능한 넓게 퍼질 수 있는 방향으로 사영을 수행
- 데이터의 분산이 가장 큰 방향으로의 선형변환을 수행
- 가장 큰 분산과 그 방향 = 공분산행렬의 최대 고유치와 고유벡터

→ 데이터의 공분산행렬의 고유치와 고유벡터를 찾아,  
고유치가 가장 큰 값부터 순서대로 이에 대응하는  $m$ 개의 고유벡터를 찾아서 행렬  $W$ 를 구성

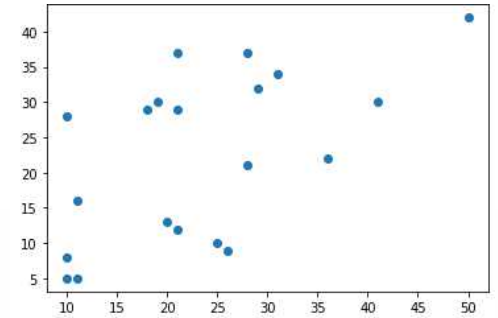


## PCA 알고리즘의 수행 단계

- ① 입력 데이터 집합  $X$ 의 평균  $\mu_x$ 와 공분산  $\Sigma_x$ 를 계산

$$\mu_x = \frac{1}{N} \sum_{i=1}^N x_i \quad \Sigma_x = \frac{1}{N} (X - M_x)(X - M_x)^T \quad M_x = \mu_x \mathbf{1}^T$$

(1은 모든 원소의 값이 1인  $N$ 차원 열벡터)



- ② 고유치 분석을 통해 공분산  $\Sigma_x$ 의 고유치행렬  $\Lambda$ 과 고유벡터행렬  $U$ 를 계산

고유치  $\lambda_1, \dots, \lambda_n$ 을 대각 성분으로 가지는 대각행렬

고유벡터  $u_1, \dots, u_n$ 을 열벡터로 가지는 행렬

$$\Sigma_x = U \Lambda U^T = [u_1, u_2, \dots, u_n] \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} [u_1, u_2, \dots, u_n]^T$$

[ 55.34980455 202.39493229]

[[-0.75801271 -0.65223979]  
[ 0.65223979 -0.75801271]]

## PCA 알고리즘의 수행 단계

- ③ 고유치가 큰 것부터 순서대로  $m$ 개의 고유치  $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$ 를 선택

```
[ 55.34980455 202.39493229]
```

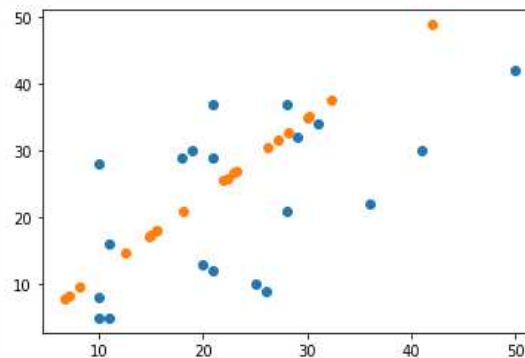
- ④ 선택한 고유치에 대응되는 고유벡터를 열벡터로 가지는 변환행렬  $W$ 를 생성

$$W = [u_1, u_2, \dots, u_m]$$

```
[[-0.75801271 -0.65223979]
 [ 0.65223979 -0.75801271]]
```

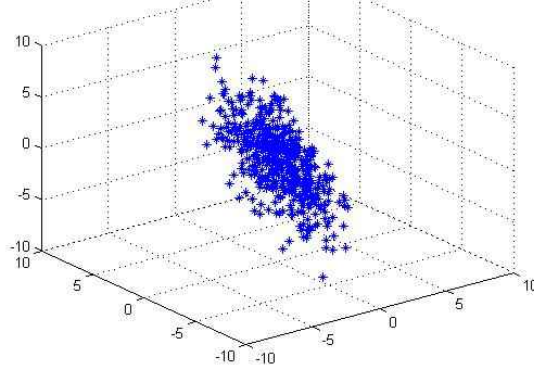
- ⑤  $W$ 에 의한 선형변환으로 특징 데이터  $Y$ 를 얻음

$$Y = W^T X$$

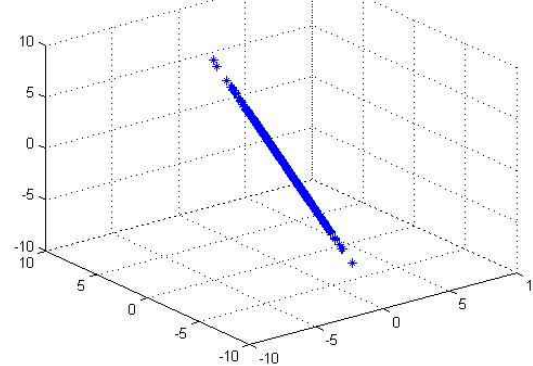


## PCA 적용의 예

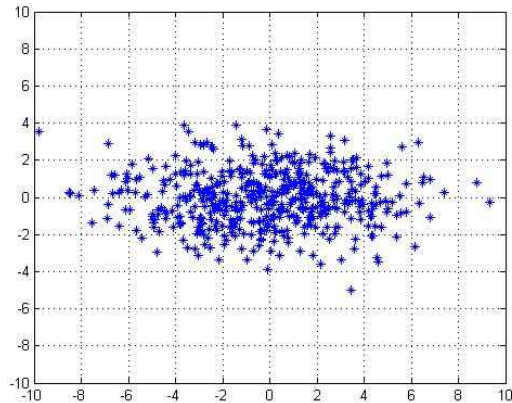
입력 데이터



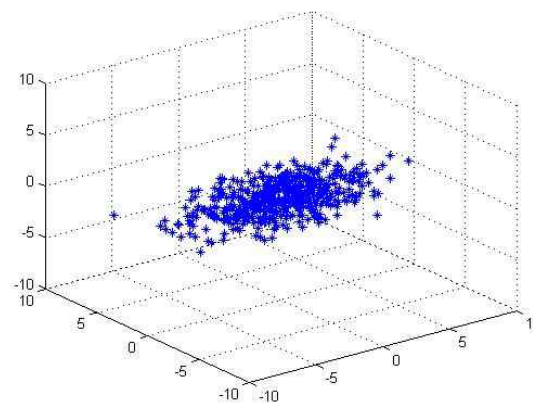
1차 주성분 벡터에 의한 특징추출



1,2차 주성분 벡터에 의한 특징추출



1,2,3차 주성분 벡터에 의한 특징추출



## PCA의 수학적 유도

### ○ 축소되는 차원 $m$ 을 선택하는 기준

손실되는 정보량의 비중  
( $m$ 개의 특징으로 표현 가능한 정보의 비율)

$$r(n, m) = \frac{\sum_{i=1}^m \lambda_i}{\sum_{i=1}^n \lambda_i}$$

선택할 고유벡터의 수(특징의 차수)  $m$

$$\sum_{i=1}^m \lambda_i / \sum_{i=1}^n \lambda_i > \theta$$

## [예] 얼굴 영상의 표현

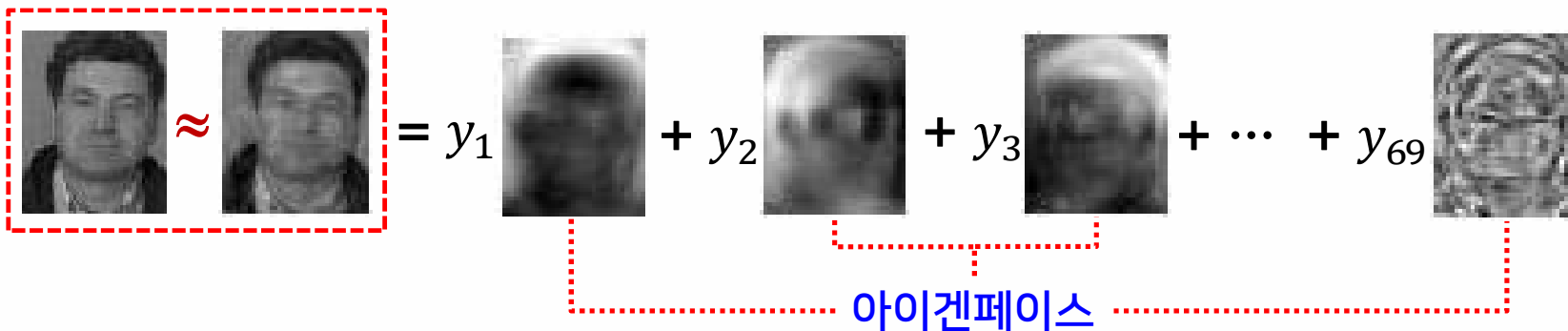
○ n차원의 얼굴 영상을 m개의 기저벡터를 사용해서 표현

□ n차원을 m차원으로 차원 축소 ( $n \gg m$ )

□ 좋은 기저벡터 → "Eigenface"

✓ Eigenface → 얼굴 영상에 PCA를 적용하여  
찾아진 고유벡터를 영상으로 표현한 것

$$x = Wy = y_1 w_1 + y_2 w_2 + \dots + y_m w_m$$

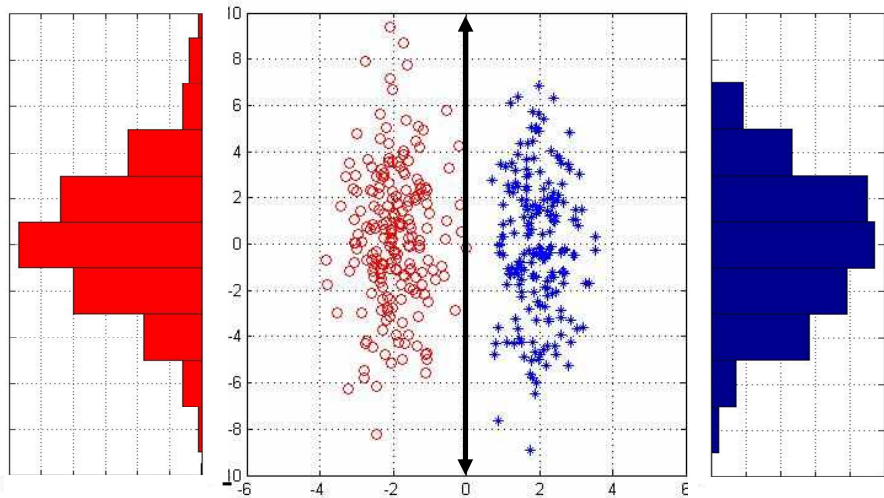


## 주성분분석법의 특성과 문제점

○ 데이터 분석에 대한 특별한 목적이 없는 경우에 가장 합리적인 차원 축소의 기준

○ 비지도학습

□ 클래스 레이블 정보를 활용하지 않음 → 분류의 핵심 정보의 손실 초래



↓  
분류에 적합한 특징추출?

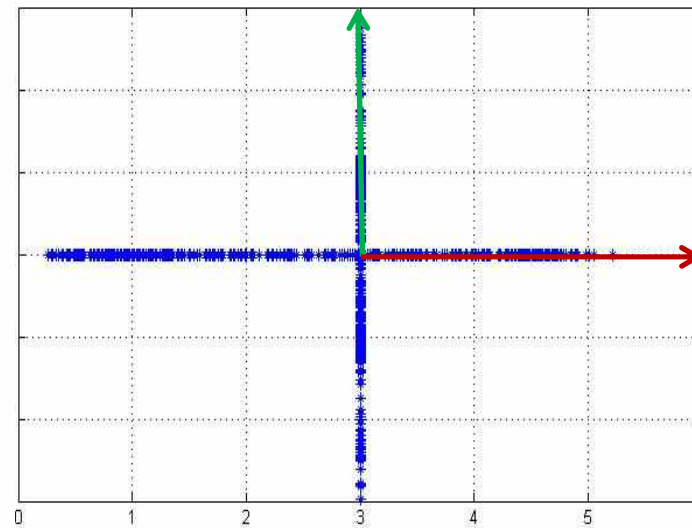
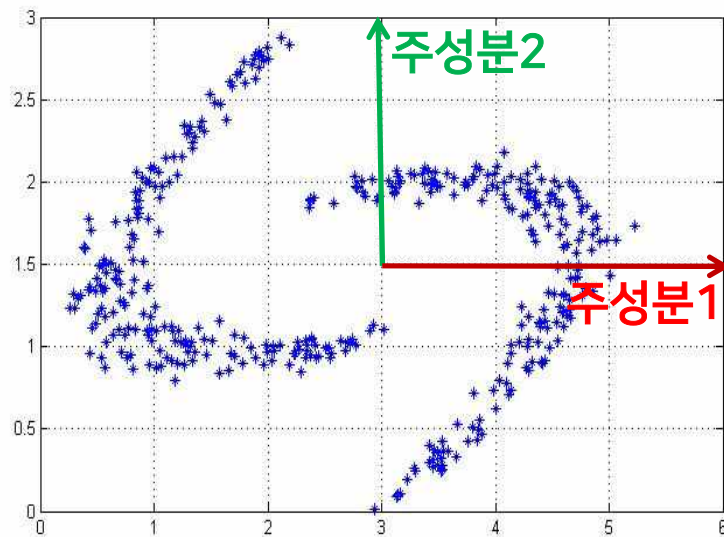
→ 데이터의 클래스 정보를 활용한 지도학습이 필요

→ "선형판별분석법(LDA)"

## 주성분분석법의 특성과 문제점

### ○ 선형변환의 한계

- 데이터의 비선형 구조를 반영하지 못함



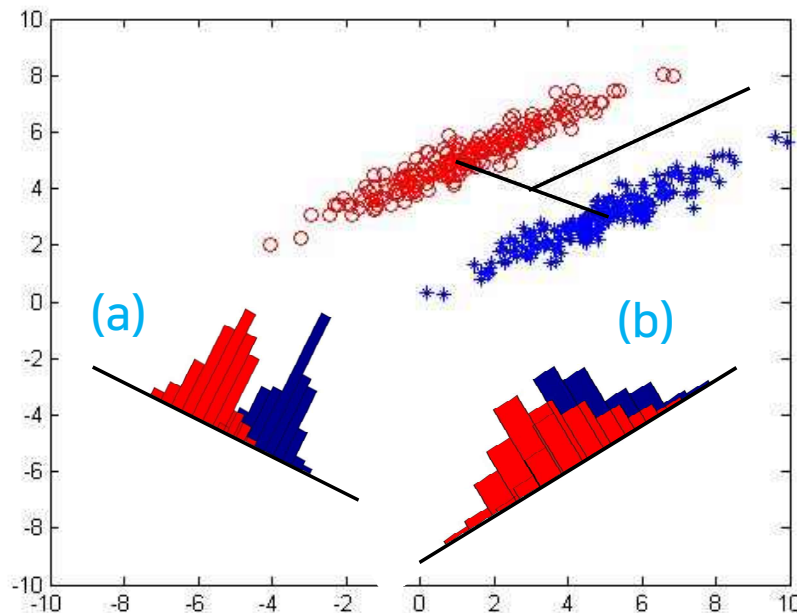
3

## 선형판별분석법



## 선형판별분석법

- 목적 → 클래스 레이블 정보를 적극 활용  
→ 클래스 간 판별이 잘 되는 방향으로 차원 축소



분류에 적합한 특징의 방향이란?

- 각 클래스가 가능한 서로 멀리 떨어질 수 있도록 거리를 유지

## LDA: 이진 분류 문제

클래스 간 거리에 대한 목적함수

$$J = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

$$m_k = \frac{1}{|C_k|} \sum_{x_i \in C_k} \mathbf{w}^T \mathbf{x}_i = \mathbf{w}^T \mathbf{m}_k$$

$$s_k = \sum_{x_i \in C_k} (\mathbf{w}^T \mathbf{x}_i - m_k)^2 = \sum_{x_i \in C_k} \mathbf{w}^T (\mathbf{x}_i - \mathbf{m}_k) (\mathbf{x}_i - \mathbf{m}_k)^T \mathbf{w}$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}}{\mathbf{w}^T \sum_{k=1}^2 \sum_{x_i \in C_k} (\mathbf{x}_i - \mathbf{m}_k) (\mathbf{x}_i - \mathbf{m}_k)^T \mathbf{w}} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

클래스 간 산점행렬  
between-scatter matrix

클래스 내 산점행렬  
within-scatter matrix

목적함수를 최대화하는 기저벡터  $\mathbf{w}$ 

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

## LDA: 다중 클래스 분류

**목적함수** → 각 클래스 내의 산점도는 작게, 클래스 간의 산점도는 크게

$$J(\mathbf{W}) = \text{Trace}\{(\mathbf{W}S_W\mathbf{W}^T)^{-1}(\mathbf{W}S_B\mathbf{W}^T)\} \quad \text{Trace} \rightarrow \text{정방행렬의 대각원소의 합을 계산하는 연산}$$

$$S_W = \sum_{k=1}^M S_k = \sum_{k=1}^M \sum_{x_i \in C_k} (x_i - \mathbf{m}_k)(x_i - \mathbf{m}_k)^T$$

$$S_B = \sum_{k=1}^M N_k(\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

$N_k \rightarrow$  클래스  $C_k$ 에 속한 데이터의 수  
 $\mathbf{m} \rightarrow$  전체 데이터 집합에 대한 평균

**목적함수를 최대화 하는 변환행렬  $\mathbf{W}$**

$S_W^{-1}S_B$ 의 고유벡터들을 열벡터로 가지는 행렬

$$S_W^{-1}S_B = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T \longrightarrow \mathbf{W} = [\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_m]$$

## 알고리즘의 수행 단계

- ① 입력데이터  $X$ 를 각 클래스 레이블에 따라  $M$ 개의 클래스로 나누어 각각 평균  $\mathbf{m}_k$ 와 클래스 간 산점행렬  $S_B$ , 그리고 클래스 내 산점행렬  $S_W$ 를 계산

$$\mathbf{m}_k = \frac{1}{N_k} \sum_{x_i \in C_k} x_i \quad (k = 1, 2, \dots, M)$$

$$S_W = \sum_{k=1}^M S_k = \sum_{k=1}^M \sum_{x_i \in C_k} (x_i - \mathbf{m}_k)(x_i - \mathbf{m}_k)^T$$

$$S_B = \sum_{k=1}^M N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

$N_k \rightarrow$  클래스  $C_k$ 에 속한 데이터의 수

$\mathbf{m} \rightarrow$  전체 데이터 집합에 대한 평균

## 알고리즘의 수행 단계

- ② 고유치 분석을 통해 행렬  $S_W^{-1}S_B$ 의 고유치행렬  $\Lambda$ 와 고유행렬벡터  $U$ 를 계산

$$S_W^{-1}S_B = U\Lambda U^T = [u_1, u_2, \dots, u_n] \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{bmatrix} [u_1, u_2, \dots, u_n]^T$$

- ③ 고유치가 큰 것부터 순서대로  $m$ 개의 고유치  $\{\lambda_1, \lambda_2, \dots, \lambda_m\}$ 를 선택

- ④ 선택한 고유치에 대응되는 고유벡터를 열벡터로 가지는 변환행렬  $W$ 를 생성

$$W = [u_1, u_2, \dots, u_m]$$

- ⑤  $W$ 에 의한 선형변환으로 특징 데이터  $Y$ 를 얻음

$$Y = W^T X$$

## 선형판별분석법의 특성과 문제점

- 지도학습 능력
- 선형변환의 한계
  - 복잡한 비선형 구조를 가진 경우에는 적절한 변환이 불가  
→ 커널법, 비선형 매니폴드 학습법 등을 이용한 접근 필요
- 선택하는 고유벡터의 개수(축소된 특징 차원)  $m$ 의 결정
  - 직접 분류를 통해 얻어지는 데이터에 대한 분류율을 기준으로 결정
  - 행렬  $S_W^{-1}S_B$ 에 의해 찾아지는 고유벡터의 개수가 제한
    - ✓ 클래스의 개수가  $M$ 이면 특징 벡터는 최대  $M - 1$ 차원으로 제한

## 선형판별분석법의 특성과 문제점

- 작은 표본집합의 문제 small sample set problem
  - 입력 데이터 수가 입력 차원보다 작은 경우
    - 클래스 내 산점행렬( $S_W$ )의 역행렬이 존재하지 않음
    - PCA로 먼저 차원 축소한 후, 이에 대해 LDA 적용

4

## 거리 기반 차원 축소 방법



## 거리 기반 차원 축소 방법

### ○ 기본 목적

□ 두 데이터 쌍 간의 거리를 최대한 유지하는 방향으로 차원 축소

✓ 원래 데이터  $\rightarrow \{x_1, x_2, \dots, x_n\}$

✓ 추출된 저차원의 특징  $\rightarrow \{y_1, y_2, \dots, y_n\}$  ( $y_i = f(x_i)$ )

✓ 원래 데이터의 거리행렬  $\rightarrow D = \{d_{ij}\}$   $d_{ij} = \text{dist}(x_i, x_j)$

✓ 특징 벡터의 거리행렬  $\rightarrow \Delta = \{\delta_{ij}\}$   $\delta_{ij} = \text{dist}(y_i, y_j)$

✓ 목적  $\rightarrow \sum_{i < j} (d_{ij} - \delta_{ij})^2$ 의 최소화

## 거리 기반 차원 축소 방법

### ○ 거리의 정의에 따라 다양한 방법이 존재

#### ☐ 다차원 척도법 Multi-Dimensional Scaling: MDS

✓ 유클리디안 거리 사용

#### ☐ t-SNE

✓ 확률밀도함수를 활용하여 거리를 정의

#### ☐ Isomap

✓ 측지 거리 geodesic distance 사용

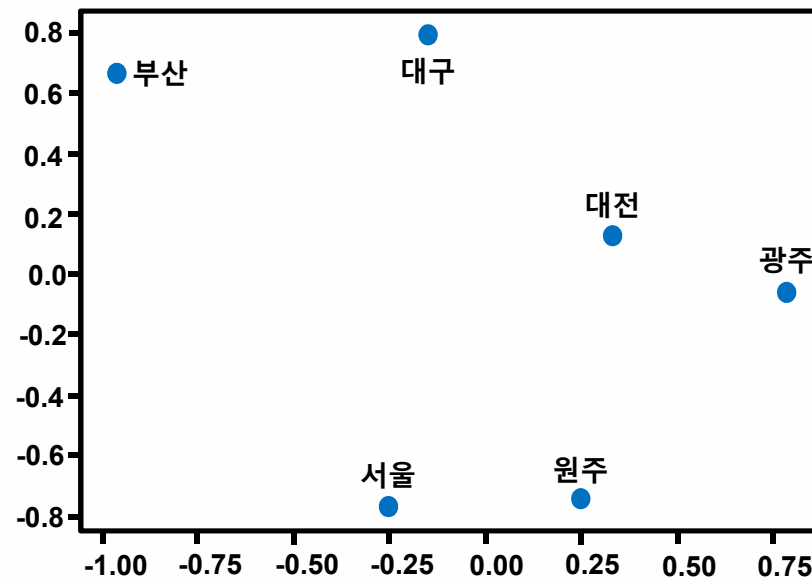
## 다차원 척도법

○ 거리행렬 D가 값으로 정의되거나 유클리디안 거리 사용

$$d_{ij} = \text{dist}(x_i, x_j) = \|x_i - x_j\|^2$$

□ 목적  $\rightarrow \sum_{i,j} (d_{ij} - \delta_{ij})^2$ 를 최소화하는 특징  $\{y_1, y_2, \dots, y_n\}$ 를 찾는 것

D	서울	대구	대전	광주	원주	부산
서울	0	350	239	330	187	480
대구	350	0	123	283	379	110
대전	239	123	0	186	135	390
광주	330	283	186	0	274	400
원주	187	379	135	274	0	450
부산	480	110	390	400	450	0



## t-SNE

○ t-통계적 이웃 임베딩 t-Stochastic Neighbor Embedding

- 데이터 간의 거리와 특징 간의 거리를 조건부확률을 이용한 유사도로 정의

$$d_{ij} = \text{similarity}(\mathbf{x}_i, \mathbf{x}_j) = p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)}$$

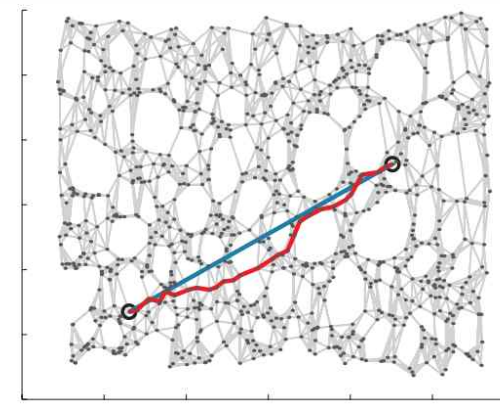
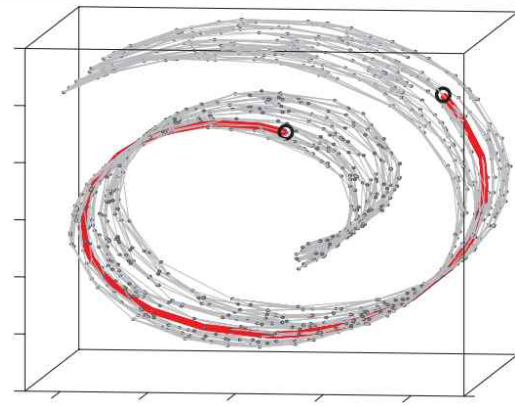
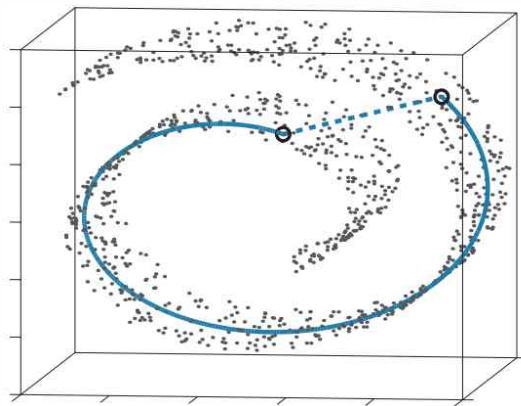
$$\delta_{ij} = \text{similarity}(\mathbf{y}_i, \mathbf{y}_j) = q_{j|i} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq i} (1 + \|\mathbf{y}_i - \mathbf{y}_k\|^2)^{-1}}$$

- 거리가 멀리 떨어진 데이터 사이의 관계를 더 잘 반영

# Isomap

## ○ 측지 거리 사용

- 데이터들을 정점으로 가지는 그래프 간의 경로를  
데이크스트라 Dijkstra 알고리즘으로 계산



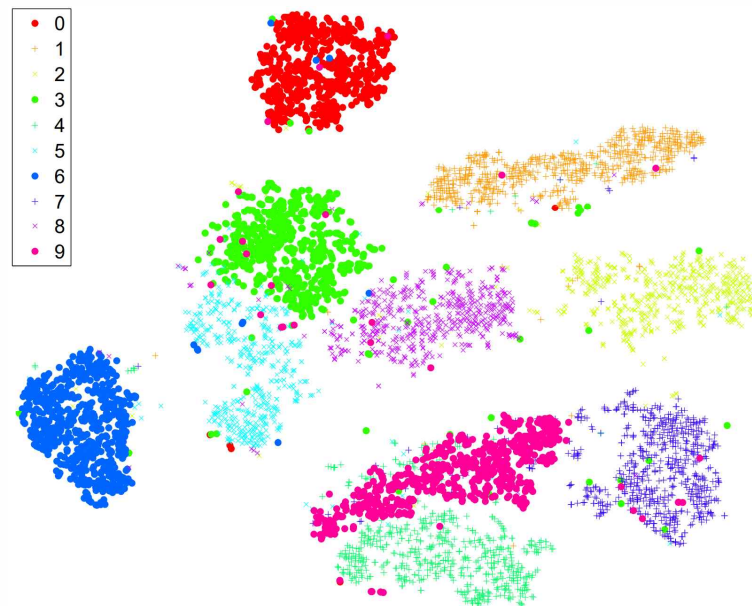
DOI: 10.1126/science.290.5500.2319

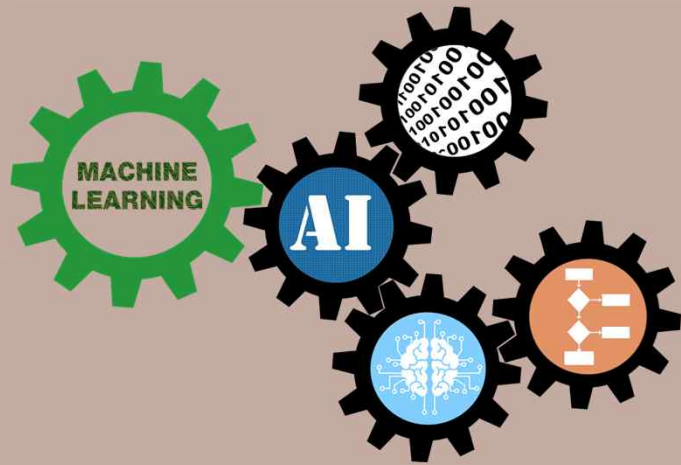
## 거리 기반 차원 축소 방법의 특징

- 입력 데이터와 특징 데이터 간의 매핑 함수를 정의하지 않음
  - 새로운 데이터에 대해서는 그에 대응하는 특징값을 찾을 수 없음
  - 데이터 시각화의 용도로 주로 사용

0 0 0 0 0 0 0 0 0 0 0 0 0 0  
 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
 2 2 2 2 2 2 2 2 2 2 2 2 2 2  
 3 3 3 3 3 3 3 3 3 3 3 3 3 3  
 4 4 4 4 4 4 4 4 4 4 4 4 4 4  
 5 5 5 5 5 5 5 5 5 5 5 5 5 5  
 6 6 6 6 6 6 6 6 6 6 6 6 6 6  
 7 7 7 7 7 7 7 7 7 7 7 7 7 7  
 8 8 8 8 8 8 8 8 8 8 8 8 8 8  
 9 9 9 9 9 9 9 9 9 9 9 9 9 9

t-SNE





다음시간안내

## 제6강

# 앙상블 학습