## 선형대수학팀

## 3팀

김다민 이지원 조성우 김수인 방건우

## CONTENTS

1.주성분 분석

2.고차원 시각화

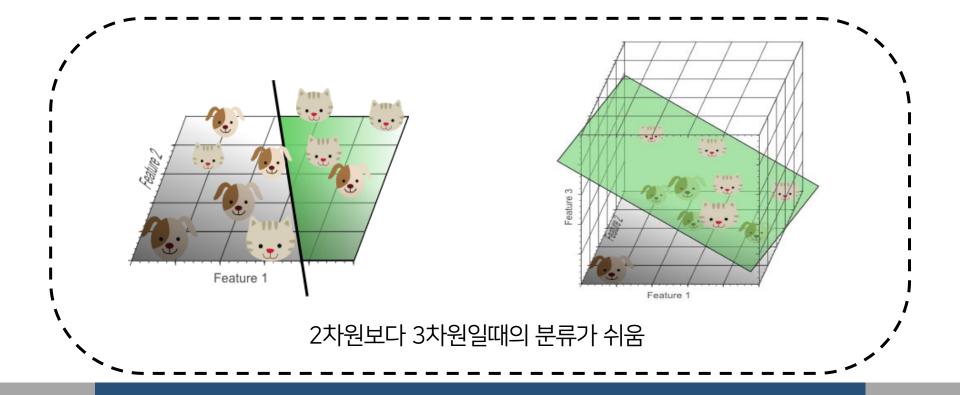
## 1

주성분 분석

## 차원의 저주

차원의 축복(?)

데이터 분석에서 변수의 추가, 즉 차원의 증가는 모델의 성능 향상에 도움이 됨



#### 차원의 저주

차원의 축복(?)



데이터 분석에서 변수의 추가 즉 차원이 증가는 모델의 성능 호상에 도움이 된 그렇다면 사용가능한 모든 변수들을 통해 모델링하면

성능이 뛰어날까?

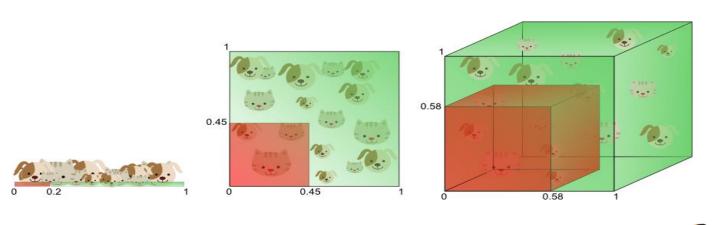


2차원보다 3차원일때의 분류가 쉬움

#### 차원의 저주

차원의 저주 개념

차원이 증가함에 따라 학습 알고리즘이 제대로 작동하지 않는 현상



차원이 증가 ▶ 공간의 크기가 증가 ▶ 데이터간 거리가 멀어짐 적은 데이터로만 이 공간을 표현하는 경우 **모델 성능 저하** & **계산량 증가** 

#### 차원의 저주



차원의 저주를 해결하기 위해서는 차원의 개수를 조절해야 함

:

회귀분석의 **변수 선택**과 데이터마이닝의 **변수 필터링** 등 다양한 **차원축소** 기법 존재

본 클린업에서는 선형대수학 기반 다변량 기법인 PCA에 대해 다툳

#### 차원의 저주



차원의 저주를 해결하기 위해서는 차원의 개수를 조절해야 함

:

회귀분석의 **변수 선택**과 데이터마이닝의 **변수 필터링** 등 다양한 **차원축소** 기법 존재

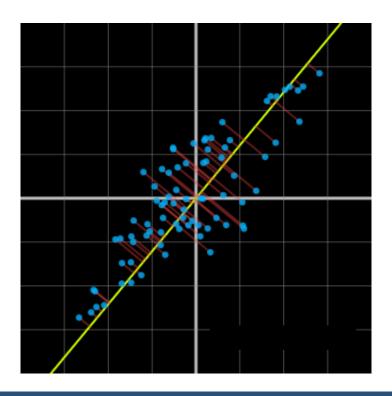


본 클린업에서는 선형대수학 기반 다변량 기법인 PCA에 대해 다룸

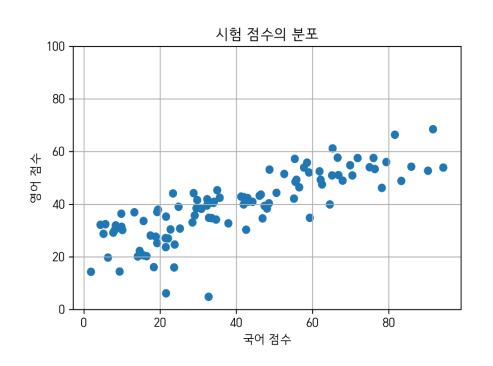
## 주성분 분석 *PCA*

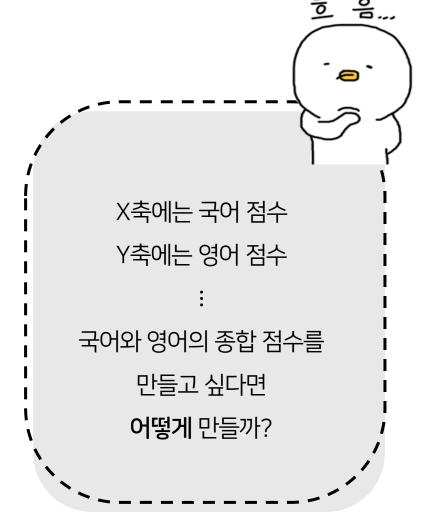
주성분 분석의 개념

데이터 분석에서 자주 사용되는 **차원 축소** 기법 원 데이터의 **분산**을 최대한 **보존**하는 방향으로 변수 변형



#### 주성분 분석 *PCA*





#### 주성분 분석 *PCA*

산술평균 (동등한 가중치)

이를 벡터의 내적으로 표현하면

$$\begin{bmatrix} 80\\60 \end{bmatrix} \cdot \begin{bmatrix} 0.5\\0.5 \end{bmatrix}$$

#### 주성분 분석 *PCA*

산술평균 (동등한 가중치)

국어80점, 수학 60점이라면 80 × 05 + 60 × 05

이를 벡터의 내적으로 표현하면

 $\begin{bmatrix} 80 \\ 60 \end{bmatrix} \cdot \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$ 

가중평균 (한 과목에 가중치)

같은 점수에

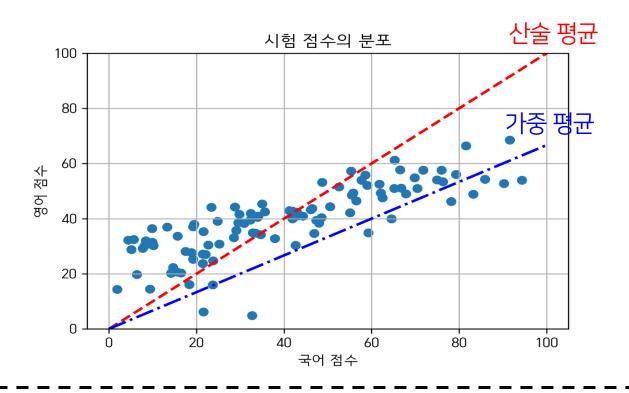
6:4 가중치를 적용하면,

 $80 \times 0.6 + 60 \times 0.4$ 

내적으로 표현하면

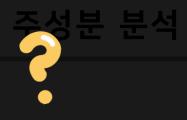
 $\begin{bmatrix} 80 \\ 60 \end{bmatrix} \cdot \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$ 

#### 주성분 분석 *PCA*



종합 점수를 얻는 방식을 수학적으로는

점수 벡터를 비율 벡터에 내적(정사영)하는 문제로 생각



주성분 분석 PCA

#### 두가지 고민

시험 점수의 분포

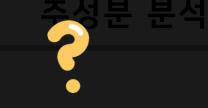
산술 평균

- ① 데이터를 어떤 벡터에 내적하는 것이 최적의 결과를 줄 것인가?
- ② 내적의 대상이 될 벡터를 찾을 때, 데이터 분포의 중심을 축으로 하는

벡터를 찾으면 좋지 않을까?

종합 점수를 얻는 방식을 수학적으로는

점수 벡터를 비율 벡터에 내적(정사영)하는 문제로 생각



주성분 분석 PCA

#### 두가지 고민

시헌 전수의 부포

산술 평균

① 데이터를 어떤 벡터에 내적하는 것이 최적의 결과를 줄 것인가?

② 내적의 대상이 될 벡터를 찾을 때, 데이터 분포의 중심을 축으로 하는 벡터를 찾으면 좋지 않을까?



공분산 행렬로부터 답을 찾을 수 있음

점수 벡터를 비율 벡터에 내적(정사영)하는 문제로 생각

#### 공분산 행렬

공분산 행렬 Covariance matrix

변수들의 공분산을 행렬로 나타낸 것

$$cov(X,Y) = E[(X - \mu_x)(Y - \mu_y)^T = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{m1} & \cdots & \sigma_{mn} \end{bmatrix}$$

#### <del>공분</del>산 행렬

공분산 행렬의 수식적 의미

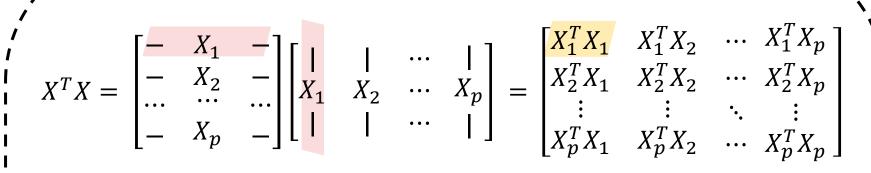
n개의 관찰값과 p개의 변수가 있는 데이터 행렬 X를 생각 행렬 X의 각 열들은 **정규화**된 상태여야 함

$$X = \begin{bmatrix} | & | & \cdots & | \\ X_1 & X_2 & \cdots & X_p \\ | & | & \cdots & | \end{bmatrix} \in R^{n \times p}$$

왜 정규화를 해야 할까?

- ① 데이터의 분포가 원점에서 뻗어나가는 벡터 공간 형태를 만들어 데이터의 추세 파악
  - ② 변수 간 조건을 동등하게 해놓고 정보량을 비교하는 것이 목적

#### 공분산 행렬



각 열벡터는 편차벡터이므로 내적은 편차 제곱의 합

:

데이터의 크기(n)가 많을 수록 편차 제곱의 합은 커지므로 위 행렬을 n으로 나누면 **공분산행렬** 완성

$$\mathbf{\Sigma} = \frac{1}{n} X^T X$$

#### 공분산 행렬

공분산 행렬의 기하학적 의미

공분산 행렬은 **데이터의 구조**를 설명하며 **변수 간 관계** 파악에 도움

행렬은 서형변화

공분산 행렬을 통해 데이터를 변환하면 분산과 공분산만큼 공간이 변회

x축 방향으로 퍼진 정도

3

41

/축 방향으로 퍼진 정도

#### 공분산 행렬

공분산 행렬의 기하학적 의미

공분산 행렬은 **데이터의 구조**를 설명하며 **변수 간 관계** 파악에 도움

#### 행렬은 **선형변환**

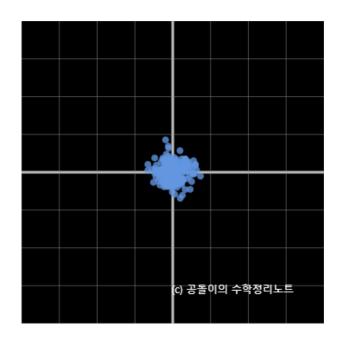
공분산 행렬을 통해 데이터를 변환하면 분산과 공분산만큼 공간이 변화



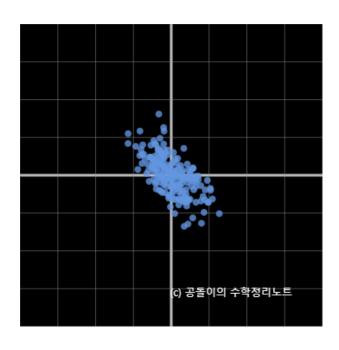
#### <del>공분</del>산 행렬

공분산 행렬  $\begin{bmatrix} 3 & -2 \\ -2 & 4 \end{bmatrix}$ 을 통해 데이터를 **변환** 

데이터가 어느 방향으로 어떻게 분포되어 있는지 알 수 있음











#### 데이터의 분산과 정보량의 관계

차원축소는 정보의 손실을 동반하므로

중요한 정보를 보존하면서 차원을 축소해야 함

분산을 통해 변수의 정보량 파악 가능

항상 정답은 아니지만

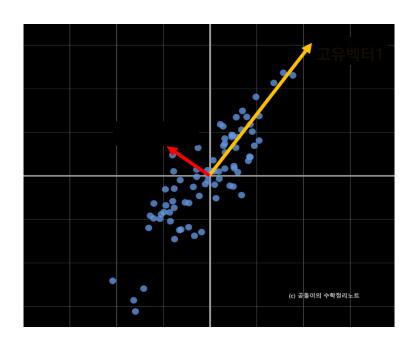
설명변수의 분산이 크면 좋은 예측이 가능

1

#### 주성분 분석

### 분산의 보존 주성분 찾기

정보를 최대한 보존하는 선에서 차원을 축소하려면 **분산이 큰 방향**으로 데이터들을 정사영



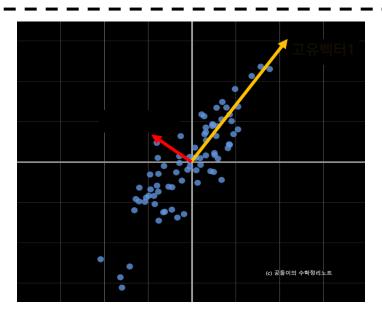
그림에서 <mark>노란색 벡터</mark>가 데이터의 분산을 제일 잘 보존



공분산 행렬을 통해 찾을 수 있음

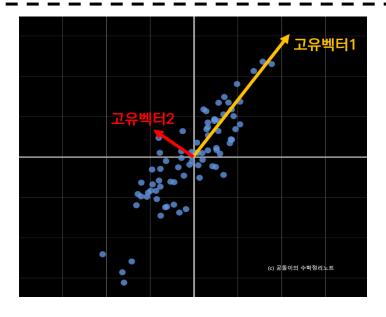
### 분산의 보존

공분산 행렬이라는 선형변환을 적용했을 때 노란색 벡터를 찾는 것은 데이터를 **어느 방향**으로 **얼마나** 잡아당겼냐의 문제



#### 분산의 보존

공분산 행렬이라는 선형변환을 적용했을 때 노란색 벡터를 찾는 것은 데이터를 **어느 방향**으로 **얼마나** 잡아당겼냐의 문제



해당 벡터는 공분산 행렬의 고유벡터

#### 분산의 보존

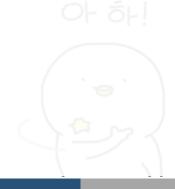
고유벡터는 선형변환 시 방향은 유지한 채로 길이만 변하는 벡터이므로 공분산 행렬의 고유값을 통해 잡아당긴 정도인

분산의 크고 작음을 알 수 있음

선대팀 2주차 클린업 참고



고유값이 큰 순서대로 고유벡터를 정렬하면 분산이 큰 순서대로 주성분을 구하는 것



#### 분산의 보존

고유벡터는 선형변환 시 방향은 유지한 채로 길이만 변하는 벡터이므로 공분산 행렬의 고유값을 통해 잡아당긴 정도인

분산의 크고 작음을 알 수 있음

선대팀 2주차 클린업 참고



고유값이 큰 순서대로 고유벡터를 정렬하면 분산이 큰 순서대로 주성분을 구하는 것! 아하!



#### 주성분 선택

# p차원 데이터를 m차원까지 선택하는 주성분 분석 진행할 때 고유값은 p개이며 $\lambda_1,\lambda_2,\dots,\lambda_p$ 로 표현 가능

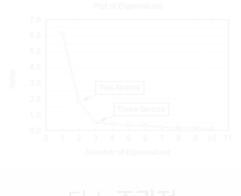
공분산 행렬이 full rank임을 가정

전체 데이터의 90% 설명

전체 데이터의 분산 중 90%만큼 설명하는 차원까지 감소

$$\frac{\sum_{j=1}^{m} \lambda_j}{\sum_{i=1}^{p} \lambda_i} \ge 0.9$$

Elbow Point Scree plot의 elbow point를 사용



다소 주관적

#### 주성분 선택

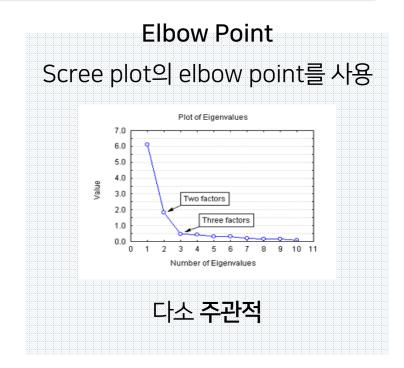
## p차원 데이터를 m차원까지 선택하는 주성분 분석 진행할 때 고유값은 p개이며 $\lambda_1,\lambda_2,\dots,\lambda_p$ 로 표현 가능

공분산 행렬이 full rank임을 가정

전체 데이터의 90% 설명

전체 데이터의 분산 중 90%만큼 설명하는 차원까지 감소

$$\frac{\sum_{j=1}^{m} \lambda_j}{\sum_{i=1}^{p} \lambda_i} \ge 0.9$$



#### 주성분 분석 *PCA* PCA 요약

데이터를 설명하는 **차원 축을 변경**하고, 그 중 **많은 정보를 가진 축**만 남기는 것

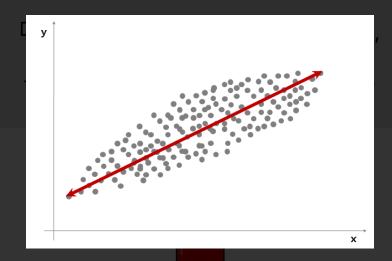


원래 데이터에서 명시되었던

변수 축의 변경으로 인해 해석력 상실

## 주성분 분석<u>그렇다면 언제 PCA를 사용</u>해야 할까?

PCA 요약



빨간색 벡터 한 축으로 데이터를 표현해도,

소실되는 정보량이 적어 PCA 사용 적합 목을 변경시켜 된래 데이터에서 명시되었던

변수들의 축의 변경으로 인한 해석력을 상실

다중공선성 문제가 심각할수록 잘 작동함

# 2

## 고차원 시각화

#### 고차원 시각화

### 고차원 시각화

고차원 시각화의 필요성

고차원의 데이터를 이용할 때, 클러스터링 결과 등의 시각화는 어려움 해석이 어렵고 그림이 복잡해지기 때문



고차원 시각화 이용!



#### 고차원 시각화



#### 고차원 시각화

#### 고차원 시각화의 필요성 고차원 시각화 주의사항

고차원의 데이터를 이용할 때, 클러스터링 결과 등의 시각화는 어려움

고차원 시각화는 고차원 데이터를 저차원으로 임베딩

PCA를 제외한 방법들은 저차원의 좌표 정보에만 관심

→ 각 데이터포인트들의 <mark>상호 거리</mark>만이 중요 고차원 시각화 이용!

앞으로 나올 x는 하나의 데이터포인트를 의미!

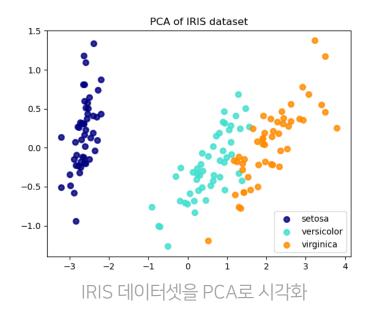


#### 고차원 시각화

#### PCA *Principal Component Analysis*

**PCA** 

가장 기초적인 고차원 시각화 기법 분산이 큰 상위 2개 축을 선정하여 그들을 축으로 하는 시각화를 진행



### PCA Principal Component Inalysis

PCA 기반 시각화의 아이디어

**PCA** 

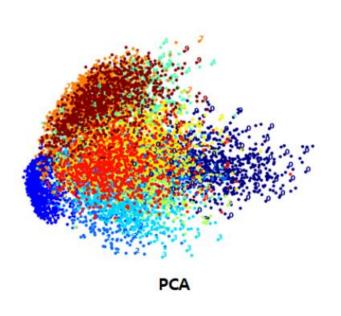
가장 기초적인 고차원 시각화 기법 시강화는 2차원 혹은 3차원 공간에서 원활히 상용됨

→ 차원축소 알고리즘인 PCA를 통해 3개 이내의 차원으로 축소하여 시각화 가능!



### PCA *Principal Component Analysis*

MNIST 데이터 시각화 with PCA



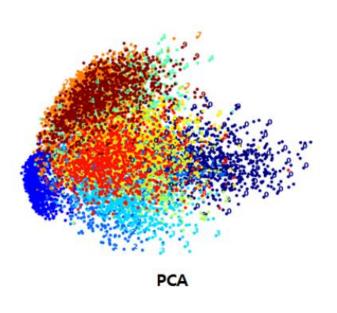
MNIST 데이터를 잘 분리하지는 못했지만, 시각화를 **행렬분해 기반**으로 가능케 함



이론적으로 기반이 단단하며, 차원 축소 아이디어를 통해 시**각화를 해냈다는 사실**에 의의가 있음

### PCA *Principal Component Analysis*

MNIST 데이터 시각화 with PCA



MNIST 데이터를 잘 분리하지는 못했지만, 시각화를 **행렬분해 기반**으로 가능케 함



이론적으로 기반이 탄탄하며, 차원 축소 아이디어를 통해 **시각화를 해냈다는 사실**에 의의가 있음

### PCA *Principal Component Analysis*

MNIST 데이터 시각화 with PCA



어쨌든 해냈다는 사실이 자랑스러운 PCA

### PCA *Principal Component Analysis*

PCA 기반 시각화의 특징

장점

**행렬분해**를 기반으로 함 차원이 크지 않고, 변수의 **분산이 데이터들의 특징을 잘 구별**할 경우 나쁘지 않은 성능 보임

단점

데이터의 유사도가 아닌, 단순히 **분산이 큰 축을 기준**으로 차원 축소 → 시각화가 잘 되지는 않음



## PCA *Principal Component Analysis*

PCA 기반 시각화의 특징

장점

**행렬분해**를 기반으로 함 차원이 크지 않고, 변수의 **분산이 데이터들의 특징을 잘 구별**할 경우 나쁘지 않은 성능 보임

단점

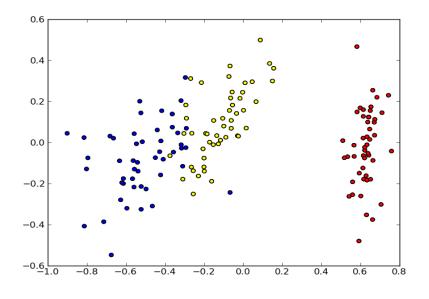
데이터의 유사도가 아닌, 단순히 **분산이 큰 축을 기준**으로 차원 축소 → 시각화가 잘 되지는 않음



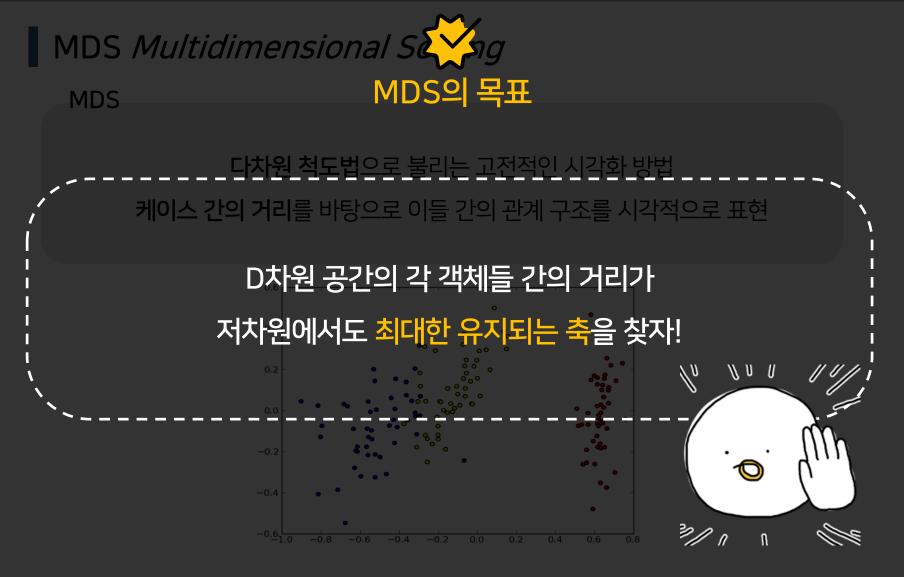
# MDS *Multidimensional Scaling*

**MDS** 

다차원 척도법으로 불리는 고전적인 시각화 방법 케이스 간의 거리를 바탕으로 이들 간의 관계 구조를 시각적으로 표현



IRIS 데이터셋을 MDS로 시각화

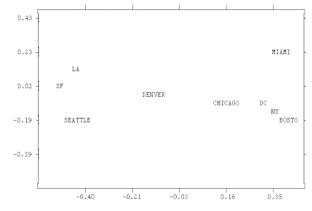


IRIS 데이터셋을 MDS로 시각화

# MDS *Multidimensional Scaling*

| Hand |





미국 도시 간 거리 행렬

MDS를 학습한 결과



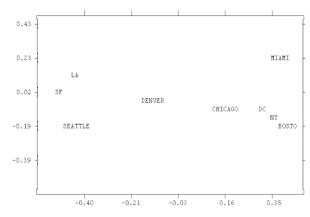
거리 행렬이 2차원 좌표로부터 계산된 정보이기 때문에 **2차원 지도가 온전히 복원**됨



단, 새로 만들어진 공간은 거리 정보만 반영했기 때문에 임의의 방향으로 **회전**될 수 있음

# MDS *Multidimensional Scaling*





미국 도시 간 거리 행렬

MDS를 학습한 결과



거리 행렬을 사용하는 방법론이지만, 실제 거리 뿐만 아니라 변수 간의 차이를 통해 거리를 계산할 수 있기 때문에 활용도가 높은 방법론!

### MDS *Multidimensional Scaling*

MDS 알고리즘



d차원 데이터 X에서 **유사도** 혹은 **거리 행렬** D를 구하고,

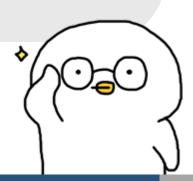


그 정보를 **내적 행렬** *B*에 넣은 뒤,



원하는 p차원 행렬  $\tilde{X}$ 를 구현

이때 행렬 D만 필요로 하기 때문에 원데이터 X는 필수가 아님



# MDS *Multidimensional Scaling*

① 유사도 혹은 거리 행렬 D 구하기

객체 간의 좌표가 존재한다면, 행렬 D의 원소는 다음과 같은 조건을 만족

$$(1) d_{ij} \ge 0$$

(2) 
$$d_{ii} = 0$$

$$(3) d_{ij} = d_{ji}$$

In addition to (1), (2), (3), it satisfies  $d_{ij} \le d_{ik} + d_{kj}$ 

#### 이때 다음의 거리 혹은 유사도 행렬을 고려 가능

거리 행렬: Euclidean, Manhattan, etc.

유사도 행렬: Correlation, Jaccard, etc.

자카드 유사5

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

교집합을 합집합으로 나눈 것

## MDS *Multidimensional Scaling*

① 유사도 혹은 거리 행렬 D 구하기

객체 간의 좌표가 존재한다면, 행렬 D의 원소는 다음과 같은 조건을 만족

$$(1) d_{ij} \ge 0$$

(2) 
$$d_{ii} = 0$$

$$(3) d_{ij} = d_{ji}$$

In addition to (1), (2), (3), it satisfies  $d_{ij} \le d_{ik} + d_{kj}$ 

#### 이때 다음의 거리 혹은 유사도 행렬을 고려 가능

거리 행렬: Euclidean, Manhattan, etc.

유사도 행렬: Correlation, Jaccard, etc.

자카드 유사도

$$J(A,B) = \frac{|A \cap B|}{|A \cup B|}$$

교집합을 합집합으로 나눈 것!

## MDS *Multidimensional Scaling*

① 유사도 혹은 거리 행렬 D 구하기



데이터 X로부터 행렬 D를 구하는 것을 그림으로 표현하면 다음과 같음

**X** (d by n)

	X <sub>1</sub>	X <sub>2</sub>	<b>X</b> <sub>3</sub>	X <sub>4</sub>	 x <sub>n</sub>
V <sub>1</sub>					
<b>V</b> <sub>2</sub>					
<b>v</b> <sub>3</sub>					
$\mathbf{v}_{d}$					

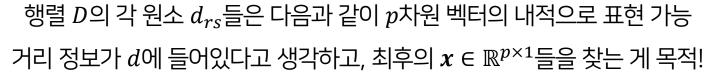
7

**D** (n by n)

	X <sub>1</sub>	X <sub>2</sub>	<b>X</b> <sub>3</sub>	X <sub>4</sub>		x <sub>n</sub>			
X <sub>1</sub>									
X <sub>2</sub>									
<b>x</b> <sub>3</sub>									
X <sub>4</sub>									
$\mathbf{x}_{n}$									

# MDS *Multidimensional Scaling*

② 내적 행렬 *B* 만들기



$$d_{rs}^2 = (\boldsymbol{x}_r - \boldsymbol{x}_s)^T (\boldsymbol{x}_r - \boldsymbol{x}_s)$$

내적 행렬 B는 거리행렬 D로부터 다음과 같이 얻음

$$[B]_{rs} = b_{rs} = \mathbf{x}_r^T \mathbf{x}_s$$

단, 위의 x는 원데이터의 열벡터가 아니라, 구현할  $d \times n$ 행렬  $\tilde{X}$ 의 열벡터

## MDS *Multidimensional Scaling*

② 내적 행렬 B 만들기

행렬 D의 각 원소  $d_{rs}$ 들은 다음과 같이 p차원 벡터의 내적으로 표현 가능거리 정보가 d에 들어있다고 생각하고, 최후의  $x\in\mathbb{R}^{p\times 1}$ 들을 찾는 게 목적!

$$d_{rs}^2 = (\boldsymbol{x}_r - \boldsymbol{x}_s)^T (\boldsymbol{x}_r - \boldsymbol{x}_s)$$

내적 행렬 B는 거리행렬 D로부터 다음과 같이 얻음

$$[B]_{rs} = b_{rs} = \boldsymbol{x}_r^T \boldsymbol{x}_s$$

단, 위의 x는 원데이터의 열벡터가 아니라, 구현할  $d \times n$ 행렬  $\tilde{X}$ 의 열벡터

# MDS Multidimensional Scaning

② 주어진 레-dro 밖에 없는데, B의 원소들을 어떻게 구할까?

행렬 D의 각 원소  $d_{rs}$ 들은 다음과 같이 p차원 벡터의 내적으로 표현 가능 거리 정보가 d물 <mark>라지, 가정광 복잡한 원소인 필요함멸</mark> 찾는 게 목적!

정보가 부족하기에 각 변수별 평균이 0이라는 가정 필요

$$\sum_{r=1}^{n} x_{ri} = 0, (i = 1, 2, \dots, p)$$

내적 행렬 B는 거리행렬 D로부터 다음과 같이 얻음

 $[B]_{rs} = b_{rs} = \mathbf{x}_r^T \mathbf{x}_s$ 

「내적행렬 B는 거리행렬 D로부터 다음과 같이 얻음



# MDS *Multidimensional Scaling*

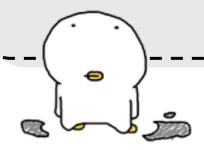
② 내적 행렬 B 만들기

일련의 연산 과정을 거치면

:

$$b_{rs} = -\frac{1}{2} \left( d_{rs}^2 - \frac{1}{n} \sum_{s=1}^n d_{rs}^2 - \frac{1}{n} \sum_{r=1}^n d_{rs}^2 - \frac{1}{n^2} \sum_{r=1}^n \sum_{s=1}^n d_{rs}^2 \right)$$

인  $n \times n$  행렬 B의 원소를 구할 수 있음



## MDS *Multidimensional Scaling*

③ 거리 정보를 보존하는 좌표 추출



구해진 행렬 B를 바탕으로  $p \times n$  행렬인  $\tilde{X}$ 를 구해보자!

$$B = \tilde{X}\tilde{X}^T$$
,  $rank(B) = rank(\tilde{X}\tilde{X}^T) = rank(\tilde{X}) = p$ 

이때 B는 대칭행렬이고, positive semi-definite이며 rank가 p임



 $\rightarrow p$ 개의 음이 아닌 고윳값, n-p개의 0인 고윳값을 구할 수 있음 (by EVD)  $B = V\Lambda V^T, \Lambda = diag(\lambda_1, \cdots, \lambda_p), V = \begin{bmatrix} \boldsymbol{v}_1, \cdots, \ \boldsymbol{v}_p \end{bmatrix}$ 

### MDS *Multidimensional Scaling*

③ 거리 정보를 보존하는 좌표 추출



구해진 행렬 B를 바탕으로  $p \times n$  행렬인  $\tilde{X}$ 를 구해보자!

$$B=\tilde{X}\tilde{X}^T,\ rank(B)=rank\big(\tilde{X}\tilde{X}^T\big)=\mathrm{rank}\big(\tilde{X}\big)=p$$
이때  $B$ 는 대칭행렬이고, positive semi-definite이며  $p$ 임



 $\rightarrow p$ 개의 음이 아닌 고윳값, n-p개의 0인 고윳값을 구할 수 있음 (by EVD)

$$B = V \Lambda V^T, \Lambda = diag(\lambda_1, \dots, \lambda_p), V = [v_1, \dots, v_p]$$

### MDS *Multidimensional Scaling*

③ 거리 정보를 보존하는 좌표 추출



구해진 행렬 B를 바탕으로  $p \times n$  행렬인  $\tilde{X}$ 를 구해보자!

0인 고윳값이 n-p개 있으므로, B는 다음과 같이 표현 가능

$$B_1 = V_1 \Lambda_1 V_1^T, \Lambda_1 = diag(\lambda_1, \dots, \lambda_p), V_1 = [\boldsymbol{v}_1, \dots, \boldsymbol{v}_p]$$



따라서 최종적으로 구한 p차원 좌표 행렬  $\tilde{X}$ 는 다음과 같음

$$\tilde{X} = V_1 \Lambda_1^{\frac{1}{2}}$$

### MDS *Multidimensional Scaling*

③ 거리 정보를 보존하는 좌표 추출



구해진 행렬 B를 바탕으로  $p \times n$  행렬인  $\tilde{X}$ 를 구해보자!

0인 고윳값이 n-p개 있으므로, B는 다음과 같이 표현 가능

$$B_1 = V_1 \Lambda_1 V_1^T, \Lambda_1 = diag(\lambda_1, \dots, \lambda_p), V_1 = [\boldsymbol{v}_1, \dots, \boldsymbol{v}_p]$$

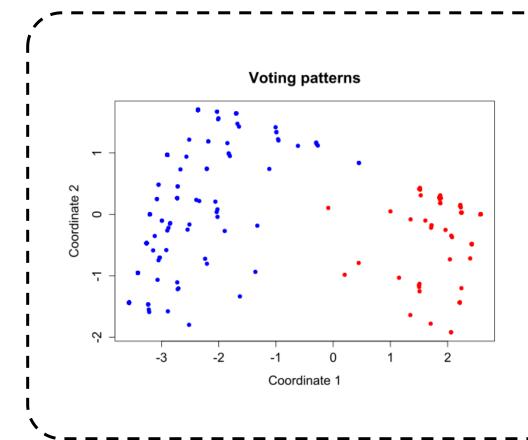


따라서 최종적으로 구한 p차원 좌표 행렬  $\tilde{X}$ 는 다음과 같음

$$\tilde{X} = V_1 \Lambda_1^{\frac{1}{2}}$$

## MDS *Multidimensional Scaling*

③ 거리 정보를 보존하는 좌표 추출

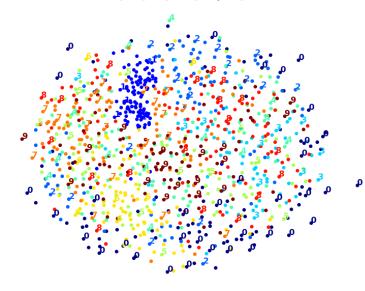


미국 하원 투표에 적용된 고전적인 MDS 예시

빨간색 점은 공화당, 파란색 점은 민주당 의원

## MDS *Multidimensional Scaling*

MNIST 데이터 시각화 with MDS



차이가 잘 보이지 않음 MDS는 **차원의 증가에 매우 취약** 

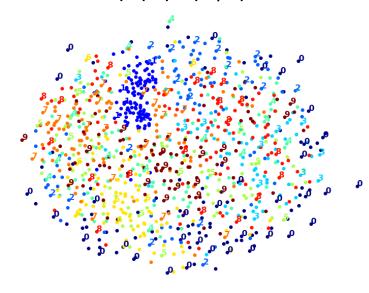
MDS는 모든 점들 간의 거리 정보 중요도가 같음

→ 가까운 점들 간의 거리 정보보다 멀리 떨어진 점들 간의 정보 영향력이 큼

무의미한 정보에 집중하는 경향이 있음

### MDS *Multidimensional Scaling*

MNIST 데이터 시각화 with MDS



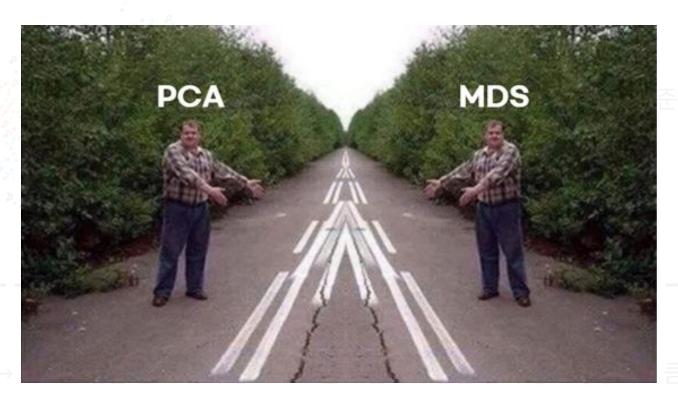
차이가 잘 보이지 않음 MDS는 **차원의 증가에 매우 취약** 

MDS는 모든 점들 간의 **거리 정보 중요도가 같음** 

→ 가까운 점들 간의 거리 정보보다 멀리 떨어진 점들 간의 정보 영향력이 큼 무의미한 정보에 집중하는 경향이 있음

# MDS *Multidimensional Scaling*

MNIST 데이터 시각화 with MDS



어쨌든 해냈다는 사실이 자랑스러운 PCA와 MDS

## MDS *Multidimensional Scaling*

MDS 기반 시각화의 특징

장점

고차원에서의 **객체 간 거리**를 저차원에서 <mark>보존</mark>하는 것에 집중 **비선형적**인 데이터에 대해서도 잘 작동하며 **다양한 거리 척도**를 사용 가능

단점

고차원에서는 무의미한 정보에 집중해 시각화가 잘 이루어지지 않음

→ 보다 저차원에서 직관적인 시각화 가능

### MDS *Multidimensional Scaling*

MDS 기반 시각화의 특징

장점

고차원에서의 **객체 간 거리**를 저차원에서 **보존**하는 것에 집중 **비선형적**인 데이터에 대해서도 잘 작동하며 **다양한 거리 척도**를 사용 가능

단점

고차원에서는 무의미한 정보에 집중해 시각화가 잘 이루어지지 않음

→ 보다 저차원에서 직관적인 시각화 가능

### **ISOMAP**

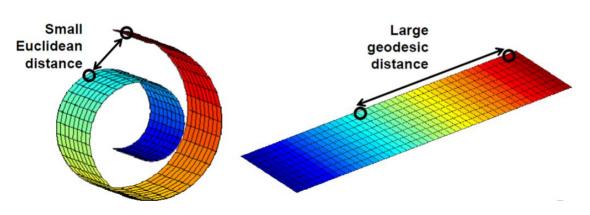
ISOMAP Isometric Mapping

MDS와 PCA를 결합한 기법

모든 점 사이의 측지 거리를 유지하는 임베딩을 추구

#### 측지 거리란?

두 측점 사이의 매니폴드를 따라 이루어진 거리



측점 사이의 유클리드 거리는 가깝지만 측지 거리는 멂

#### **ISOMAP**



### ISOMAP Isometric Map SOMAP의 목표

MDS와 PCA를 결합한 기법

<u>모든 점 사이의 **측지 거리를 우지**하는 임베딕은 추</u>근

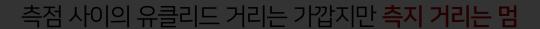
#### 내재된 매니폴드 구조를 <mark>형면으로 풀어내어 거리를 구하자!</mark> 두 측점 사이의 매니폴드를 따라 이루어진 거리

 $\rightarrow$  MDS에서 거리 행렬 D를 다르게 정의하여 해결

Euclidean distance



가보자고!



# ISOMAP 프로세스

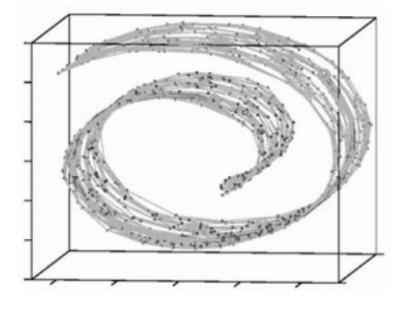
① 인접 그래프 구축

ε-ISOMAP

두 점 사이의 거리가 ε**보다 작은 점**을 이웃으로 선택

k-ISOMAP

k-nearest neighbor에 해당하는 점을 이웃으로 선택



### ISOMAP 프로세스

② 최단 경로 구하기



두 점들 간의 **최단 경로를 탐색**해 거리 행렬 D를 만듦 최단 경로 탐색을 위해 다익스트라 or 플로이드 이용

다익스트라 Dijkstra

한 점에서 모든 점 사이의

최단거리 탐색

시간복잡도: O(ElogV)

플로이드 Floyd

모든 점에서 모든 점 사이의

최단거리 탐색

시간복잡도:  $O(n^3)$ 

V: 노드, E: 간선

### ISOMAP 프로세스

② 최단 경로 구하기



두 점들 간의 **최단 경로를 탐색**해 거리 행렬 D를 만듦 최단 경로 탐색을 위해 다익스트라 or 플로이드 이용

#### 다익스트라 Dijkstra

한 점에서 모든 점 사이의

최단거리 탐색

시간복잡도: O(ElogV)

플로이드 Floyd

모든 점에서 모든 점 사이의

최단거리 탐색

시간복잡도:  $O(n^3)$ 

V: 노드, E: 간선

### ISOMAP 프로세스

② 최단 경로 구하기



두 점들 간의 **최단 경로를 탐색**해 거리 행렬 *D*를 만듦 최단 경로 탐색을 위해 다익스트라 or 플로이드 이용

다익스트라 Dijkstra

한 점에서 모든 점 사이의

최단거리 탐색

시간복잡도: O(ElogV)

플로이드 Floyd

모든 점에서 모든 점 사이의

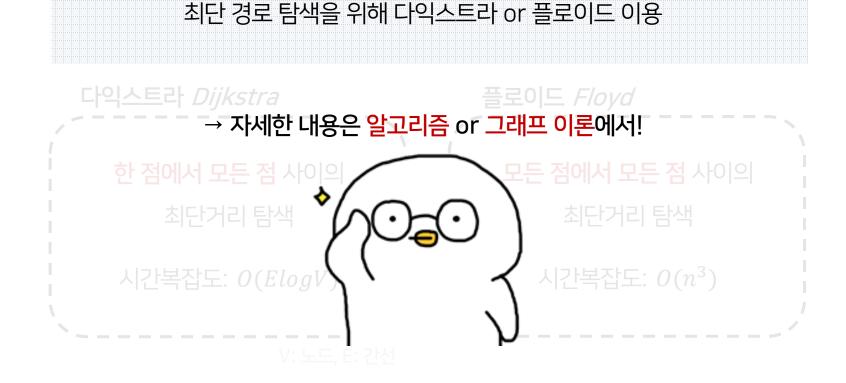
최단거리 탐색

시간복잡도:  $O(n^3)$ 

V: 노드, E: 간선

### ISOMAP 프로세스

② 최단 경로 구하기

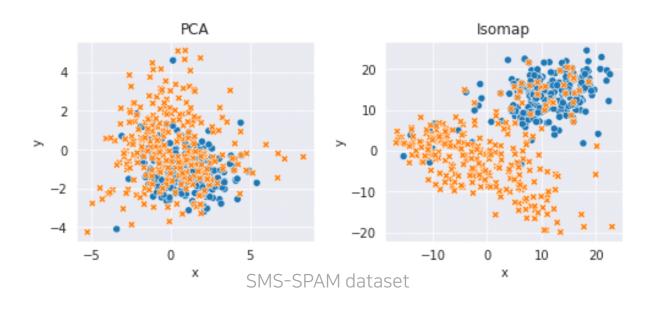


두 점들 간의 **최단 경로를 탐색**해 거리 행렬 D를 만듦

### ISOMAP 프로세스

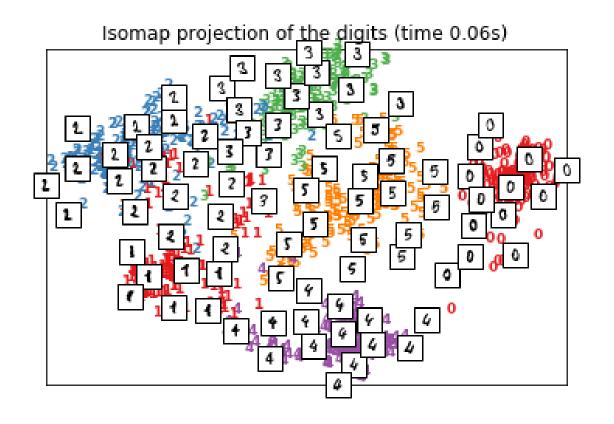
③ MDS 알고리즘 수행하기

재정의한 거리 행렬 D를 이용해 똑같이 MDS를 수행



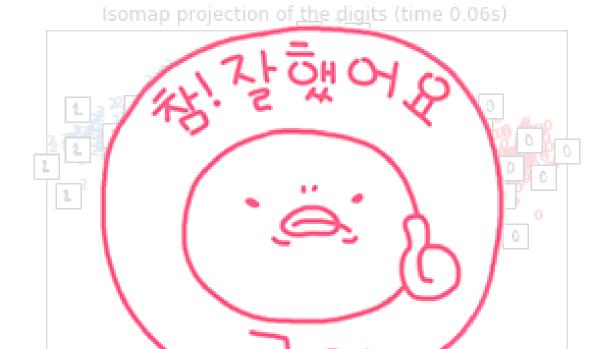
→ PCA에 비해 스팸을 매우 잘 분류

# MNIST 데이터 시각화 with ISOMAP



→ 7개의 숫자만 사용하긴 했지만 이전의 방법보다 훨씬 좋은 성능을 보임

### MNIST 데이터 시각화 with ISOMAP



→ 7개의 숫자만 사용하긴 했지만 이전의 방법보다 훨씬 좋은 성능을 보임

## ISOMAP의 특징

장점

비선형 차원 축소기법으로, 내재된 <mark>매니폴드</mark>를 잘 반영함 비선형의 경우에 좋은 성능을 보임

단점

매니폴드가 없다면, 성능이 좋지 않음 계산 비용이 매우 큼 (다익스트라, 플로이드)

→ 데이터의 형태를 잘 파악하고 사용해야 함





## ISOMAP의 특징

장점

비선형 차원 축소기법으로, 내재된 매니폴드를 잘 반영함 비선형의 경우에 좋은 성능을 보임

단점

매니폴드가 없다면, 성능이 좋지 않음 계산 비용이 매우 큼 (다익스트라, 플로이드)

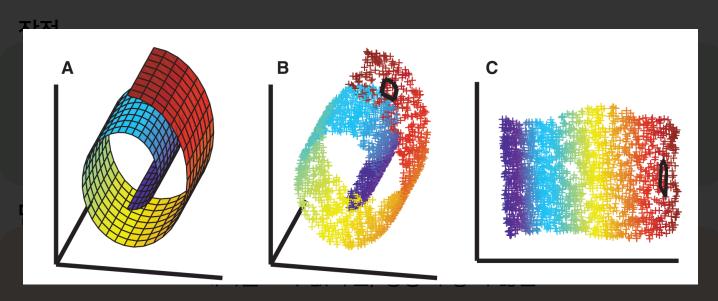
→ 데이터의 형태를 잘 파악하고 사용해야 함







# ISOMAP의 특<mark>웡</mark>마나 매니폴드를 잘 반영하나요?



계산 비용이 매우 큼 (다익스트라, 플로이드) 매니폴드를 <mark>완벽하게</mark> 펴낸 모습을 볼 수 있음

→ 데이터의 형태를 잘 파악하고 사용해야 함

LLE

LLE Locally Linear Embedding

고차원 공간에서 인접한 데이터들 간의 선형적 구조를 보존하며 임베딩

그럼 LLE는 선형 모델인가요?

LLE는 **좁은 범위**(Locally)에서 **선형 모델**(Linear)을 연결하며 매니폴드를 표현하는 알고리즘

→ 좁은 범위에서 선형적으로 표현했으므로 넓게 보면 **비선형 모델**!

LLE



LLE Locally Linear Embedding 목표

고차원 공간에서 인접한 데이터를 간의 선형적 구조를 모존하며 임베랑

#### 고차원에서 이웃관계인 점들

그럼 LLE는 선형 모델인가요? 저차원에서도 이웃관계여야 한다!

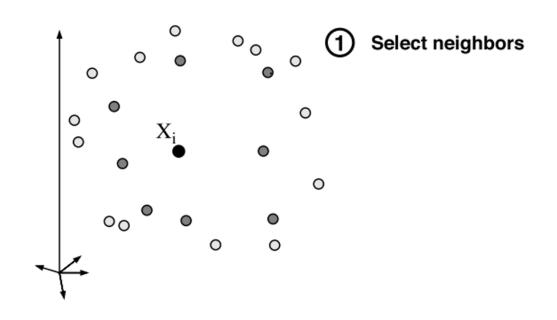
LLE는 <mark>좁은 범위</mark>(Locally)에서 **선형 모델**(Linear)을 연결하며 **하** 

매니폴드를 표현하는 알고리즘

→ 좁은 범위에서 선형적으로 표현했으므로 넓게 보면 비<mark>선형</mark>&

## LLE 프로세스

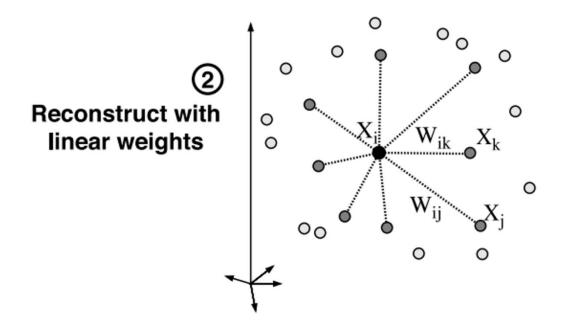
① 가장 가까운 이웃 탐색



각 데이터포인트에서 **k개의 이웃**을 선정

## LLE 프로세스

② 가중치 매트릭스 구성 (by local linearity)



 ${m k}$ 개의 이웃의 선형결합으로  $X_i$ 를 최대한 구현하는  $W_{ij}$ 를 학습

#### LLE 프로세스

② 가중치 매트릭스 구성 (by local linearity)

k개의 이웃의 선형결합으로  $X_i$ 를 최대한 구현하는  $W_{ij}$ 를 학습

$$E(W) = \sum_{i} \left| x_i - \sum_{j} W_{ij} x_j \right|^2$$

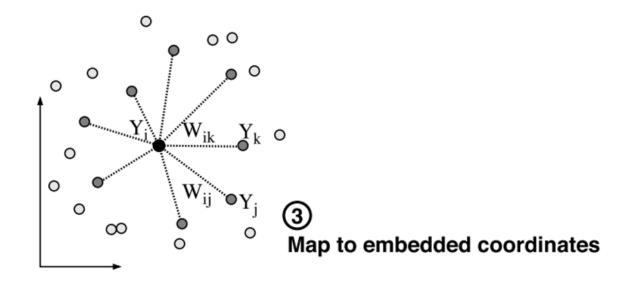
st.  $W_{ij} = 0$  does not belong to the neighbor of  $x_i$ 

$$\forall_i$$
 ,  $\sum_i W_{ij} = 1$ 

손실 함수로서 이웃 점들에게 둘러 쌓이면 0, 아니면 0보다 커져 총합을 줄이는 방향으로 최적화

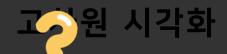
## LLE 프로세스

③ 부분 고윳값 분해



가중치를 보장하며 **차원을 축소**하여 최적의 임베딩을 찾음

$$\phi(W) = \sum_{i} \left| \mathbf{y}_{i} - \sum_{j} W_{ij} \mathbf{y}_{j} \right|^{2}$$



## LLE 프로세스 어떻게 SVD로 임베딩을 찾나요?



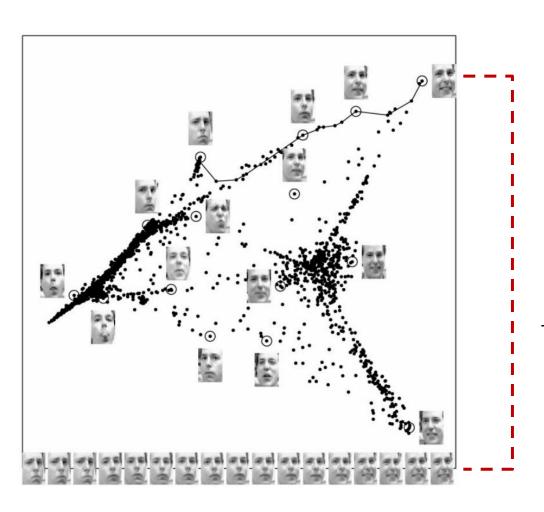
DXN NXN NXN NXd

Sparse Matrix  $M = (I - W)^T (I - W)$ 를 구한 후

가중치 SVD를 통해 하위 n개 고유벡터 추출을 찾음

By Rayleitz-Ritz Theorem (꽤나 심오하여 깊게 다루지 않음) 궁금하면 팀장에게 개인적으로 질문

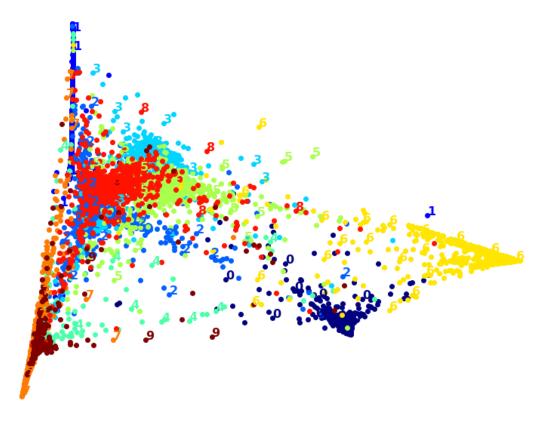
### LLE 프로세스



이웃 간의 구조를 보존하였더니 이웃들끼리 한 축으로 이동하며 학습

→ **연속적인 정보를 보존**할 수 있게 됨

## MNIST 데이터 시각화 with LLE



→ 저차원으로 매핑된 후 직선의 형태들로 임베딩됨



## LLE의 특징

장점

Local Minimum 문제가 **발생하지 않음** (SVD로 최적화) 비선형 임베딩 생성이 가능

단점

ISOMAP보다는 매니폴드를 잘 못잡아냄 국소 선형성 가정에 부합하지 않으면 성능 저하

→ 이웃관계 유지가 중요한 경우 사용을 추천함

## LLE의 특징

장점

Local Minimum 문제가 **발생하지 않음** (SVD로 최적화) 비선형 임베딩 생성이 가능

단점

ISOMAP보다는 매니폴드를 잘 못잡아냄 국소 선형성 가정에 부합하지 않으면 성능 저하

→ 이웃관계 유지가 중요한 경우 사용을 추천함

#### t-SNE

#### t-SNE와 LLE의 공통점

고차원에서 **이웃 관계**인 점은 저차원에서도 가깝게 임베딩하는 것이 목표

t-SNE 와 LLE의 차이점

t-SNE는 이웃보다 조금 **더 떨어져 있는 점들**까지 반영함

#### t-SNE

t-SNE와 LLE 목표의 공통점



t-SNE는 t분포를 사용한 SNE 기법이므로 SNE를 먼저 알아보자!

이웃보다 조금 **더 떨어져 있는 점들**까지 반영함

#### t-SNE

SNE Stochastic Neighbor Embedding

SNE는 고차원 공간에서의 유클리드 거리를 포인트들 간의 **유사성**을 표현하는 **조건부 확률**로 변환하는 방법

이 때, 확률적으로 이웃을 선택하며 그 확률은 다음과 같음

고차원에서 j번째 데이터를 이웃으로 선택할 확률

$$p_{i|j} = \frac{e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{\frac{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}$$

저차원에서 j번째 데이터를 이웃으로 선택할 확률

$$p_{i|j} = \frac{e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}}$$

#### t-SNE

SNE Stochastic Neighbor Embedding

SNE는 고차원 공간에서의 유클리드 거리를 포인트들 간의 **유사성**을 표현하는 **조건부 확률**로 변환하는 방법

이 때, 확률적으로 이웃을 선택하며 그 확률은 다음과 같음



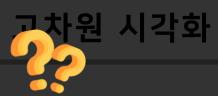
고차원에서 j번째 데이터를 이웃으로 선택할 확률

$$p_{i|j} = \frac{e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}{\sum_{k \neq i} e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma_i^2}}}$$



저차원에서 j번째 데이터를 이웃으로 선택할 확률

$$p_{i|j} = \frac{e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}}{\sum_{k \neq i} e^{-\|\mathbf{x}_i - \mathbf{x}_j\|^2}}$$

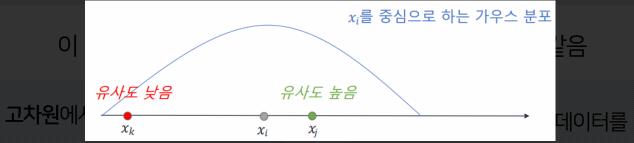


## t-SNE고차원에서 근처 점들을 이웃으로 선택하는 기준은?

SNE Stochastic Neighbor Embedding

Gaussian radius를 기준으로 사용하기 때문에

포인트들 간의 유더 강까이 있을수록 확률로 기환하는 방법



이웃으로 선택할 확률

이웃으로 선택할 확률

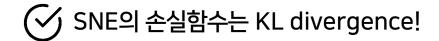
 $\sigma$ 가 포함된 하이퍼파라미터 perplexity는 아래와 같이 계산됨

$$p_{i|j} = \frac{e^{-2\sigma_i^2}}{\frac{\|P_i erplexity(P_i) - 2^{\sum_i p_{j|i} \log_2 p_{j|i}} e^{-\|X_i - X_j\|^2}}{\sum_{k \neq i} e^{-\|X_i - X_j\|^2}}}$$

 $\sum_{k\neq i} e^{-zo_i}$ 

t-SNE

SNE



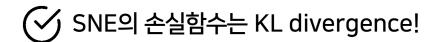
KL divergence Kullback-Leibler divergence

KL divergence는 **두 분포의 유사성을 파악**하게 해주며, 스칼라 값으로 나오기에 유사도 행렬로 사용할 수는 없음

$$Cost = \sum_{i} KL(P_i||Q_i) = \sum_{i} \sum_{j} p_{j|i} log \frac{p_{j|i}}{q_{j|i}}$$

t-SNE

SNE



KL divergence Kullback-Leibler divergence

KL divergence는 **두 분포의 유사성을 파악**하게 해주며, 스칼라 값으로 나오기에 유사도 행렬로 사용할 수는 없음

$$Cost = \sum_{i} KL(P_i||Q_i) = \sum_{i} \sum_{j} p_{j|i} log \frac{p_{j|i}}{q_{j|i}}$$

두 분포 p와 q가 유사할수록 log의 진수 부분이 작아짐



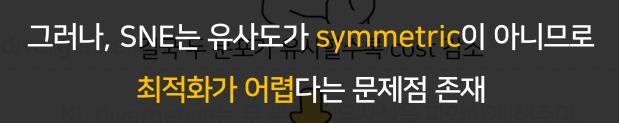


KL divergence 결국 두 분포가 유사할수록 cost 감소



고차원에서 이웃으로 뽑을 확률과 저차원에서 이웃으로 뽑을 확률을

최대한 <mark>일치</mark>시키는 것이 목표



고차원에서 어웃으로 뽑을 확률과 저차원에서 이웃으로 뽑을 확률을 최대한 일치시키는 것이 목표 ★ t-SNE 의 필요성!

#### t-SNE

t-SNE의 필요성1: 최적화의 어려움

symmetric SNE

일반적으로 상호동일하지 않은 조건부 확률을 점들간 상호동일하게 만들어주기 위해 조정함으로써 다음과 같이 cost function을 간략화

$$Cost = \sum_{i} KL(P_i||Q_i) = \sum_{i} \sum_{j} p_{ij} log \frac{p_{ij}}{q_{ij}}$$

#### t-SNE

t-SNE의 필요성1: 최적화의 어려움

symmetric SNE

일반적으로 상호동일하지 않은 조건부 확률을 점들간 상호동일하게 만들어주기 위해 조정함으로써 다음과 같이 cost function을 간략화



$$p_{ij} = \frac{p_{j|i} + p_{i|j}}{2n}$$

$$Cost = \sum_{i} KL(P_i||Q_i) = \sum_{i} \sum_{j} p_{ij} log \frac{p_{ij}}{q_{ij}}$$



♣ 유사도를 symmetric으로 만들어주어 최적화가 쉬운 형태!

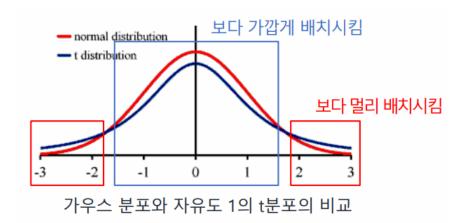
#### t-SNE

t-SNE의 필요성2: Crowding problem



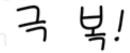
SNE의 Crowding problem

SNE에서 가정하는 가우시안 분포는 양쪽 꼬리가 충분히 두텁지 않음 따라서 중간정도 가까운 거리와 **비교적 먼 거리를 구별하지 못하는 문제** 



t-SNE

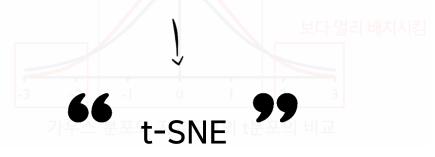
t-SNE의 필요성2: Crowding problem





이를 완화하기 위해 가우시안 분포와 유사하지만

꼬리가 더 두터운 자유도가 1인 t분포 사용



#### t-SNE

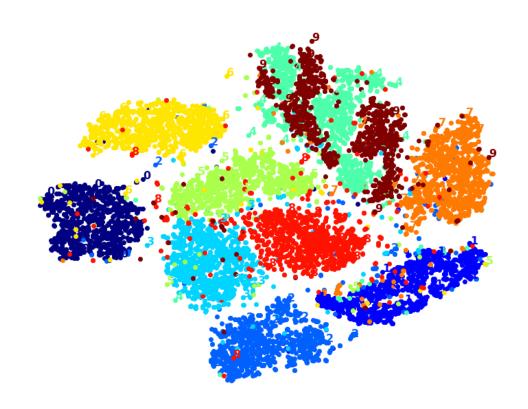
t-SNE t-Stochastic Neighbor Embedding

SNE의  $p_{ij}$ 는 동일하게 사용하되,  $q_{ij}$ 에만 t분포를 적용한 방법으로,  $t분포를 사용한 임베딩 공간의 두 점 사이 유사도 <math>q_{ij}$ 는 다음과 같이 정의됨

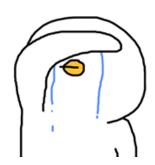
$$q_{ij} = \frac{\left(1 + |y_i - y_j|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + |y_i - y_j|^2\right)^{-1}}$$

t-SNE는  $\theta$ 공간에서 가까운 좌표 값으로 만든 q와 p가 다르다면 gradient descent를 통해 더 비슷해지는 방향으로 y의 점들을 이동

## MNIST 데이터 시각화 with t-SNE



이거 하려고 지금까지 고생했습니다..



클러스터가 잘 생성되어 있을 뿐만 아니라, 유사한 숫자끼리 인접!

#### t-SNE

t-SNE 기반 시각화 특징

장점

PCA와 달리 군집이 중복되지 않아 시각화에 유용

전역적 특징도 잘 잡아내 클러스터링 구조를 잘 보존

단점

**계산 시간**이 매우 많이 걸리며, **랜덤 프로세스**에 의존적

→보통 50차원 이내로 압축 뒤 t-SNE

#### t-SNE

t-SNE 기반 시각화 특징

장점

PCA와 달리 군집이 중복되지 않아 시각화에 유용 전역적 특징도 잘 잡아내 클러스터링 구조를 잘 보존

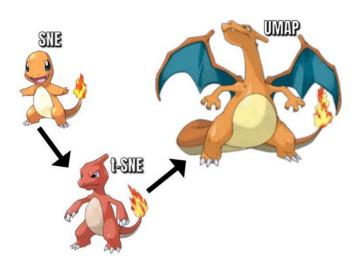
단점

계산 시간이 매우 많이 걸리며, 랜덤 프로세스에 의존적

## UMAP

UMAP 소개

## 고차원 시각화 기법 중 가장 <mark>최신</mark>에 개발된 알고리즘 대규모 데이터셋에 대해 효과적



SNE, t-SNE, UMAP의 관계

**UMAP** 

UMAP 소개

고차원 시각화 기법 중 가장 <mark>최신</mark>에 개발된 알고리즘 대규모 데이터셋에 대해 효과적

> 클린업의 범위를 <mark>한참</mark> 웃도는 알고리즘 본 클린업에서는 간략하게 소개

> > SNE, t-SNE, UMAP의 관계

#### **UMAP**

UMAP 알고리즘의 직관적 이해





Fuzzy simplical complex

고차원 그래프를 만드는 단계에서 고안된 개념으로,

두 점이 연결될 가능도를 나타내는 가중치가 있는 그래프 표현

수학적 이론으로 들어가면 상당히 어려우니 직관적인 이해만 해보자!



#### UMAP

UMAP 알고리즘의 직관적 이해



두 포인트의 연결 여부를 판단하기 위해 각 포인트를 중심으로 원을 그리고, 원이 겹쳐지는 점들에 대해 연결되어 있다고 생각



원의 반지름이 커짐에 따라 그래프를 fuzzy하게 만들어 점 사이의 연결 정도를 줄임

#### **UMAP**

UMAP 알고리즘의 직관적 이해



두 포인트의 연결 여부를 판단하기 위해 각 포인트를 중심으로 원을 그리고, 원이 겹쳐지는 점들에 대해 연결되어 있다고 생각

r----원이 너무 **작으면** 크기가 작고 **고립된 클러스터**가 생성 원의 크기 -----원이 너무 크면 모든 점들이 **다 연결**되도록 만듬



N-th nearest neighbor로 해결!

#### **UMAP**

UMAP 알고리즘의 직관적 이해



두 포인트의 연결 여부를 판단하기 위해 각 포인트를 중심으로 원을 그리고, 원이 겹쳐지는 점들에 대해 연결되어 있다고 생각

? N-th nearest neighbor

n번째로 가까운 점까지 반지름을 확장하는 방식 데이터가 <mark>충분</mark>한 지역은 <mark>반지름을 작게</mark>, 희소한 지역은 <mark>반지름을 크게</mark> 설정 가능

#### UMAP

UMAP 알고리즘의 직관적 이해



두 포인트의 연결 여부를 판단하기 위해 각 포인트를 중심으로 원을 그리고, 원이 겹쳐지는 점들에 대해 연결되어 있다고 생각



원의 반지름이 커짐에 따라 그래프를 fuzzy하게 만들어 점 사이의 연결 정도를 줄임



## UMAP

UMAP 알고리즘의 직관적 이해



각 포인트들이 **가장 가까운 점들과는 연결**되어야 한다고 규정함으로써 local strucutre 보존, global structure와 밸런스를 맞추도록 함



이렇게 생성된 **고차원 그래프**를 가지고, 저차원 그래프의 레이아웃을 최대한 비슷하게 최적화하여 매핑



## **UMAP**

UMAP 알고리즘의 직관적 이해



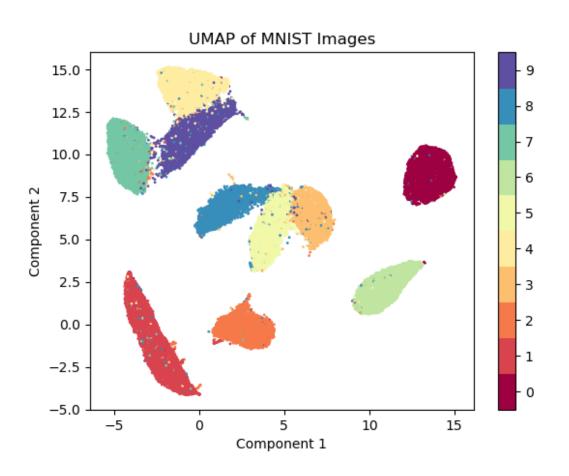
각 포인트들이 **가장 가까운 점들과는 연결**되어야 한다고 규정함으로써 local strucutre 보존, global structure와 밸런스를 맞추도록 함



이렇게 생성된 **고차원 그래프**를 가지고, 저차원 그래프의 레이아웃을 최대한 비슷하게 최적화하여 매핑



## MNIST 데이터 시각화 with UMAP



#### MNIST 데이터 시각화 with UMAP



고차원에서의 **이웃 관계**가 잘 **유지**됨

## UMAP 사용 가이드

하이퍼파라미터



n\_neighbors

초기 고차원 그래프 생성 시에 사용되는 nearest neighbor의 수

```
크면 global structure에 집중
--- 작으면 local structure에 집중
```

## UMAP 사용 가이드

하이퍼파라미터



n\_neighbors

초기 고차원 그래프 생성 시에 사용되는 nearest neighbor의 수

```
·--- 크면 global structure에 집중
--- 작으면 local structure에 집중
```



저차원 공간에서 포인트들 간의 최소 거리 UMAP이 얼마나 점들을 **촘촘하게** 묶을 지 조절

```
크면 포인트들이 상대적으로 느슨하게 퍼지게 됨
작으면 포인트들이 촘촘하게 무리지어 있게 됨
```

#### **UMAP**

장점

t-SNE보다 **더 빠른 속도** 퍼지 위상학과 리만 기하학에 기반한 **강력한 수학적 기반** 지역적(local) 정보를 넘어 **전역적(global) 정보**까지 잘 보존

단점

**하이퍼파라미터**의 영향을 크게 받음 정보손실에 의한 **데이터 왜곡** 

→ **다양한 하이퍼 파라미터 조합**을 시도해봐야 함

#### **UMAP**

장점

t-SNE보다 **더 빠른 속도** 퍼지 위상학과 리만 기하학에 기반한 **강력한 수학적 기반** 지역적(local) 정보를 넘어 **전역적(global) 정보**까지 잘 보존

단점

**하이퍼파라미터**의 영향을 크게 받음 정보손실에 의한 **데이터 왜곡** 

→ **다양한 하이퍼 파라미터 조합**을 시도해봐야 함

## **UMAP**

UMAP 시연



UMAP 시연에 관심 있는 분들은 다음 링크를 들어가보세요

https://pair-code.github.io/understanding-umap/

## 고차원 시각화 총정리

변수들이 전반적으로 정규분포 형태를 띄고, 해석력을 조금이나마 챙기고 싶다. 비교적 저차원이다. 데이터들이 선형일 것이다. 최대한 빨라야 한다.

데이터 포인트 간의 거리 정보만이 중요하다. 비교적 저차원이다.





## 고차원 시각화 총정리

데이터셋에 매니폴드가 내재하는 것 같다. 단순 유클리드 거리가 아닌, 실제 매니폴드 상의 거리 정보만을 반영하여 시각화하고 싶다.



데이터셋에 매니폴드가 내재하는 것 같다. 데이터 포인트 간의 **이웃 관계를 저차원에서도 유지**하고 싶다.

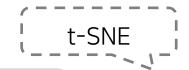




## 고차원 시각화 총정리

고차원 데이터 시각화가 하고싶고 시간은 많다.

군집별로 구별이 잘 되어있으면 좋겠다.



용도에 맞는 전문가용 모델이 필요하다.

군집별로 구별이 잘 되어있으면 좋겠고, 시간이 적게 들었으면 좋겠다.



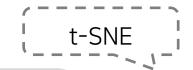


#### 선대팀 클린업

## 고차원 시각화 총정리

고차원 데이터 시각화가 하고싶고 시간은 많다.

군집별로 구별이 잘 되어있으면 좋겠다.



용도에 맞는 전문가용 모델이 필요하다.

군집별로 구별이 잘 되어있으면 좋겠고, 시간이 적게 들었으면 좋겠다.





# THANK YOU