

# CH06 차원 축소

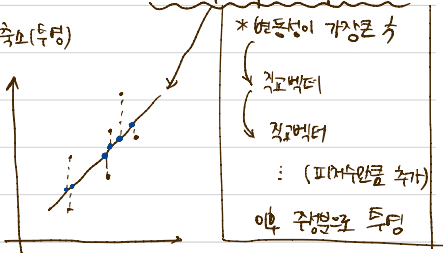
## 01. 차원 축소 dimension reduction

- 많은 피쳐 다차원 → 축소 → 저차원 구조
  - sparse : data간 거리 ↑
  - 피쳐간 상관관계 ↑ : 다중공선성 문제
- 차원 축소 = 피쳐 수 ↓
- 정보 손실 (격화) : 해석성 ↑

- 피쳐 선택 Selection : 제거
  - 한번에 고려, Latent factor
- 피쳐 추출 extraction : 양적 → 완전히 새로 고
- 이미지, 텍스트 처리에 Good.
- PCA, SVD, NMF
  - Semantic Topic modeling

## 02. PCA : Principal Component Analysis

- 피쳐간 상관관계 → 주성분 principal components (축)
  - 축소 (특성)



(변형 매수) 주성분 = 고유벡터

$$C = P \Sigma P^T$$

임베딩 → 고유벡터 + 고유값

이진 입력 data를 실형 변환

1. input → 공분산 행렬
2. 고유벡터, 고유값 계산
3. ↓ 큰 순서
4. input (변환)

- PCA. 여러 속성 계산,  $\Delta$  > 0 필요

-  $PCA(n\_component = )$   
→ 변환할 피쳐 수

↳ fit + transform, fit\_transform (input\_data)

- PCA.explained\_variance\_ratio\_

: 전체 변동성에서 결관된 변동성 비율

n 개  
↓  
2 개  
2 개 피쳐 비율

## 03 LDA : Linear Discriminant Analysis

: 선형 판별 분석법

- 분류 목적 → 분기 기준 최대화 유지

PCA : 최대 변동성 가지는 축

LDA : 최대한 분리의 하계하는 축

↑ 클래스간 분산 & 클래스내 분산

- 분산행렬 → 고유 vector

PCA : 공분산행렬 → 고유 vector

- sklearn.discriminant\_analysis.LinearDiscriminantAnalysis

## 04 SVD : Singular Value Decomposition

- PCA :  $n \times n$  행렬 → 고유 vector

SVD :  $m \times n$  → 고유 vector

$$A = U \Sigma V^T$$

$m \times n \quad m \times m \quad n \times n$     U, V는 등-1 벡터, 서로 직교  
Σ는 대각행렬.  
A의 특이값 = Σ의 대각원소

- Truncated SVD : Σ의 대각원소 중 상위 일부 데이터 사용  
음분해. Δ 내적은 2 대용 불필요

- scipy.sparse.linalg.svds

scipy.linalg.svd : 최소행렬만 지원

- 비선형, 이미지 압축에 사용

## 05 NMF: Non-negative Matrix Factorization

- $\approx$  SVD. Low-Rank Approximation, 항렬 분해
- 원본 항렬의 모든 원소  $> 0$  이면, 간단한 항렬 분해 가능
- 이미지 압축, 추천시스템, 추천에 사용

## 06

- PCA의 개념을 고려 데이터를 잘 설명할 수 있는

잠재적인 요소를 추출하는 것

- PCA: 변동성  $\rightarrow$  축  $\rightarrow$  투영

고유벡터  $\rightarrow$  선형변환

- LDA: 불가분성  $\rightarrow$  축

- SVD/NMF: 고차원성질  $\rightarrow$  항렬분해  $\rightarrow$  저차원성질  
항렬화 가능  $\leftarrow$