

# SVD(Singular Value Decomposition), PCA(Principle Components Analysis), and LSA(Latent Semantic Indexing)

Hyopil Shin(Seoul National University)

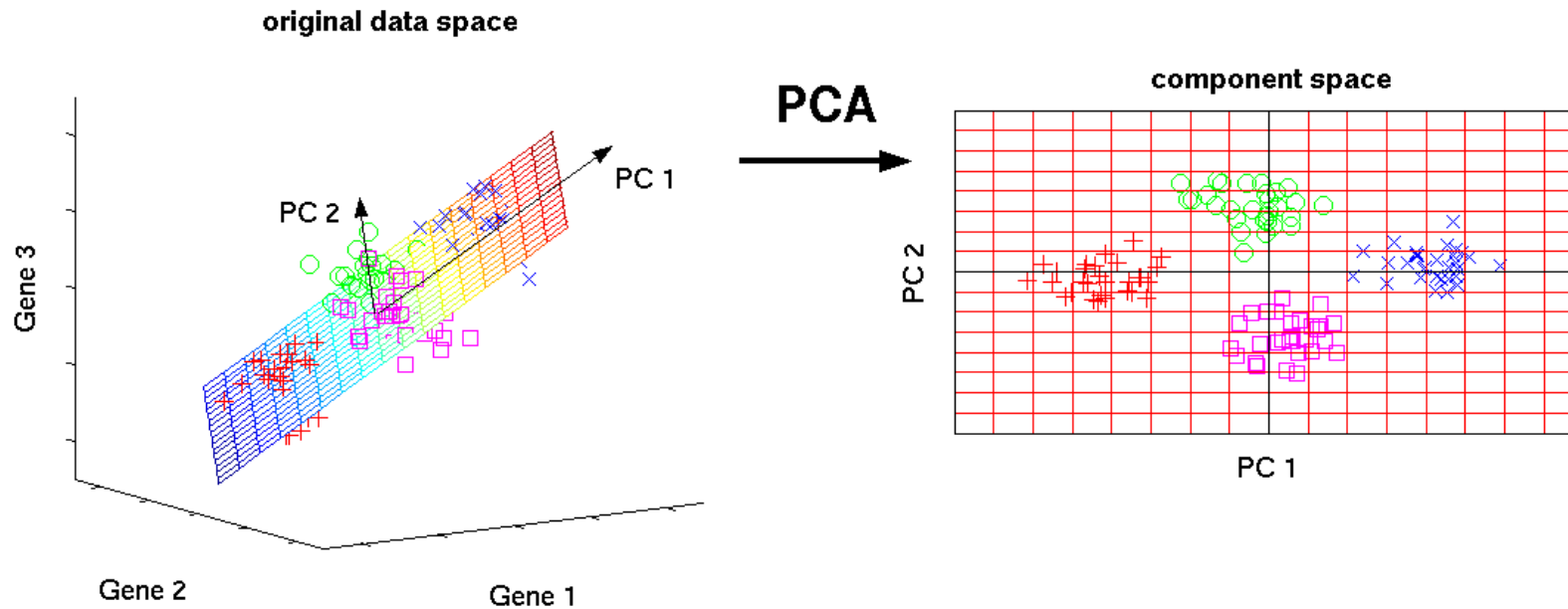
Computational Linguistics

#Semantics with Dense Vectors

Supplement Data

# PCA(Principle Components Analysis)

- 데이터의 분산(variance)을 최대한 보존하면서 서로 직교하는 새 축을 찾아 고차원 공간의 표본들이 선형 연관성이 없는 저차원 공간으로 변환하는 기법



# PCA(Principle Components Analysis)

- 공분산 행렬

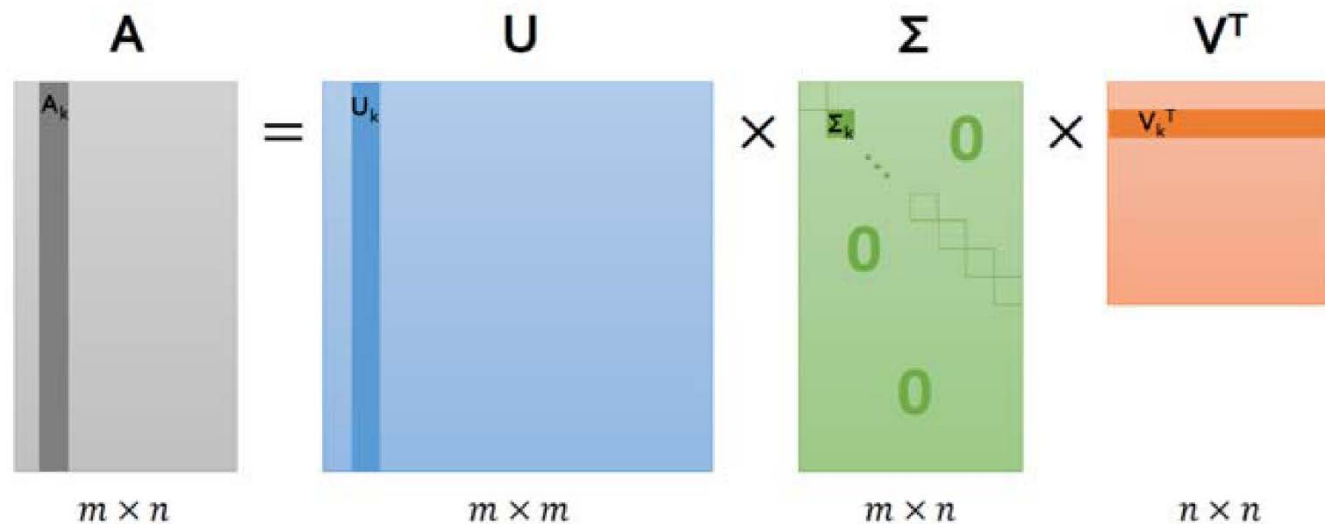
$$\text{cov}(A) = \frac{1}{n-1} A^T A \propto A^T A$$

- 고유분해 (eigen decomposition)
- $A^T A = V \Lambda V^T$

# 특이값 분해(singular value decomposition)

- $m \times n$  크기의 데이터 행렬  $A$ 를 다음과 같이 분해

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



# 특이값 분해(singular value decomposition)

- 특이벡터 – 행렬  $U$ 와  $V$ 에 속한 열벡터
- 모든 특이벡터는 서로 직교

$$U = [\vec{u_1} \quad \vec{u_2} \quad \dots \quad \vec{u_m}]$$

$$V = [\vec{v_1} \quad \vec{v_2} \quad \dots \quad \vec{v_n}]$$

$$\vec{u_k} = \begin{bmatrix} u_{k1} \\ u_{k2} \\ \dots \\ u_{km} \end{bmatrix} \quad \vec{v_k} = \begin{bmatrix} v_{k1} \\ v_{k2} \\ \dots \\ v_{kn} \end{bmatrix}$$

$$\vec{u_k}^T \vec{u_k} = 1, \quad U^T U = I$$

$$\vec{v_k}^T \vec{v_k} = 1, \quad V^T V = I$$

# 특이값 분해(singular value decomposition)

- 행렬  $\Sigma$ 의 특이값은 모두 0보다 크거나 같으며 내림차순으로 정렬
- 행렬  $\Sigma$ 의  $k$ 번째 대각원소 해당하는  $\Sigma_k$ 는 행렬  $A$ 의  $k$ 번째 고유값에 제곱근을 취한 값과 같음
- 특이값 분해를 주성분 분석과 비교하기 위해 행렬  $A$ 를 제공
  - 대각행렬의 거듭제곱은 대각원소들만 거듭제곱 해 준 결과와 동일

$$\Sigma_k = \sqrt{\lambda_k}$$

$$\begin{aligned} A^T A &= (U \Sigma V^T)^T U \Sigma V^T \\ &= V \Sigma U^T U \Sigma V^T \\ &= V \Sigma^2 V^T \\ &= V \Lambda V^T \end{aligned}$$

# 특이값 분해의 종류

- thin SVD
- compact SVD
- truncated SVD

thin SVD

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

compact SVD

$$A = U_r \Sigma_r V_r^T$$

truncated SVD

$$A' = U_t \Sigma_t V_t^T$$

# 특이값 분해 예시

- SVD 예

$$A = U\Sigma V^T$$
$$\begin{bmatrix} 2 & 3 \\ 1 & 4 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.82 & -0.58 & 0 & 0 \\ 0.58 & 0.82 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 5.47 & 0 & 0 & 0 \\ 0 & 0.37 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0.40 & 0.91 \\ -0.91 & 0.40 \end{bmatrix}$$

$$A' = U_1 \Sigma_1 V_1^T$$

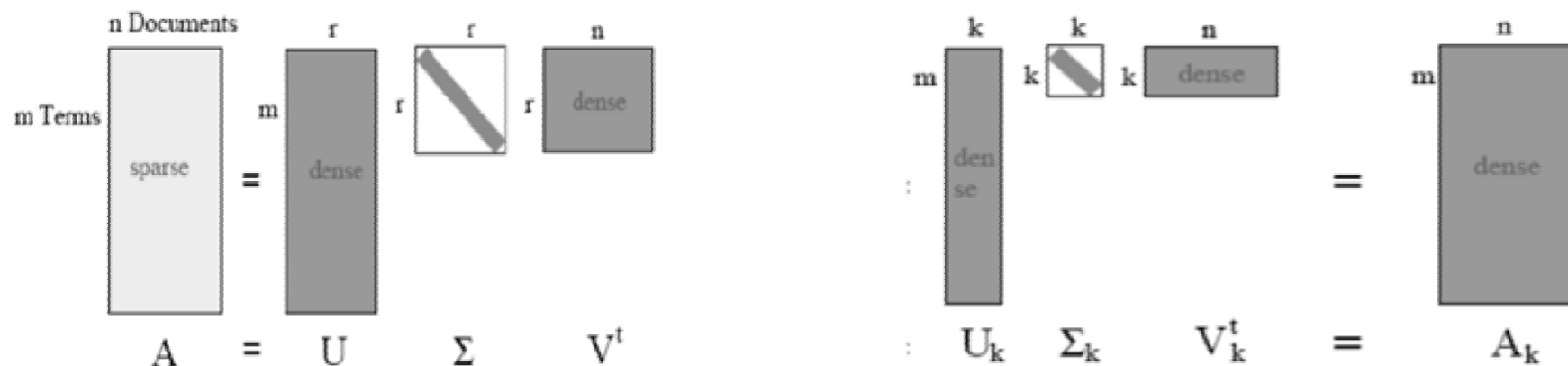
- Truncated SVD 예

$$\begin{bmatrix} 1.79 & 4.08 \\ 1.27 & 2.89 \\ 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0.82 \\ 0.58 \\ 0 \\ 0 \end{bmatrix} [5.47] \begin{bmatrix} 0.40 & 0.91 \end{bmatrix}$$



# Latent Semantic Analysis

- 단어-문서행렬(word-document matrix), 단어-문맥행렬(window-based co-occurrence matrix) 등 입력 데이터에 특이값 분해를 수행해 데이터의 차원수를 줄여 계산 효율성을 키우는 한편 행간에 숨어있는 의미를 이끌어 내기 위한 방법



# Latent Semantic Analysis

$$A_k = U_k \Sigma_k V_k^T$$

$$U_k^T A_k = U_k^T U_k \Sigma_k V_k^T = I \Sigma_k V_k^T = \Sigma_k V_k^T = X_1$$

$$A_k V_k = U_k \Sigma_k V_k^T V_k = U_k \Sigma_k^T I = U_k V_k^T = X_2$$

# Latent Semantic Analysis 예시

doc1 : 나,는,학교,에,가,ㄴ,다

doc2 : 학교,에,가,는,영희

doc3 : 나,는,영희,는,좋,다

-	doc1	doc2	doc3
나	1	0	0
는	1	1	2
학교	1	1	0
에	1	1	0
가	1	1	0
ㄴ	1	0	0
다	1	0	1
영희	0	1	1
좋	0	0	1

# Latent Semantic Analysis 예시

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

$$\begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 2 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} -0.17 & 0.27 & -0.40 \\ -0.63 & -0.41 & -0.03 \\ -0.32 & 0.37 & 0.21 \\ -0.32 & 0.37 & 0.21 \\ -0.32 & 0.37 & 0.21 \\ -0.17 & 0.27 & -0.40 \\ -0.33 & -0.12 & -0.52 \\ -0.30 & -0.29 & 0.49 \\ -0.15 & -0.39 & -0.13 \end{bmatrix} \begin{bmatrix} 3.61 & 0 & 0 \\ 0 & 2.04 & 0 \\ 0 & 0 & 1.34 \end{bmatrix} \begin{bmatrix} -0.63 & -0.53 & -0.57 \\ 0.56 & 0.20 & -0.80 \\ -0.54 & 0.83 & -0.17 \end{bmatrix}$$

$$A' = U_2\Sigma_2V_2^T$$

$$\begin{bmatrix} 0.71 & 0.44 & -0.09 \\ 0.97 & 1.04 & 1.99 \\ 1.15 & 0.76 & 0.04 \\ 1.15 & 0.76 & 0.04 \\ 1.15 & 0.76 & 0.04 \\ 0.71 & 0.45 & -0.09 \\ 0.62 & 0.58 & 0.88 \\ 0.36 & 0.45 & 1.11 \\ -0.09 & 0.14 & 0.97 \end{bmatrix} = \begin{bmatrix} -0.17 & 0.27 \\ -0.63 & -0.41 \\ -0.32 & 0.37 \\ -0.32 & 0.37 \\ -0.32 & 0.37 \\ -0.17 & 0.27 \\ -0.33 & -0.12 \\ -0.30 & -0.29 \\ -0.15 & -0.39 \end{bmatrix} \begin{bmatrix} 3.61 & 0 \\ 0 & 2.04 \end{bmatrix} \begin{bmatrix} -0.63 & -0.53 & -0.57 \\ 0.56 & 0.20 & -0.80 \end{bmatrix}$$

# Latent Semantic Analysis 예시

---

$$X_1 = \begin{bmatrix} -2.28 & -1.90 & -2.07 \\ 1.14 & 0.42 & -1.64 \end{bmatrix}$$

$$\text{round}(A') = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 2 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix}, \quad A = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 2 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 1 \end{bmatrix}$$

$$X_2 = \begin{bmatrix} -0.63 & 0.56 \\ -2.30 & -0.84 \\ -1.16 & 0.76 \\ -1.16 & 0.76 \\ -1.16 & 0.76 \\ -0.63 & 0.56 \\ -1.20 & -0.24 \\ -1.10 & -0.60 \\ -0.57 & -0.80 \end{bmatrix}$$