

# Latent Semantic Indexing (LSI)

## Motivating Example

- In search application, we use a word-document vector for efficient indexing and retrieval.

	$w_1$	$w_2$	$w_3$	---	$w_m$	$\rightarrow$
$d_1$	2.9					
$d_2$		3.4				
:						
$d_n$						

A word-document vector is a  $n \times m$  vector. Every cell indicates "importance" of the word in the document.

→ TF-IDF

(Term Frequency

- Invers Document  
frequency )

$D_1 = \text{Dog says Bow}$

$D_2 = \text{Cat says Mew}$

$D_3 = \text{Bow Bow Mew Mew}$

	Dog	says	Bow	Cat	Mew
$D_1$	1	1	1	0	0
$D_2$	0	1	0	1	1
$D_3$	0	0	2	0	2

Query : "Cat Mew"

Query Vector: [0 0 0 1 1]

Use similarity measures: Cosine?

$\cos(\vec{q}_1, \vec{d}_1)$

$\cos(\vec{q}_1, \vec{d}_2)$

$\cos(\vec{q}_1, \vec{d}_3)$

Rank 3

Rank 1

Rank 2

We want to reduce dimensions of word-document vector

- Number of documents & words are very high.
- Without reducing dimensions, most of the feature (dimension, word count) for a document will be zero. (Sparse)

① Reduce dimension to avoid sparsity

② We want to relate word-word & doc-doc.

		Financial Inst.		
		<u><math>w_1 = \text{Bank}</math></u>	<u><math>w_2 = \text{Money}</math></u>	<u><math>w_3 = \text{River}</math></u>
$d_1$		1	1	0
	$d_2$	1	0	1
$D_1 = \text{Bank money}$				
$D_2 = \text{Bank River}$				

	good	nice	movie
D <sub>1</sub>	1	0	1
D <sub>2</sub>	0	1	1

D<sub>1</sub> = good movie

D<sub>2</sub> = nice movie

q = "nice movie"

## Latent Semantic Indexing

An m × n dimension matrix  
can be divided it into 3 matrices

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

# Singular Value Decomposition (SVD)

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

$C$  = word document vector  
( $m \times n$ )

$m$  = no. of docs

$n$  = no. of words

$U$  =  $m \times r$  matrix ( $\overset{(doc)}{doc}$  matrix)

$V^T$  =  $r \times m$  matrix ( $\overset{(word)}{word}$  matrix)

$\Sigma$  = Diagonal matrix  
( $\sqrt{\text{Eigen values}}$ ) (Singular Values)

(Ref: Diagonalization)

$$\boxed{P^{-1} S A S^{-1} = D}$$

Eigen Values

## SVD Example

$$C = \begin{bmatrix} 5 & 5 \\ -1 & 7 \end{bmatrix}$$

Find  $U, \Sigma, V^T$

$$\underline{C = U \Sigma V^T}$$

$$C^T C = (V \Sigma^T U^T) \cdot (U \Sigma V^T)$$

$$\boxed{C^T C = V \Sigma^T \Sigma V^T \quad \textcircled{P}}$$

$$C V = U \Sigma \quad \textcircled{2}$$

① Find  $C^T C$

$$C^T C = \begin{bmatrix} 5 & -1 \\ 5 & 7 \end{bmatrix} \begin{bmatrix} 5 & 5 \\ -1 & 7 \end{bmatrix}$$
$$= \begin{bmatrix} 26 & 18 \\ 18 & 74 \end{bmatrix}$$

② Find Eigen values & vectors  
of  $C^T C$ ,

$$\lambda_1 = 20, \quad \lambda_2 = 80$$

$$V_1 = \begin{bmatrix} -3/\sqrt{10} \\ 1/\sqrt{10} \end{bmatrix}$$

---

$$V_2 = \begin{bmatrix} 1/\sqrt{10} \\ 3/\sqrt{10} \end{bmatrix}$$

③ Define  $V$  &  $\Sigma$

$\Sigma$  = square root of Eigen  
value, diagonal matrix

$$\Sigma = \begin{bmatrix} 2\sqrt{5} & 0 \\ 0 & 4\sqrt{5} \end{bmatrix}$$

$$V = \begin{bmatrix} -3/\sqrt{10} & 1/\sqrt{10} \\ 1/\sqrt{10} & 3/\sqrt{10} \end{bmatrix}$$

How to find U?

$$CC^T = (U\Sigma V^T)(V\Sigma^T U^T)$$

$$= U\Sigma\Sigma^T U^T$$

$$CC^T = \begin{bmatrix} 5 & 5 \\ -1 & 7 \end{bmatrix} \begin{bmatrix} 5 & -1 \\ 5 & 7 \end{bmatrix}$$

$$= \begin{bmatrix} 50 & 30 \\ 30 & 50 \end{bmatrix}$$

$$\lambda_1 = 20, \quad \lambda_2 = 80$$

$$U_1 = \begin{bmatrix} 1/\sqrt{2} \\ -1/\sqrt{2} \end{bmatrix}$$

$$U_2 = \begin{bmatrix} 1/\sqrt{2} \\ 1/\sqrt{2} \end{bmatrix}$$

$$U = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}$$

$$\begin{bmatrix} 5 & 5 \\ -1 & 7 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \circ \begin{bmatrix} 2\sqrt{5} & 0 \\ 0 & 4\sqrt{5} \end{bmatrix}$$

$$\circ \begin{bmatrix} -\sqrt{3}/10 & \sqrt{1}/10 \\ \sqrt{1}/10 & \sqrt{3}/10 \end{bmatrix}$$

# Word document matrix A

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

Annotations:

- $U$ : doc-doc matrix  $m \times r$
- $\Sigma$ : Singular Value ( $\sqrt{\text{Eigen values}}$ )
- $V^T$ : word-word matrix  $r \times n$

"Ignore low singular Values"  
First  $r$  singular values non zero