

# Latent Semantic Indexing

1

## Summary: Singular Value Decomposition (SVD)

$M$  - number of terms                       $N$  - number of documents

$C$  the term×document **matrix of weights** (our old  $w_{ij}$ )

– of **rank  $r$**  ( $r \leq \min(M,N)$ )

$CC^T$  and  $C^TC$ : symmetric, share the same **eigenvalues**  $\lambda_1, \lambda_2, \dots$

–  $\lambda_1, \lambda_2, \dots$  are indexed in **decreasing order**

$$\text{SVD: } C = U\Sigma V^T$$

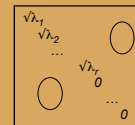
$U$  is  $M \times M$ ; columns are **orthogonal eigenvectors of  $CC^T$**

$V$  is  $N \times N$ ; columns are **orthogonal eigenvectors of  $C^TC$**

$\Sigma$  is  $M \times N$  **diagonal matrix**:  $\Sigma(i,i) = \sqrt{\lambda_i}$  for  $1 \leq i \leq r$

$\Sigma(i,j) = 0$  otherwise

$\sqrt{\lambda_i}$  called **singular values**

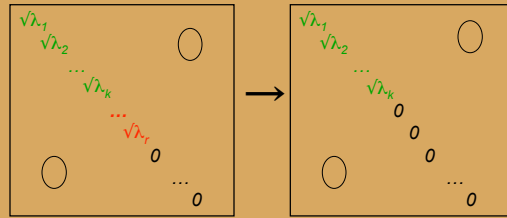


2

## Summary: Reduced Rank Approximation of C

Reduce rank of  $\Sigma$  from  $r$  to  $k$ : keep only  $k$  largest singular values:

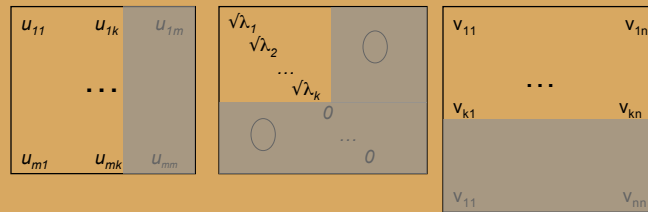
$\Sigma_k$  is  $M \times N$  diagonal matrix:  $\Sigma(i,i) = \sqrt{\lambda_i}$  for  $1 \leq i \leq k$   
 $\Sigma(i,j) = 0$  otherwise



Approximation:  $C_k = U \Sigma_k V^T$

3

## Summary: Reduce dimension matrices Query approximation



$$C_k = U'_k \Sigma'_k V_k'^T$$

$M \times N$        $M \times k$        $k \times k$        $k \times N$

View  $V_k'^T$  as a representation of documents in a  $k$ -dimensional space.

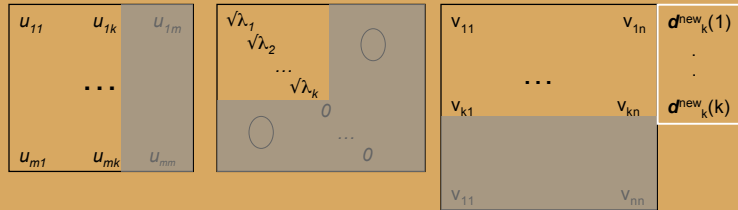
$$V_k'^T = (\Sigma'_k)^{-1} (U'_k)^T C_k \Rightarrow (\Sigma'_k)^{-1} (U'_k)^T q = q_k$$

$$C_k^T C_k = V_k' (\Sigma'_k)^2 (V_k')^T \Rightarrow V_k' (\Sigma'_k)^2 q_k = C_k^T q$$

compares documents      compares doc. to query

4

## Adding a new document to approximation



$$\mathbf{C}_k = \mathbf{U}'_k \mathbf{\Sigma}'_k \mathbf{V}'_k{}^T$$

$M \times (N+1) \quad M \times k \quad k \times k \quad k \times (N+1)$

View as a **representation** of documents in a **k-dimensional space**.

$$\mathbf{V}'_k{}^T = (\mathbf{\Sigma}'_k)^{-1} (\mathbf{U}'_k)^T \mathbf{C}_k \quad \Rightarrow \quad (\mathbf{\Sigma}'_k)^{-1} (\mathbf{U}'_k)^T \mathbf{d}^{new} = \mathbf{d}^{new}_k$$

$\mathbf{d}_{new}$  is new document to add to  $\mathbf{C}_k$   $\Rightarrow$  add column  $\mathbf{d}^{new}_k$  to  $\mathbf{V}'_k{}^T$

5

## Alternative query approximation

**Have:**  $\mathbf{C}_k = \mathbf{U}'_k \mathbf{\Sigma}'_k \mathbf{V}'_k{}^T$

**Just used:**  $\mathbf{V}'_k{}^T = (\mathbf{\Sigma}'_k)^{-1} (\mathbf{U}'_k)^T \mathbf{C}_k$

as a representation of documents in a k-dimensional space

$\Rightarrow$  use  $(\mathbf{\Sigma}'_k)^{-1} (\mathbf{U}'_k)^T \mathbf{q} = \mathbf{q}_k$

- Appropriate for **adding vectors to  $\mathbf{V}'_k{}^T$**  and **comparing queries to docs**.

**Consider again:**  $\mathbf{C}_k{}^T \mathbf{C}_k = \mathbf{V}'_k (\mathbf{\Sigma}'_k)^2 (\mathbf{V}'_k)^T$

gives a second k-dimensional representation of documents:

$$\mathbf{\Sigma}'_k \mathbf{V}'_k{}^T = (\mathbf{U}'_k)^T \mathbf{C}_k$$

$\Rightarrow$  approximate  $\mathbf{q}$  as  $(\mathbf{U}'_k)^T \mathbf{q} = \mathbf{q}'_k$

- Appropriate for **comparing queries to docs**

$$\mathbf{V}'_k (\mathbf{\Sigma}'_k) \mathbf{q}'_k = \mathbf{V}'_k (\mathbf{\Sigma}'_k) (\mathbf{U}'_k)^T \mathbf{q} = \mathbf{C}_k{}^T \mathbf{q}$$

6

## Original LSI paper:

Deerwester, Dumais, et. al.  
***Indexing by Latent Semantic Analysis***  
Journal of the Society for Information Science,  
41(6), 1990, 391-407.

Example from that paper follows

7

Deerwester, Dumais et. al. Table:

Terms	Documents					m1	m2	m3	m4
	c1	c2	c3	c4	c5				
<i>human</i>	1	0	0	1	0	0	0	0	0
<i>interface</i>	1	0	1	0	0	0	0	0	0
<i>computer</i>	1	1	0	0	0	0	0	0	0
<i>user</i>	0	1	1	0	1	0	0	0	0
<i>system</i>	0	1	1	2	0	0	0	0	0
<i>response</i>	0	1	0	0	1	0	0	0	0
<i>time</i>	0	1	0	0	1	0	0	0	0
<i>EPS</i>	0	0	1	1	0	0	0	0	0
<i>survey</i>	0	1	0	0	0	0	0	0	1
<i>trees</i>	0	0	0	0	0	1	1	1	0
<i>graph</i>	0	0	0	0	0	0	1	1	1
<i>minors</i>	0	0	0	0	0	0	0	1	1

8

Deerwester, Dumais et. al. example, cont.:

## Matrix $V^T$ for $k=2$

0.20	0.61	0.46	0.54	0.28	0.00	0.02	0.02	0.08
0.06	0.17	-0.13	-0.23	0.11	0.19	0.44	0.62	0.53

9

Deerwester, Dumais, et al Figure 1

