

Latent Semantic Indexing

1

Introduction

- Vector model => use of theory of linear algebra
- Look at **matrix formulation**:

$$\begin{pmatrix} w_{11} & \dots & w_{t1} \\ \vdots & \ddots & \vdots \\ w_{1N} & \dots & w_{tN} \end{pmatrix} \bullet \begin{pmatrix} w_{1q} \\ w_{tq} \end{pmatrix} = \begin{pmatrix} s_{1q} \\ s_{Nq} \end{pmatrix}$$

document vector query vector scores

$$s_{xq} = \sum_{i=1}^t (w_{ix} * w_{iq})$$

2

1

Goals

- # terms t large - large dimension
 ⇒ reduce dimension
 - find some semantic relationship
 - correlate terms to find structure
 - synonymy
 - polysomy
- “people choose same main terms <20% time”

3

Summary: Singular Value Decomposition (SVD) Review set-up

M - number of terms N - number of documents
C the MxN (term×doc.) matrix of weights (our old w_{ij}) ≥ 0

- of rank r ($r \leq \min(M, N)$)
- documents are columns of C
- compare query to all documents: $C^T q$
 $[N \times M][M \times 1]$

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
- $CC^T(i,j)$ measures strength co-occurrence terms i and j
- $C^TC(i,j)$ measures similarity documents i and j

4

Summary: Singular Value Decomposition (SVD) Theorem

matrix C has a *singular value decomposition*

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

Where:

U M×M matrix

with columns = orthogonal eigenvectors of CC^T

V N×N matrix

with columns = orthogonal eigenvectors of C^TC

Σ M×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*

5

Summary: Reduced Rank Approximation of C

- Reduce rank of Σ from r to **k**
keep only **k** largest singular values

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

6

Summary: Reduced Rank Approximation of C

- Approximation:

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

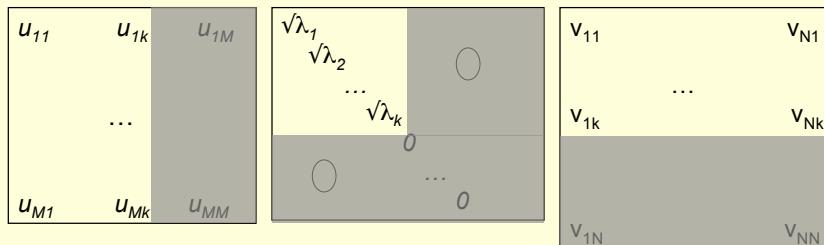
- Theorem:

C_k is the best rank-k approximation to C under the least square fit (Frobenius) norm

$$= \sqrt{\sum_{i=1}^M \sum_{j=1}^N (C(i,j) - C_k(i,j))^2}$$

7

Query Approximation using reduce dimension matrices



$$\mathbf{C}_k = \begin{matrix} \mathbf{U}'_k \\ M \times N \end{matrix}$$

$$\begin{matrix} \mathbf{\Sigma}'_k \\ k \times k \end{matrix} \quad \begin{matrix} \mathbf{V}'^T_k \\ k \times N \end{matrix}$$

8

Using the Approximation

- View V'_k^T as a representation of documents in a k -dimensional space
– a “concept space”?
- Transform query vector \mathbf{q} into that space:

$$C_k^T C_k = (U'_k \Sigma'_k V'_k^T)^T (U'_k \Sigma'_k V'_k^T) = (V'_k \Sigma'_k^T U'_k^T) (U'_k \Sigma'_k V'_k^T)$$

$$= V'_k (\Sigma'_k)^2 (V'_k)^T \quad \text{compares documents}$$

$$\Rightarrow C_k^T \mathbf{q} = V'_k (\Sigma'_k)^2 \mathbf{q}_k \quad \text{compares doc. to query}$$

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

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

recalling $(V'_k^T)(V'_k) = (U'_k^T)(U'_k) = I_9$

Adding a new document

add new document d^{new} to $C_k \Rightarrow$ add column d_k^{new} to V'_k^T

Transform d^{new} into the k -dimensional space version d_k^{new}

$$V'_k^T = (\Sigma'_k)^{-1} (U'_k)^T C_k \Rightarrow (\Sigma'_k)^{-1} (U'_k)^T d^{new} = d_k^{new}$$

$u_{11} \quad u_{1k} \quad u_{1M}$	\vdots	$u_{M1} \quad u_{Mk} \quad u_{MM}$	$\begin{matrix} \sqrt{\lambda}_1 & & \\ & \sqrt{\lambda}_2 & \\ & \dots & \\ & & \sqrt{\lambda}_k \\ & & 0 \\ & \circlearrowleft & \end{matrix}$	$\begin{matrix} v_{11} & & & d_k^{new}(1) \\ \dots & & & \vdots \\ v_{1k} & & & d_k^{new}(k) \\ & \dots & & \\ v_{1N} & & & v_{NN} \end{matrix}$
------------------------------------	----------	------------------------------------	--	--

$$C_k = \begin{matrix} U'_k \\ M \times (N+1) \end{matrix} \quad U'_k \quad M \times k$$

$$\Sigma'_k \quad k \times k$$

$$V'_k^T \quad k \times (N+1)$$

10

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

11

Deerwester, Dumais et. al. Table:

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

12

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

13

Deerwester, Dumais, et al Figure 1

2-D Plot of Terms and Docs from Example

