

Latent Semantic Indexing

1

Introduction

- Vector model => use of theory of linear algebra
- Look at **matrix formulation**
 - M - number of terms in lexicon
 - N - number of documents in collection
 - C the $M \times N$ (term \times doc.) **matrix of weights** ≥ 0 (our old w_{ij})

$$\begin{pmatrix} c_{11} & \dots & c_{M1} \\ \dots & \dots & \dots \\ c_{1N} & \dots & c_{MN} \end{pmatrix} \bullet \begin{pmatrix} w_{1q} \\ \dots \\ w_{Mq} \end{pmatrix} = \begin{pmatrix} s_{1q} \\ \dots \\ s_{Nq} \end{pmatrix}$$

document vector query vector scores

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

2

Goals

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

3

Set-up

- C the $M \times N$ (term \times doc.) matrix of non-negative weights
- of rank r ($r \leq \min(M, N)$)
 - documents are columns of C

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

4

Use Singular Value Decomposition (SVD)

Theorem:

$M \times N$ matrix C of rank r has a

singular value decomposition $C = U\Sigma V^T$

Where:

U $M \times M$ matrix

with columns = *orthogonal eigenvectors of CC^T*

V $N \times N$ matrix

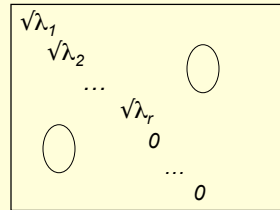
with columns = *orthogonal eigenvectors of C^TC*

Σ $M \times N$ *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*

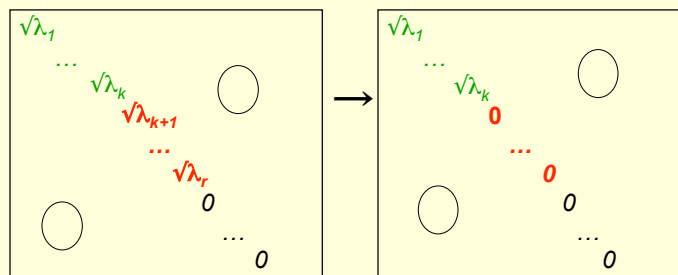


5

Reduce Rank

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

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



6

Reduced Rank Approximation of C

- Approximation:

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

[M×N] [M×M] [M×N] [N×N]

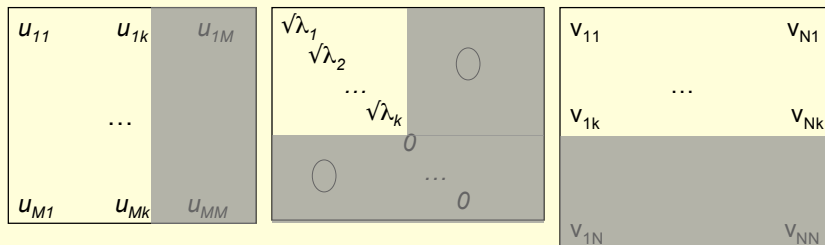
- Theorem:

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

Reduced dimension matrices



$$C_k =$$

M×N

$$U'_k$$

M×k

$$\Sigma'_k$$

k×k

$$V'_k{}^T$$

k×N

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:

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

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

$$\Rightarrow \mathbf{q}_k = (\Sigma'_k{}^{-1})^2 V'_k{}^T C_k{}^T \mathbf{q} = (\Sigma'_k{}^{-1})^2 V'_k{}^T V'_k \Sigma'_k{}^T U'_k{}^T \mathbf{q} \\ = (\Sigma'_k{}^{-1}) (U'_k)^T \mathbf{q}$$

$$\text{recalling } (V'_k{}^T)(V'_k) = (U'_k{}^T)(U'_k) = I$$

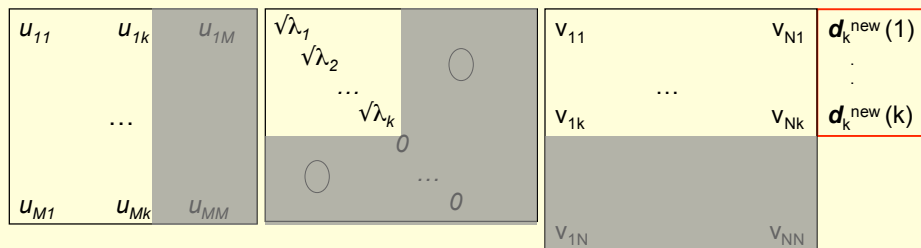
9

Adding a new document

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

Transform \mathbf{d}^{new} into the **k-dimensional space** version \mathbf{d}_k^{new}

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



$$C_k = \begin{matrix} M \times (N+1) \\ U'_k \\ M \times k \\ \Sigma'_k \\ k \times k \\ V'_k{}^T \\ k \times (N+1) \end{matrix}$$

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					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

12

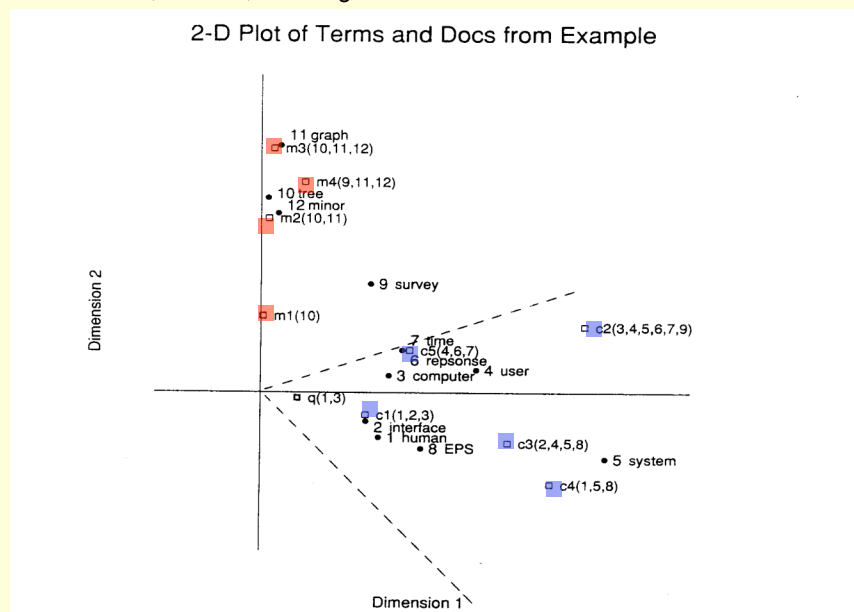
Deerwester, Dumais et. al. example, cont.:

Matrix $V'_k{}^T$ 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



Summary

- LSI uses SVD to get a **reduced-rank** and **reduced-size** approximation to C
- LSI can be viewed as a **preprocessor** for
 - query evaluation
 - clustering
- SVD **computation** can be **costly**
 - do once (or rarely)