# 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×N (term×doc.) matrix of weights ≥ 0 *(our old $w_{ij}$ )*

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

document vector    query vector    scores

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

2

---

## 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×N (term×doc.) matrix of weights ≥ 0 *(our old $w_{ij}$ )*

$$\begin{bmatrix} c_{11} & ... & c_{M1} \\ \mathbf{C^T} & & \\ c_{1N} & ... & c_{MN} \end{bmatrix} \bullet \begin{bmatrix} w_{1q} \\ \mathbf{q} \\ w_{Mq} \end{bmatrix} = \begin{bmatrix} s_{1q} \\ \\ s_{Nq} \end{bmatrix}$$

document vector    query vector    scores

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

3

---

## Goals

- # terms M large - large dimension
  - ⇒reduce dimension

- find some semantic relationship
  - correlate terms to find structure
    - synonomy
    - polysomy
  "people choose same main terms <20% time"

4

---

## Set-up

C  the M×N (term×doc.) matrix of non-negative weights
  - of rank r  ( r ≤ min(M,N) )
  - documents are *columns* of C

consider $CC^T$ and $C^TC$:
- symmetric,
- share the same eigenvalues $\lambda_1, \lambda_2,...$
  - $\lambda_1, \lambda_2, ...$ 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

5

---

## Use Singular Value Decomposition (SVD)

**Theorem:**

M×N matrix C of rank r has a

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

Where:

U  M×M matrix
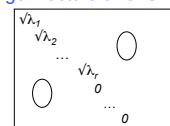  with columns = orthogonal eigenvectors of $CC^T$

V  N×N matrix
  with columns = orthogonal eigenvectors of $C^TC$

$\Sigma$  M×N diagonal matrix:
  $\Sigma(i,i) = \sqrt{\lambda_i}$  for $1 \le i \le r$
  $\Sigma(i,j) = 0$  otherwise

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

$$\begin{bmatrix} \sqrt{\lambda_1} & & & & \\ & \sqrt{\lambda_2} & & & \\ & & ... & & \\ & & & \sqrt{\lambda_r} & \\ & & & & 0 \\ & & & & ... \\ & & & & 0 \end{bmatrix}$$

6

## Reduce Rank

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

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



7

## 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
$$= \sqrt{\sum^M_{i=1} \sum^N_{j=1} (C(i,j) - C_k(i,j))^2}$$

8

## Reduced dimension matrices



$C_k =$    $U'_k$      $\Sigma'_k$      $V'_k{}^T$
M×N    M×k      k×k      k×N

9

## Using the Approximation

- View $V'_k{}^T$ as a representation of documents in a k-dimensional space
  – a "concept space"?

- Transform query vector **q** into that space:

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

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

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

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

10

## Adding a new document

add new document $\boldsymbol{d}^{new}$ to $C_k$ => add column $\boldsymbol{d}_k{}^{new}$ to $V'_k{}^T$

Transform $\boldsymbol{d}^{new}$ into the k-dimensional space version $\boldsymbol{d}_k{}^{new}$

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



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

11

## 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

12

Deerwester, Dumais et. al. Table:

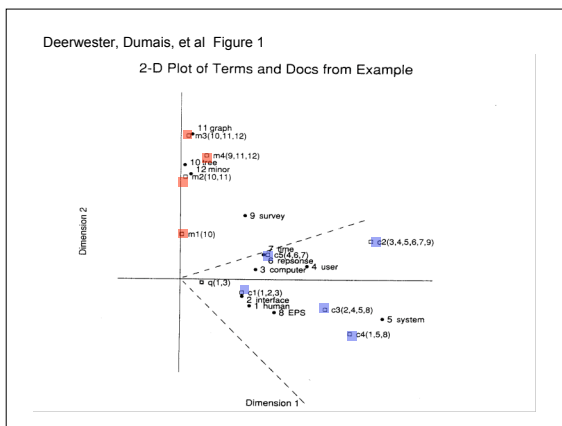| Terms | c1 | c2 | c3 | Documents c4 | c5 | m1 | m2 | m3 | m4 |
|---|---|---|---|---|---|---|---|---|---|
| human | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| interface | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| computer | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| user | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| system | 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 |

13

---

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

14

---

Deerwester, Dumais, et al  Figure 1

2-D Plot of Terms and Docs from Example



16

# 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)

16

---

# Another application of SVD: collaborative filtering

Modeling Relationships at Multiple Scales to Improve Accuracy of Large Recommender Systems,  Robert M. Bell, Yehuda Koren and Chris Volinsky, *KDD 07*

•one of methods used for Netflix challenge
•find factors to describe user and items
•use to fill in vaues of unknown ratings of items by users

17