## Latent Semantic Indexing: Introduction

- Analysis of term-document interaction for corpus of text documents
- Standard vector model:
  - document vector of term weights
- Goals:
  - reduce dimension of document vectors
  - uncover latent factors:
    - document as vector of factor weights
- uses of theory of linear algebra

1

## 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 $w_{ij}$ )*

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

$C^T$ • q =

document vector      query vector      scores

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

2

## 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 $\lambda_1, \lambda_2,\dots$
  - $\lambda_1, \lambda_2, \dots$ 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

3

## 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
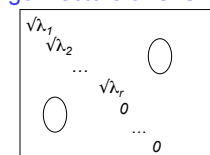   with columns = orthogonal eigenvectors of $CC^T$
V  N×N matrix
   with columns = orthogonal eigenvectors of $C^TC$
Σ M×N diagonal matrix:
   $\Sigma(i,i) = \sqrt{\lambda_i}$   for *1 ≤ i ≤ r*
   $\Sigma(i,j) = 0$  otherwise
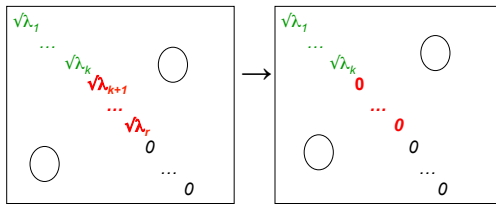$\sqrt{\lambda_i}$ *called singular values*

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

4

## 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



5

---

## Reduced Rank Approximation of C

- Approximation:

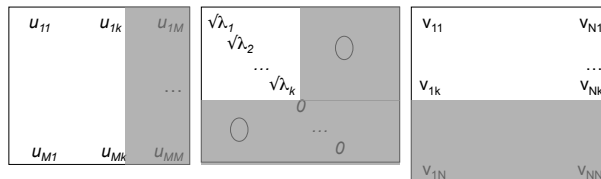  $$C_k = U\Sigma_k V^T$$

  [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

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

6

---

## Reduced dimension matrices



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

7

---

## Semantic Interpretation

- remaining k dimensions:   k factors
- View $V'_k{}^T$ as a representation of documents in the k-dimensional space
- View $U'_k{}^T$ as a representation of terms in the k-dimensional space
- $\Sigma_k$ scales between them

- find some semantic relationship?
  - "concept space"?
  - correlating terms to find structure
    - synonomy
    - polysomy
    "people choose same main terms <20% time"

8

2

## Using the Approximation

- $V'_k{}^T$ as a representation of documents in a k-dimensional space

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

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

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

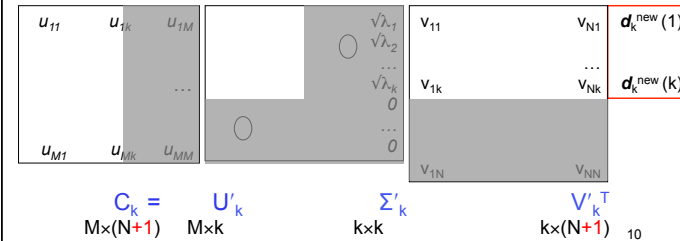  Then $(\Sigma'_k)^{-1} (U'_k)^T \boldsymbol{q} = \boldsymbol{q}_k$

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

9

## 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 = $    $U'_k$    $\Sigma'_k$    $V'_k{}^T$
$M \times (N+1)$   $M \times k$    $k \times k$    $k \times (N+1)$   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 | 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 |

12

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

## Matrix $V'^T_k$ 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

# 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)

15

4