

---

BLG 540E  
TEXT RETRIEVAL SYSTEMS

Latent Semantic Indexing

Arzucan Özgür

# Today's topic

---

- ▶ **Latent Semantic Indexing**
  - ▶ Term-document matrices are very large
  - ▶ But the number of topics that people talk about is small (in some sense)
    - ▶ Clothes, movies, politics, ...
  - ▶ Can we represent the term-document space by a lower dimensional latent space?



# Vector Space Model: Pros

---

- ▶ **Automatic** selection of index terms
- ▶ **Partial matching** of queries and documents (*dealing with the case where no document contains all search terms*)
- ▶ **Ranking** according to **similarity score** (*dealing with large result sets*)
- ▶ **Term weighting** schemes (*improves retrieval performance*)
- ▶ Various extensions
  - ▶ Document clustering
  - ▶ Relevance feedback (modifying query vector)
- ▶ Geometric foundation



# Problems with Lexical Semantics

---

- ▶ Ambiguity and association in natural language
  - ▶ **Polysemy**: Words often have a **multitude of meanings** and different types of usage (*more severe in very heterogeneous collections*).
    - ▶ bank, jaguar, hot
  - ▶ The vector space model is unable to discriminate between different meanings of the same word.

$$\text{sim}_{\text{true}}(d, q) < \cos(\angle(\vec{d}, \vec{q}))$$



# Problems with Lexical Semantics

---

- ▶ **Synonymy**: Different terms may have an **identical or a similar meaning**
  - ▶ Large/big, Spicy/hot, Car/automobile
- ▶ No associations between words are made in the vector space representation.

$$\text{sim}_{\text{true}}(d, q) > \cos(\angle(\vec{d}, \vec{q}))$$



# Latent Semantic Indexing (LSI)

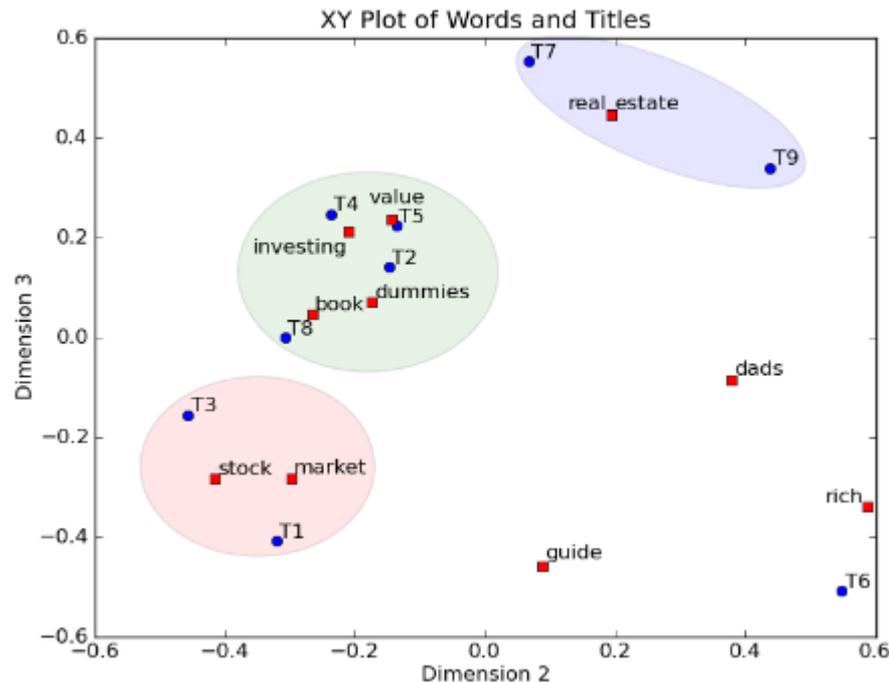
---

- ▶ Perform a **low-rank approximation** of **document-term matrix** (typical rank **100-300**)
  - ▶ General idea
    - ▶ Map documents (*and* terms) to a **low-dimensional representation**.
    - ▶ Design a mapping such that the low-dimensional space reflects **semantic associations** (latent semantic space, identification of hidden (latent) concepts).
    - ▶ Compute document similarity based on the **inner product** in this **latent semantic space**
- 



# Latent Semantic Analysis

- ▶ **Latent semantic space:** illustrating example
- ▶ Similar words and documents mapped to similar locations in the lower dimensional latent space.



---

# Linear Algebra Background

---



# Eigenvalues & Eigenvectors

- ▶ **Eigenvectors** (for a square  $m \times m$  matrix  $\mathbf{S}$ )

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$$

(right) eigenvector      eigenvalue  
 $\mathbf{v} \in \mathbb{R}^m \neq \mathbf{0}$        $\lambda \in \mathbb{R}$

*Example*

$$\begin{pmatrix} 6 & -2 \\ 4 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

- ▶ **How many eigenvalues** are there at most?

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v} \iff (\mathbf{S} - \lambda\mathbf{I})\mathbf{v} = \mathbf{0}$$

only has a non-zero solution if  $|\mathbf{S} - \lambda\mathbf{I}| = 0$

This is a  $m$ th order equation in  $\lambda$  which can have **at most  $m$  distinct solutions** (roots of the characteristic polynomial) - can be complex even though  $\mathbf{S}$  is real.

# Eigenvectors and eigenvalues

---

▶ **Example:**

$$S = \begin{pmatrix} -1 & 3 \\ 2 & 0 \end{pmatrix} \quad S - \lambda I = \begin{pmatrix} -1 - \lambda & 3 \\ 2 & -\lambda \end{pmatrix}$$

- ▶  $|S - \lambda I| = (-1 - \lambda)(-\lambda) - 3 \cdot 2 = 0$
- ▶ Then:  $\lambda + \lambda^2 - 6 = 0$ ;  $\lambda_1 = 2$ ;  $\lambda_2 = -3$
- ▶ For  $\lambda_1 = 2$ :

$$\begin{pmatrix} -3 & 3 \\ 2 & -2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = 0$$

- ▶ Solutions:  $x_1 = x_2$



# Matrix-vector multiplication

---

$$S = \begin{bmatrix} 30 & 0 & 0 \\ 0 & 20 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

has eigenvalues 30, 20, 1 with corresponding eigenvectors

$$v_1 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad v_2 = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \quad v_3 = \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$$

Any vector (say  $x = \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix}$ ) can be viewed as a combination of the eigenvectors:  
$$x = 2v_1 + 4v_2 + 6v_3$$



# Matrix vector multiplication

---

- ▶ Thus a matrix-vector multiplication such as  $Sx$  ( $S, x$  as in the previous slide) can be rewritten in terms of the eigenvalues/vectors:

$$Sx = S(2v_1 + 4v_2 + 6v_3)$$

$$Sx = 2 Sv_1 + 4 Sv_2 + 6 Sv_3 = 2\lambda_1 v_1 + 4\lambda_2 v_2 + 6\lambda_3 v_3$$

$$Sx = 60v_1 + 80v_2 + 6v_3$$

- ▶ Even though  $x$  is an arbitrary vector, the action of  $S$  on  $x$  is determined by the eigenvalues/vectors.
- 



# Matrix vector multiplication

---

- ▶ Suggestion: the effect of “small” eigenvalues is small.
- ▶ If we ignored the smallest eigenvalue (1), then instead of

$$\begin{pmatrix} 60 \\ 80 \\ 6 \end{pmatrix} \quad \text{we would get} \quad \begin{pmatrix} 60 \\ 80 \\ 0 \end{pmatrix}$$

- ▶ These vectors are similar (in cosine similarity, etc.)



# Eigenvalues & Eigenvectors

For symmetric matrices, eigenvectors for distinct eigenvalues are **orthogonal**

$$Sv_{\{1,2\}} = \lambda_{\{1,2\}} v_{\{1,2\}}, \text{ and } \lambda_1 \neq \lambda_2 \Rightarrow v_1 \cdot v_2 = 0$$

All eigenvalues of a real symmetric matrix are **real**.

$$\text{if } |S - \lambda I| = 0 \text{ and } S = S^T \Rightarrow \lambda \in \mathfrak{R}$$

All eigenvalues of a **positive semidefinite matrix** are **non-negative**

$$\forall w \in \mathfrak{R}^n, w^T S w \geq 0, \text{ then if } S v = \lambda v \Rightarrow \lambda \geq 0$$

# Example

---

▶ Let

$$S = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \longleftarrow \boxed{\text{Real, symmetric.}}$$

▶ Then

$$S - \lambda I = \begin{bmatrix} 2 - \lambda & 1 \\ 1 & 2 - \lambda \end{bmatrix} \Rightarrow$$

$$|S - \lambda I| = (2 - \lambda)^2 - 1 = 0.$$

▶ The eigenvalues are 1 and 3 (nonnegative, real).

▶ The eigenvectors are orthogonal (and real):

$$\begin{pmatrix} 1 \\ -1 \end{pmatrix} \quad \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

Plug in these values and solve for eigenvectors.

# Eigen/diagonal Decomposition

---

Let  $\mathbf{S} \in \mathbb{R}^{m \times m}$  be a **square** matrix with  **$m$  linearly independent eigenvectors**

**Theorem:** There exists an **eigen decomposition**

(cf. matrix diagonalization theorem)

$$\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}$$

*diagonal*

- ▶ Columns of  $\mathbf{U}$  are **eigenvectors** of  $\mathbf{S}$
- ▶ Diagonal elements of  $\mathbf{\Lambda}$  are **eigenvalues** of  $\mathbf{S}$

$$\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_m), \quad \lambda_i \geq \lambda_{i+1}$$



# Diagonal decomposition: why/how

---

Let  $\mathbf{U}$  have the eigenvectors as columns: 
$$\mathbf{U} = \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{bmatrix}$$

Then,  $\mathbf{S}\mathbf{U}$  can be written

$$\mathbf{S}\mathbf{U} = \mathbf{S} \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{bmatrix} = \begin{bmatrix} \lambda_1 \mathbf{v}_1 & \dots & \lambda_n \mathbf{v}_n \end{bmatrix} = \begin{bmatrix} \mathbf{v}_1 & \dots & \mathbf{v}_n \end{bmatrix} \begin{bmatrix} \lambda_1 & & \\ & \dots & \\ & & \lambda_n \end{bmatrix}$$

Thus  $\mathbf{S}\mathbf{U} = \mathbf{U}\mathbf{\Lambda}$ , or  $\mathbf{U}^{-1}\mathbf{S}\mathbf{U} = \mathbf{\Lambda}$

And  $\mathbf{S} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^{-1}$ .

---



# Diagonal decomposition - example

---

Recall  $S = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}; \lambda_1 = 1, \lambda_2 = 3.$

The eigenvectors  $\begin{pmatrix} 1 \\ -1 \end{pmatrix}$  and  $\begin{pmatrix} 1 \\ 1 \end{pmatrix}$  form  $U = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$

Inverting, we have  $U^{-1} = \begin{bmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{bmatrix}$

← Recall  
 $UU^{-1} = I.$

Then,  $S = U\Lambda U^{-1} = \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix} \begin{bmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{bmatrix}$



## Example continued

---

Let's divide  $\mathbf{U}$  (and multiply  $\mathbf{U}^{-1}$ ) by  $\sqrt{2}$

$$\text{Then, } \mathbf{S} = \underbrace{\begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}}_{\mathbf{Q}} \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 3 \end{bmatrix}}_{\mathbf{\Lambda}} \underbrace{\begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} \\ 1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}}_{(\mathbf{Q}^{-1} = \mathbf{Q}^T)}$$



# Symmetric Eigen Decomposition

---

- ▶ If  $S \in \mathbb{R}^{m \times m}$  is a **symmetric** matrix:
- ▶ **Theorem:** There exists a (unique) **eigen decomposition**

$$S = Q\Lambda Q^T$$

- ▶ where **Q** is **orthogonal:**
  - ▶  $Q^{-1} = Q^T$
  - ▶ Columns of **Q** are normalized eigenvectors
  - ▶ Columns are orthogonal.
  - ▶ (everything is real)



everything so far needs square matrices

---

- ▶ Recall  $M \times N$  term-document matrices ...



# Singular Value Decomposition

---

For an  $M \times N$  matrix  $\mathbf{A}$  of rank  $r$  there exists a factorization (Singular Value Decomposition = **SVD**) as follows:

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

$M \times M$

$M \times N$

$V \text{ is } N \times N$

The columns of  $\mathbf{U}$  are orthogonal eigenvectors of  $\mathbf{A}\mathbf{A}^T$ .  
 The columns of  $\mathbf{V}$  are orthogonal eigenvectors of  $\mathbf{A}^T\mathbf{A}$ .  
 Eigenvalues  $\lambda_1 \dots \lambda_r$  of  $\mathbf{A}\mathbf{A}^T$  are the eigenvalues of  $\mathbf{A}^T\mathbf{A}$ .

$$\sigma_i = \sqrt{\lambda_i}$$

$$\Sigma = \text{diag}(\sigma_1 \dots \sigma_r)$$

Singular values.

▶ In Matlab, use `[U,S,V] = svd (A)`

# SVD example

---

$$\text{Let } A = \begin{bmatrix} 1 & -1 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$$

Thus  $M=3$ ,  $N=2$ . Its SVD is

$$\begin{bmatrix} 0 & 2/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & -1/\sqrt{6} & 1/\sqrt{3} \\ 1/\sqrt{2} & 1/\sqrt{6} & -1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & \sqrt{3} \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}$$

---

► Typically, the singular values arranged in decreasing order.

# Low-rank Approximation

---

- ▶ SVD can be used to compute optimal **low-rank approximations**.
- ▶ Approximation problem: Find  $\mathbf{A}_k$  of rank  $k$  such that

$$A_k = \min_{X : \text{rank}(X) = k} \|A - X\|_F \longleftarrow \text{Frobenius norm}$$

$$\|A\|_F \equiv \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2}.$$

$A_k$  and  $X$  are both  $m \times n$  matrices.

Typically, want  $k \ll r$ .

---



# Low-rank Approximation

## ► Solution via SVD

$$A_k = U \operatorname{diag}(\sigma_1, \dots, \sigma_k, \underbrace{0, \dots, 0}_{\substack{\text{set smallest } r-k \\ \text{singular values to zero}}}) V^T$$

$$\underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{A_k} = \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & & & \\ & \bullet & & & \\ & & & & \\ & & & & \\ & & & & \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T}$$

# Reduced SVD

- ▶ If we retain only  $k$  singular values, and set the rest to 0, then we don't need the matrix parts in brown
- ▶ Then  $\Sigma$  is  $k \times k$ ,  $U$  is  $M \times k$ ,  $V^T$  is  $k \times N$ , and  $A_k$  is  $M \times N$
- ▶ This is referred to as the reduced SVD
- ▶ It is the convenient (space-saving) and usual form for computational applications

$$\underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{A} = \underbrace{\begin{bmatrix} * & * & * \\ * & * & * \\ * & * & * \end{bmatrix}}_U \underbrace{\begin{bmatrix} \bullet & & & & \\ & \bullet & & & \\ & & & & \end{bmatrix}}_{\Sigma} \underbrace{\begin{bmatrix} * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \\ * & * & * & * & * \end{bmatrix}}_{V^T}$$

# Approximation error

---

- ▶ How good (bad) is this approximation?
- ▶ It's the best possible, measured by the Frobenius norm of the error:

$$\min_{X : \text{rank}(X)=k} \|A - X\|_F = \|A - A_k\|_F = \sigma_{k+1}$$

where the  $\sigma_i$  are ordered such that  $\sigma_i \geq \sigma_{i+1}$ .

Suggests why Frobenius error drops as  $k$  increased.



# SVD Low-rank approximation

---

- ▶ Whereas the term-doc matrix  $A$  may have  $M=50000$ ,  $N=10$  million (and rank close to 50000)
- ▶ We can construct an approximation  $A_{100}$  with rank 100.
  - ▶ Of all rank 100 matrices, it would have the lowest Frobenius error.
- ▶ Great ... but why would we??
- ▶ Answer: *Latent Semantic Indexing*

---

# Latent Semantic Indexing via the SVD



## What it is

---

- ▶ From term-doc matrix  $A$ , we compute the approximation  $A_k$ .
- ▶ There is a row for each term and a column for each doc in  $A_k$
- ▶ Thus docs live in a space of  $k \ll r$  dimensions
  - ▶ These dimensions are not the original axes



# Example

---

- ▶ Query: *gold silver truck*
- ▶ Documents:
  - ▶ d1: *Shipment of gold damaged in a fire.*
  - ▶ d2: *Delivery of silver arrived in a silver truck.*
  - ▶ d3: *Shipment of gold arrived in a truck.*

Terms	d1	d2	d3	q
a	1	1	1	0
arrived	0	1	1	0
damaged	1	0	0	0
delivery	0	1	0	0
fire	1	0	0	0
gold	1	0	1	1
in	1	1	1	0
of	1	1	1	0
shipment	1	0	1	0
silver	0	2	0	1
truck	0	1	1	1



# Compute the SVD of A

---

$$\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

$$\mathbf{U} = \begin{bmatrix} -0.4201 & 0.0748 & -0.0460 \\ -0.2995 & -0.2001 & 0.4078 \\ -0.1206 & 0.2749 & -0.4538 \\ -0.1576 & -0.3046 & -0.2006 \\ -0.1206 & 0.2749 & -0.4538 \\ -0.2626 & 0.3794 & 0.1547 \\ -0.4201 & 0.0748 & -0.0460 \\ -0.4201 & 0.0748 & -0.0460 \\ -0.2626 & 0.3794 & 0.1547 \\ -0.3151 & -0.6093 & -0.4013 \\ -0.2995 & -0.2001 & 0.4078 \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} 4.0989 & 0.0000 & 0.0000 \\ 0.0000 & 2.3616 & 0.0000 \\ 0.0000 & 0.0000 & 1.2737 \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} -0.4945 & 0.6492 & -0.5780 \\ -0.6458 & -0.7194 & -0.2556 \\ -0.5817 & 0.2469 & 0.7750 \end{bmatrix} \quad \mathbf{V}^T = \begin{bmatrix} -0.4945 & -0.6458 & -0.5817 \\ 0.6492 & -0.7194 & 0.2469 \\ -0.5780 & -0.2556 & 0.7750 \end{bmatrix}$$



# Rank 2 Approximation

---

$$\mathbf{U} \approx \mathbf{U}_k = \begin{bmatrix} -0.4201 & 0.0748 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix}$$

$k = 2$

$$\mathbf{S} \approx \mathbf{S}_k = \begin{bmatrix} 4.0989 & 0.0000 \\ 0.0000 & 2.3616 \end{bmatrix}$$

$$\mathbf{V} \approx \mathbf{V}_k = \begin{bmatrix} -0.4945 & 0.6492 \\ -0.6458 & -0.7194 \\ -0.5817 & 0.2469 \end{bmatrix}$$

$$\mathbf{V}^T \approx \mathbf{V}_k^T = \begin{bmatrix} -0.4945 & -0.6458 & -0.5817 \\ 0.6492 & -0.7194 & 0.2469 \end{bmatrix}$$



# Computing the query vector

---

$$A^T = (USV^T)^T = VSU^T$$

$$A^TUS^{-1} = VSU^TUS^{-1}$$

$$V = A^TUS^{-1}$$

$$d = d^TUS^{-1}$$

$$q = q^TUS^{-1}$$

Thus, in the reduced  $k$ -dimensional space we can write

$$d = d^T U_k S_k^{-1}$$

$$q = q^T U_k S_k^{-1}$$



# Computing the query vector

---

$$\mathbf{q} = \mathbf{q}^T \mathbf{U}_k \mathbf{S}_k^{-1}$$

$$\mathbf{q} = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} -0.4201 & 0.0748 \\ -0.2995 & -0.2001 \\ -0.1206 & 0.2749 \\ -0.1576 & -0.3046 \\ -0.1206 & 0.2749 \\ -0.2626 & 0.3794 \\ -0.4201 & 0.0748 \\ -0.4201 & 0.0748 \\ -0.2626 & 0.3794 \\ -0.3151 & -0.6093 \\ -0.2995 & -0.2001 \end{bmatrix} \begin{bmatrix} \frac{1}{4.0989} & 0.0000 \\ 0.0000 & \frac{1}{2.3616} \end{bmatrix}$$

$$\mathbf{q} = \begin{bmatrix} -0.2140 & -0.1821 \end{bmatrix}$$



# Computing the similarity

---

d1(-0.4945, 0.6492)

d2(-0.6458, -0.7194)

d3(-0.5817, 0.2469)

$$\text{sim}(q, d) = \frac{q \bullet d}{|q| |d|}$$

$$\text{sim}(q, d_1) = \frac{(-0.2140)(-0.4945) + (-0.1821)(0.6492)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.4945)^2 + (0.6492)^2}} = -0.0541$$

$$\text{sim}(q, d_2) = \frac{(-0.2140)(-0.6458) + (-0.1821)(-0.7194)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.6458)^2 + (-0.7194)^2}} = 0.9910$$

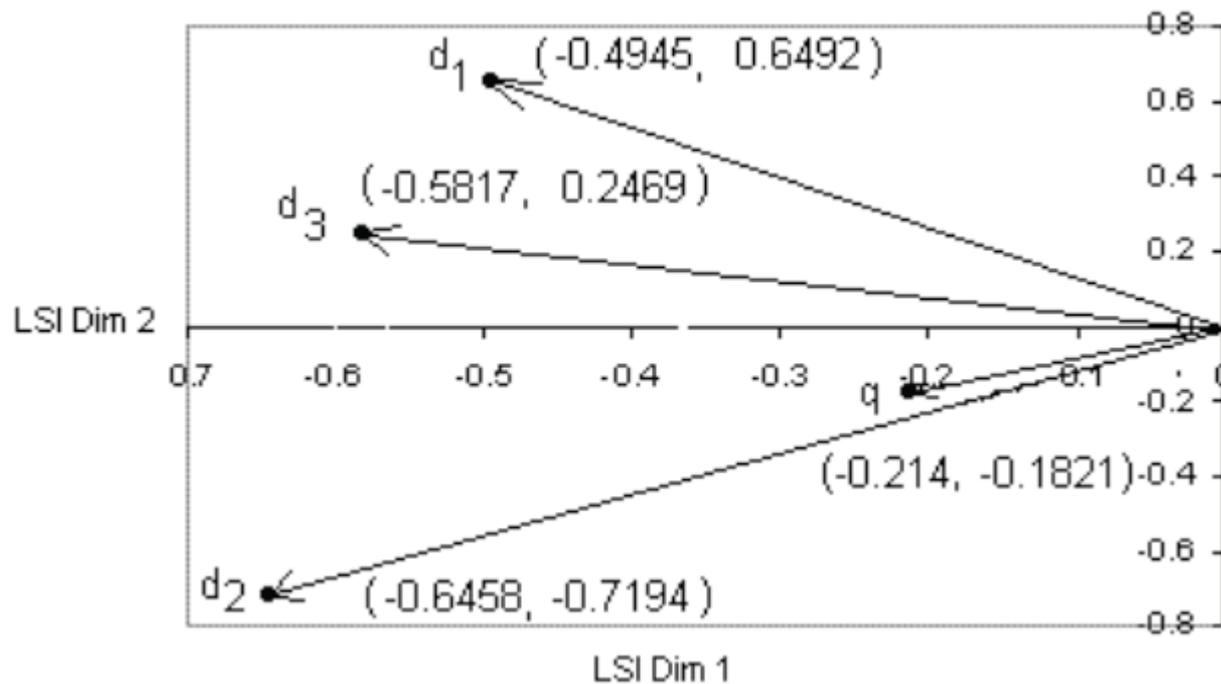
$$\text{sim}(q, d_3) = \frac{(-0.2140)(-0.5817) + (-0.1821)(0.2469)}{\sqrt{(-0.2140)^2 + (-0.1821)^2} \sqrt{(-0.5817)^2 + (0.2469)^2}} = 0.4478$$

Ranking documents in descending order

$$d_2 > d_3 > d_1$$



# LSI Query Document Vectors



# Resources

---

- ▶ *Introduction to Information Retrieval*, chapter 18.
- ▶ Some slides were adapted from
  - ▶ Prof. Dragomir Radev's lectures at the University of Michigan:
    - ▶ <http://clair.si.umich.edu/~radev/teaching.html>
  - ▶ the book's companion website:
    - ▶ <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>
  - ▶ SVD and LSI Tutorial:
    - ▶ <http://www.miislita.com/information-retrieval-tutorial/svd-lsi-tutorial-4-lsi-how-to-calculations.html>

