

Intelligent Data Analysis

Week 5 – Lecture 9

# **Latent Semantic Analysis (LSA)**

Martin Russell

# Objectives

- To understand, intuitively, how **Latent Semantic Analysis (LSA)** can
  - Discover latent **topics** in a corpus and represent them in terms of words
  - Achieve dimension reduction for document vectors
  - Represent words in terms of topics
- To understand the relationship between LSA and PCA applied to a set of document vectors

# Vector Notation for Documents

- Suppose that we have a set of documents

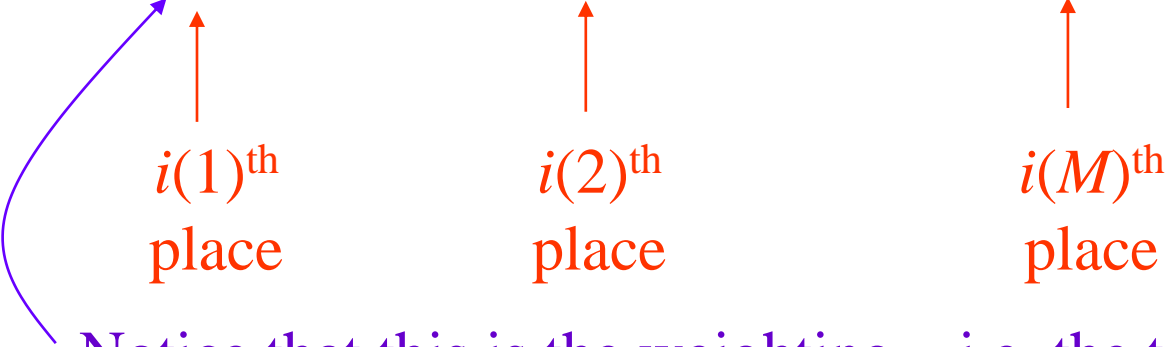
$$D = \{d_1, d_2, \dots, d_N\}$$

think of this as the corpus for IR

- Suppose that the number of **different** words in the **whole corpus** is  $V$  (vocabulary size)
- Now suppose a document  $d$  in  $D$  contains  $M$  different terms:  $\{t_{i(1)}, t_{i(2)}, \dots, t_{i(M)}\}$
- Finally, suppose term  $t_{i(m)}$  occurs  $f_{i(m)}$  times

# Vector Notation

- The **vector representation**  $vec(d)$  of  $d$  is the  $V$  dimensional vector:

$$(0, \dots, 0, w_{i(1),d}, 0, \dots, 0, w_{i(2),d}, 0, \dots, 0, w_{i(M),d}, 0, \dots, 0)$$


$i(1)^{\text{th}}$   
place

$i(2)^{\text{th}}$   
place

$i(M)^{\text{th}}$   
place

Notice that this is the weighting – i.e. the term frequency times the inverse document frequency  $w_{i(1)} = f_{i(1)d} \times IDF(i(1))$  from text IR

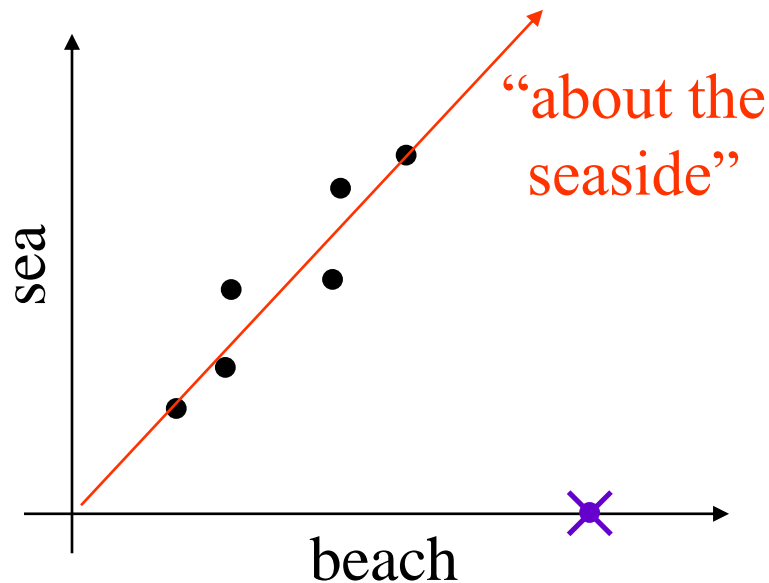
- $vec(d)$  is the **document vector** for  $d$

# Latent Semantic Analysis (LSA)

- Suppose we have a real corpus with a large number of documents
- For each document  $d$  the dimension of the vector  $vec(d)$  is potentially several thousands
- Let's focus on just 2 of these dimensions, corresponding, say, to the words 'sea' and 'beach'
- Intuitively, often, when a document  $d$  includes 'sea' it will also include 'beach'

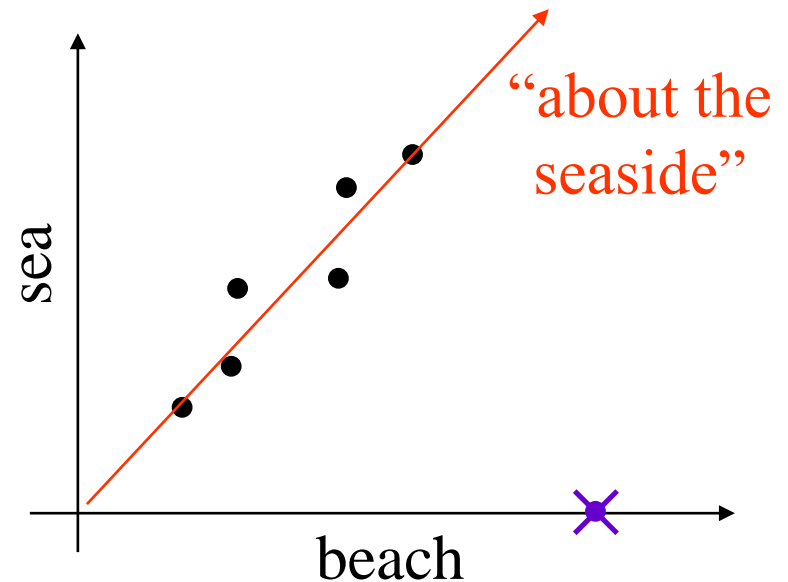
# LSA continued

- Equivalently, if  $\text{vec}(d)$  has a non-zero entry in the 'sea' component, it will often have a non-zero entry in the 'beach' component



# Latent Semantic Classes

- If we can detect this type of structure, then we can discover relationships between words **automatically**, from **data**
- In the example we have found an equivalence set of terms, including ‘beach’ and ‘sea’, which is **about the seaside**



# Finding Latent Semantic Classes

- LSA involves some advanced linear algebra - the description here is just an outline
- First construct the 'word-document' matrix  $A$
- Then decompose  $A$  using **Singular Value Decomposition (SVD)**
  - SVD is a standard technique from matrix algebra
  - Packages such as MATLAB have SVD functions:  
$$>> [U, S, V] = \text{svd}(A)$$



# Singular Value Decomposition

- Recall **eigenvector decomposition**
- An eigenvector of a square matrix  $A$  is a vector  $e$  such that  $Ae = \lambda_e e$ , where  $\lambda$  is a scalar.
- For certain matrices  $A$  we can write  $A = UDU^T$ , where  $U$  is an **orthogonal matrix** and  $D$  is **diagonal (Spectral theorem)**
  - The elements of  $D$  are the eigenvalues
  - The columns of  $U$  are the eigenvectors
- You can think of SVD as a more general version of eigenvector decomposition, which works for **general** matrices

# Word-Document Matrix

- The **word-document matrix** is the  $N \times V$  matrix whose  $n$ th row is the document vector for the  $n$ th document

$$A = \begin{bmatrix} w_{t_1 d_1} & w_{t_2 d_1} & \cdots & w_{t_m d_1} & \cdots & w_{t_V d_1} \\ w_{t_1 d_2} & w_{t_2 d_2} & \cdots & w_{t_m d_2} & \cdots & w_{t_V d_2} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ w_{t_1 d_n} & w_{t_2 d_n} & \cdots & w_{t_m d_n} & \cdots & w_{t_V d_n} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ w_{t_1 d_N} & w_{t_2 d_N} & \cdots & w_{t_m d_N} & \cdots & w_{t_V d_N} \end{bmatrix}$$

Weighting for term  $t_m$  in  $d_n$

# Singular Value Decomposition (SVD)

$$A = USV^T$$

Direction of most significant correlation

$$A = \begin{bmatrix} u_{11} & u_{12} & \dots & u_{1N} \\ u_{21} & u_{22} & \dots & u_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ u_{N1} & u_{N2} & \dots & u_{NN} \end{bmatrix} \begin{bmatrix} s_1 & 0 & 0 & \dots & 0 \\ 0 & s_2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 & \vdots \\ 0 & 0 & \dots & s_N & 0 \end{bmatrix} \begin{bmatrix} v_{11} & v_{21} & \dots & v_{m1} & \dots & v_{V1} \\ v_{12} & v_{22} & \dots & v_{m2} & \dots & v_{V2} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ v_{1m} & v_{2m} & \dots & v_{mn} & \dots & v_{Vm} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ v_{1V} & v_{2V} & \dots & v_{mV} & \dots & v_{VV} \end{bmatrix}$$

‘Strength’ of most significant correlation

$m = \frac{1}{N} \sum_{n=1}^N \text{vec}(d_i)$   
 ↑ average document vector  
 $W = W - m$   
 $W^T W$ : covariance matrix  
 $W = USV^T$

In PCA,  
 $C = VDV^T$   
 $D \propto S^2$   
 $C = W^T W = (USV^T)^T (USV^T)$   
 $= (VSU^T) \cdot USV^T$   
 $= V \cdot S^2 \cdot V^T$

# Interpretation of LSA

- The matrices  $U$  and  $V$  are **orthogonal matrices**
  - Their entries are real numbers
  - $U$  is  $N \times N$  ( $N$  is the number of documents) and  $V$  is  $V \times V$  ( $V$  is the vocabulary size)
  - They satisfy  $UU^T = I = U^TU$ ,  $VV^T = I = V^TV$
- The **singular values**  $s_1, \dots, s_N$  are positive and satisfy  $s_1 \geq s_2 \geq \dots \geq s_N$
- The off-diagonal entries of  $S$  are all zero

# Interpretation of LSA (continued)

- Focusing on  $V$ :
- The columns of  $V$   $\{v_1, \dots, v_V\}$  are  $V$  dimensional unit vectors orthogonal to each other
- They form a new **orthonormal basis** (coordinate system) for the document vector space
- Each column of  $V$  is a document vector *Super vector* corresponding to a **semantic class** (topic) in the corpus
- The importance of the topic  $v_n$  is indicated by the magnitude of the singular vector  $s_n$ .

# Interpretation of LSA (continued)

- Since  $v_n$  is a document vector, its  $j^{\text{th}}$  value corresponds to TF-IDF weight for  $j^{\text{th}}$  term in the vocabulary for the corresponding document/topic
- This can be used to interpret the topic corresponding to  $v_n$  – a large value of  $v_{nj}$  indicates that the  $j^{\text{th}}$  term in the vocabulary is significant for the topic.

# Interpretation of LSA (continued)

- Now consider  $U$
- It is easy to show that

$$Av_n = USV^T v_n = s_{nn} u_n$$

- While  $v_n$  describes the  $n^{\text{th}}$  topic as a combination of terms/words,  $u_n$  describes it as a combination of documents

# Topic-based representation

- Columns of  $V$ ,  $v_1, \dots, v_V$  are an **orthonormal basis** (coordinate system) for the document vector space
- If  $d$  is a document  $vec(d) \cdot v_n$  is the magnitude of the component of  $vec(d)$  in the direction of  $v_n$
- ..the component of  $vec(d)$  corresponding to topic  $n$

- Hence the vector  $top(d) = \begin{bmatrix} vec(d) \cdot v_1 \\ vec(d) \cdot v_2 \\ \vdots \\ vec(d) \cdot v_n \\ \vdots \\ vec(d) \cdot v_V \end{bmatrix}$

is a **topic-based** representation of  $d$  in terms of  $v_1, \dots, v_V$



# Topic-based dimension reduction

- Since the singular value  $s_n$  indicates the importance of topic  $v_n$ , we can choose to **truncate** the vector  $top(d)$  when  $s_n$  becomes small:

$$top(d) \approx \begin{bmatrix} vec(d) \cdot v_1 \\ vec(d) \cdot v_2 \\ \vdots \\ vec(d) \cdot v_n \end{bmatrix} = V_{(n)}^T vec(d)$$

where  $V_{(n)}$  is the  $V \times n$  matrix comprising the first  $n$  columns of  $V$

- $top(d)$  is a **reduced ( $n$ ) dimensional vector representation** of document  $d$

# Topic-based word representation

- Suppose  $w$  is the  $i^{th}$  word/term in the vocabulary
- The **one-hot vector**  $h(w)$  is the vector:

$$h(w) = \begin{bmatrix} 0 \\ \vdots \\ \cdot \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ \vdots \\ 0 \end{bmatrix} \longleftarrow i^{\text{th}} \text{ entry}$$

- $h(w)$  is like the document vector for a document consisting of just the word  $w$

# Topic-based word representation

- We can use  $h(w)$  to obtain a vector  $top(w)$  that describes  $w$  in terms of the topics that it contributes to:

$$top(w) = \begin{bmatrix} h(w) \cdot v_1 \\ h(w) \cdot v_2 \\ \vdots \\ h(w) \cdot v_n \end{bmatrix} = V_{(n)}^T top(w) =$$

*Words in similar direction  
=> represent similar topics*

- where  $V_{(n)}$  is the  $V \times n$  matrix comprising the first  $n$  columns of  $V$

# Topic-based word representation

- Intuitively, if two words  $v$  and  $w$  are **synonyms** they will contribute in a similar way to similar topics and the vectors  $top(v)$  and  $top(w)$  will point in similar directions
- A vector like  $top(w)$  is sometimes referred to as a **word embedding**
- We will re-visit this idea later

# More information about LSA

- See:

Landauer, T.K. and Dumais, S.T., “A solution to Platos problem: The Latent Semantic Analysis theory of the acquisition, induction and representation of knowledge”, *Psychological Review* 104(2), 211-240 (1997)

# Topic based document analysis

- There are other approaches to identifying a set of topics that represent a collection of documents – **Latent Dirichlet Allocation (LDA)**

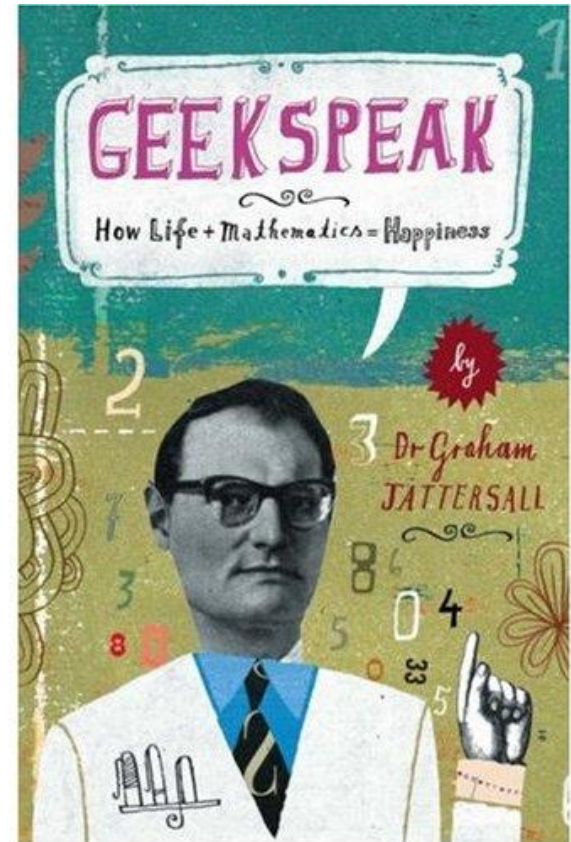
*topics ~ distribution of words*

# Thoughts on document vectors

- Once  $d$  is replaced by  $vec(d)$  it becomes a point in a vector space
- How does the structure of the vector space reflect the properties of the documents in it?
- Do clusters of vectors correspond to semantically related documents? ✓
- Can we partition the vector space into semantically different regions?
- These ideas are a link between IR and Data Mining

# For an alternative perspective...

- Chapter 14: “The cunning fox”
- Application of LSA to ‘dating agency’ personal adverts
- LSA suggests that the meaning of a personal advert can be expressed as a weighted combination of a few basic ‘concepts’



*Dr Graham Tattersall, “Geekspeak: How life + mathematics = happiness”, 2007*



# Relationship between LSA and PCA

- What is the relationship between LSA and PCA?

# Summary

- Latent Semantic Analysis
- Interpretation of LSA
- Topic-based representation of documents
- Topic-based dimension reduction

# Homework – this is challenging!

- Find the C program `doc2vec.c` on the course webCT page
- Use this to convert the BEng project specs from Laboratory 1 into a matrix of document vectors - this is the Word-Document matrix  $A$  from the notes
- Load this matrix into MATLAB and perform SVD
- Can you interpret the resulting Singular Vectors (columns of  $V$ )?