

Data Mining and Machine Learning

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Latent Semantic Analysis (LSA)

Add WeChat edu_assist_pro

Peter Jančovič

Objectives

- To understand, intuitively, how Latent Semantic Analysis (LSA) can discover latent topics in a corpus

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Vector Notation

- The vector representation $\text{vec}(d)$ of d is the V dimensional vector:

Assignment Project Exam Help

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

Add WeChat edu_assist_pro

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

$$w_{i(1),d} = f_{i(1),d} \times \text{IDF}(i(1)) \text{ from text IR}$$

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 $vec(d)$ will be tens of thousands
- Let's focus on just 2 of the dimensions, corresponding, say, to the words 'sea' and 'beach'
- Intuitively, often, when a document d includes 'sea' it will also include 'beach'

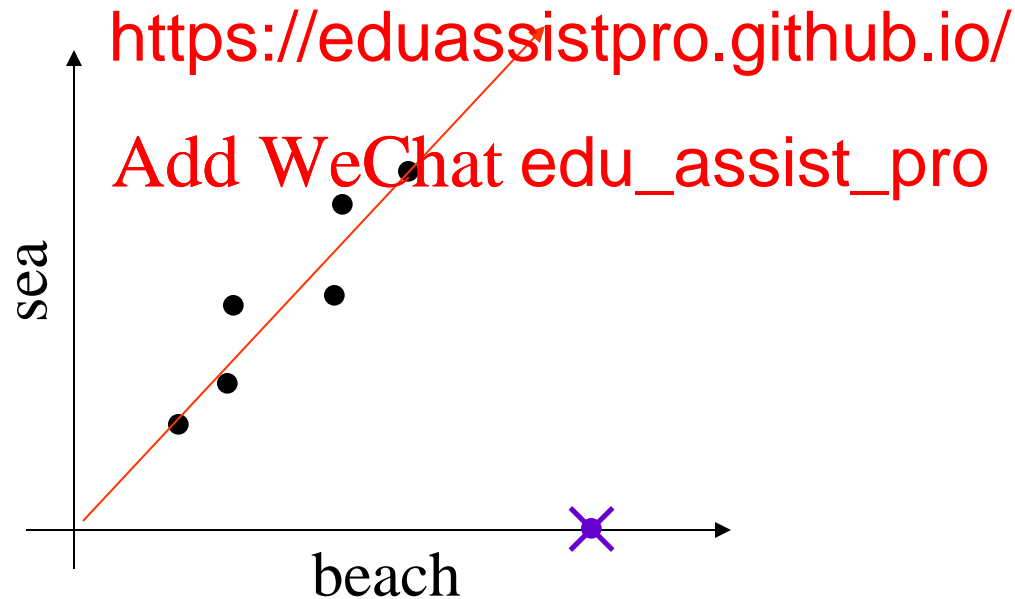
Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

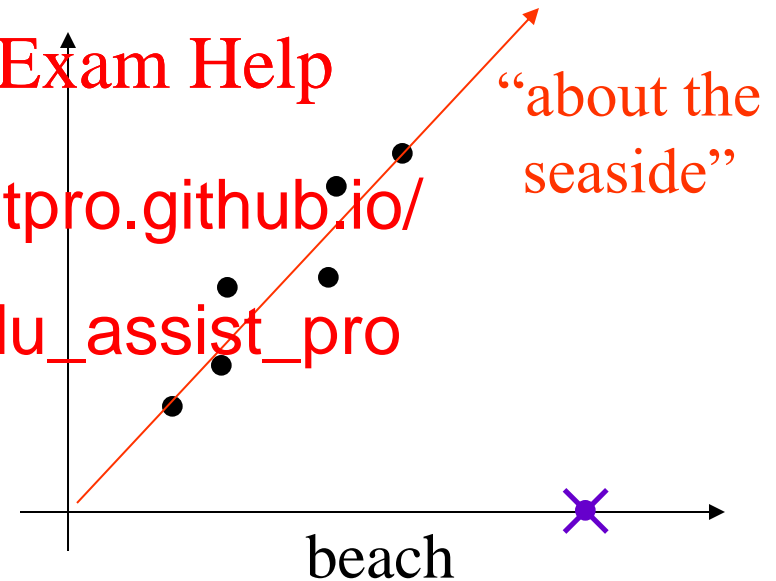
LSA continued

- 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
- In the example we have found an equivalence set of terms, including 'beach' and 'sea', which is 'about the seaside'



Finding Latent Semantic Classes

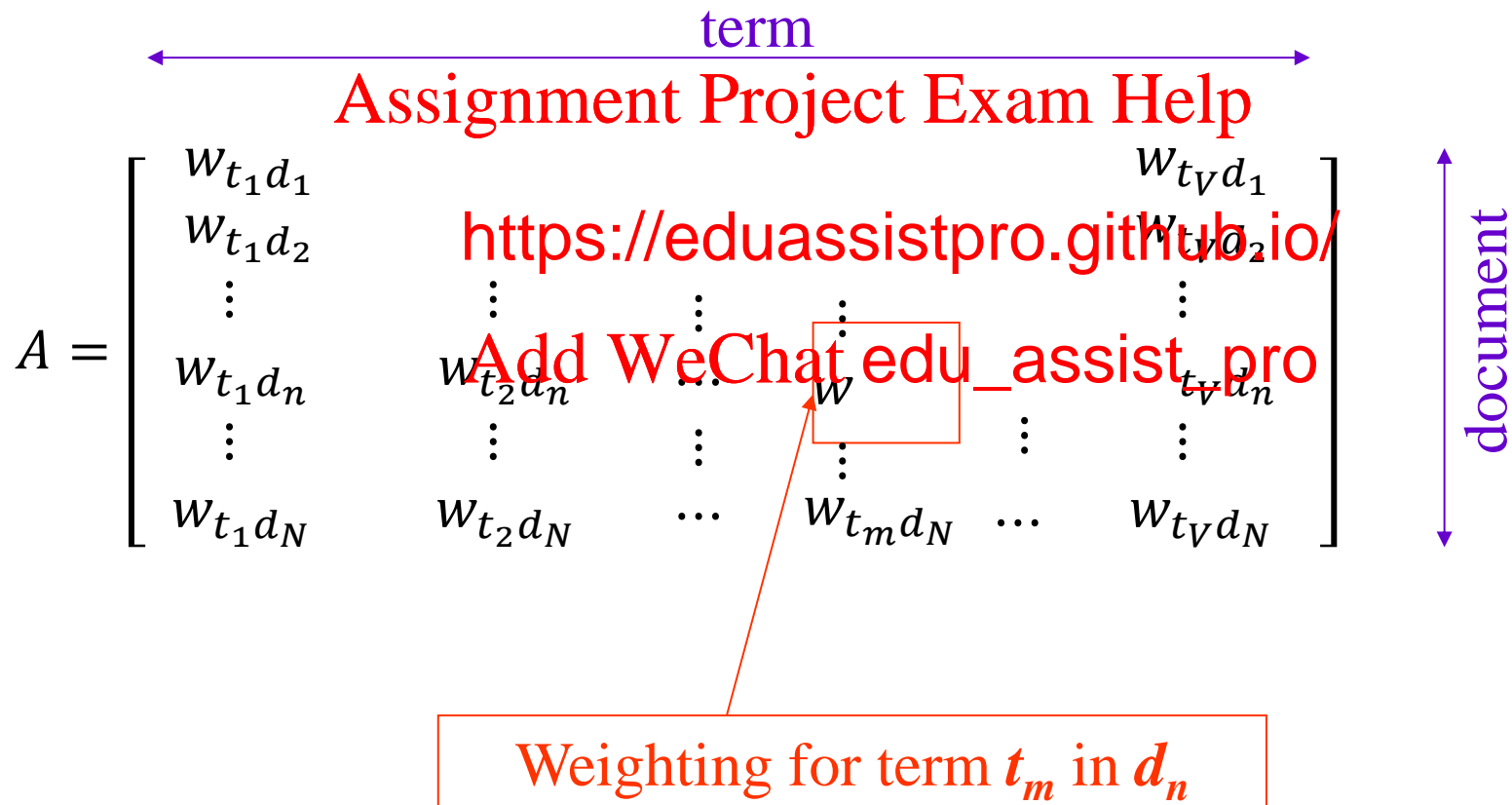
- LSA involves some advanced linear algebra – the description here is just an outline
- First construct the “word-document” matrix A
- Then decompose A using Singular Value Decomposition (SVD)
 - SVD is a standard technique in linear algebra
 - Packages such as MATLAB have SVD functions:
`>>[U,S,V]=svd(A)`

Singular Value Decomposition

- Remember eigenvector decomposition?
- An eigenvector of a square matrix A is a vector e such that $Ae = \lambda e$, where λ is a scalar
- For certain $A = UDU^T$, where U is an **orthogonal** (rotation) and D is **diagonal**
 - 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 a $N \times V$ matrix whose n^{th} row is $vec(d_n)$



Singular Value Decomposition (SVD)

N =number of docs, V =vocabulary size

$$A = USV^T$$

Direction of most significant correlation

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

$$A = \begin{bmatrix} u_{11} & u_{12} & \cdot & u_{1N} \\ u_{21} & u_{22} & \cdot & u_{2N} \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ u_{N1} & u_{N2} & \cdot & u_{NN} \end{bmatrix} \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & s_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_N \end{bmatrix} \begin{bmatrix} v_{11} & v_{21} & \cdot & v_{i1} & \cdot & v_{V1} \\ v_{12} & v_{22} & \cdot & v_{i2} & \cdot & v_{V2} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ v_{1V} & v_{2V} & \cdot & v_{iV} & \cdot & v_{VV} \end{bmatrix}$$

N (rows of U)
 V (columns of V)

‘Strength’ of most significant correlation

Interpretation of LSA

- The matrices U and V are orthogonal matrices
 - Their entries are real numbers
 - U is $N \times (\text{documents})$ and V is $V \times (\text{size})$
 - They satisfy $UU^T = I = V^TV$
- 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)

- Focussing on V :
 - The columns of V , $\{v_1, \dots, v_V\}$ are unit vectors and orthogonal to each other
 - They form a basis (coordinate system) for V space
 - Each column of V is a vector corresponding to a semantic class (topic) in the corpus
 - The importance of the topic corresponding to v_n is indicated by the size of the singular value s_n

Interpretation of LSA (continued)

- Since v_n is a document vector, its j^{th} value corresponds to TF-IDF weight for j^{th} term in the vocabulary for the corresponding document/topic
- This can be interpreted as the degree of association between the document v_n and the topic j^{th} term in the vocabulary is significant

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

Interpretation of LSA (continued)

- Now consider U
- It is easy to show that

$$Av_n = USVT^T v_n = \sum_n s_n u_n$$

- While v_n describes a combination of terms/words, u_n describes a combination of documents

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

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 $vec(d)$ in the direction of v_n

<https://eduassistpro.github.io/>

Add WeChat [edu_assist_pro](#)

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

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

More information about LSA

- See:

Landauer, T.K. and Dumais, S.T., “A solution to Platos problem: The Latent Semantic Analysis theory of the <https://eduassistpro.github.io/> and representation of knowled *review*
104(2), 211-240 (1997)

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? <https://eduassistpro.github.io/>
- Do clusters of vectors correspond to semantically related documents? [Add WeChat edu_assist_pro](#)
- 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 suggest <https://eduassistpro.github.io/> 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

Summary

- Latent Semantic Analysis

- Interpretatio

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro