

Information Retrieval

WS 2017 / **2018**

Lecture 10, Tuesday January 9th, 2018
(Latent Semantic Indexing)

Prof. Dr. Hannah Bast
Chair of Algorithms and Data Structures
Department of Computer Science
University of Freiburg

Overview of this lecture

■ Organizational

- Your experiences with ES9 Clustering

■ Contents

- Latent Semantic Indexing (LSI) Motivation + examples
- Singular Value Decomp. (SVD) Definition + how to compute
- Using LSI for retrieval Three variants
- **ES10: another pencil and paper sheet, for understanding and practicing the linear algebra (in particular, the SVD)**

In past years, we always had an implementation exercise for LSI. This was nice, because one could see the results. But most people coded just it without understanding how it really works.

■ Summary / excerpts

- Many of you liked this type of exercise ("math with meaning" + the kind of task you also get in the exam)
- Some of you don't like theory ... **there is no escape from it!**
- Several of you are unsure about how much to write and when it is enough ... **in case of doubt, ask on the forum!**
- Many of you were busy with your holidays and did not attend the last lecture
- Minor bug in the assumptions for Exercise 3, which was quickly discussed in the forum and fixed though
- "Cookie monster reminded me of missing out on the cookies in Lecture 7 due to watching the live stream"

■ Motivation

- Let's look at our example toy collection from L8 again:

D_1 and D_2 and D_3 are "about" surfing the web

D_5 and D_6 are "about" surfing on the beach

internet and web are **synonyms**, surfing is a **polysem**
= means different things in different context

	D_1	D_2	D_3	D_4	D_5	D_6
internet	1	1	0	1	0	0
web	1	0	1	1	0	0
surfing	1	1	1	2	1	1
beach	0	0	0	1	1	1

■ Motivation

- Let's look at the query **web surfing** again, using dot-product similarity as explained in L8
- Then $\text{sim}(D_3, Q) > \text{sim}(D_2, Q) = \text{sim}(D_5, Q)$

But D_2 seems just as relevant for the query as D_3 , only that the word "internet" is used instead of "web"

	REL REL REL (REL) NOT NOT						
	D ₁	D ₂	D ₃	D ₄	D ₅	D ₆	Q
internet	1	1	0	1	0	0	0
web	1	0	1	1	0	0	1
surfing	1	1	1	2	1	1	1
beach	0	0	0	1	1	1	0
	2	1	2	3	1	1	

Latent Semantic Indexing 3/8

■ Conceptual solution

	<small>REL</small> D₁	<small>REL</small> D₂	<small>REL</small> D₃	<small>(REL)</small> D₄	<small>NOT</small> D₅	<small>NOT</small> D₆	Q
internet	1	1	1	1	0	0	0
web	1	1	1	1	0	0	1
surfing	1	1	1	2	1	1	1
beach	0	0	0	1	1	1	0
	<small>2</small>	<small>2</small>	<small>2</small>	<small>3</small>	<small>1</small>	<small>1</small>	

Add the missing synonyms to the documents

Then indeed: $\text{sim}(D_1, Q) = \text{sim}(D_2, Q) = \text{sim}(D_3, Q)$

The goal of LSI is to do something like this **automagically**

Latent Semantic Indexing 4/8

■ A simple but powerful observation

	B_1 D_1	B_1 D_2	B_1 D_3	B_1+B_2 D_4	B_2 D_5	B_2 D_6	B_1	B_2
internet	1	1	1	1	0	0	1	0
web	1	1	1	1	0	0	1	0
surfing	1	1	1	2	1	1	1	1
beach	0	0	0	1	1	1	0	1

The modified matrix has **column rank 2**

That is, we can write each column as a (different) linear combination of the same two base columns (B_1 and B_2)

Note 1: the original matrix had column rank 4

Note 2: one can prove that **column rank = row rank**

Latent Semantic Indexing 5/8

■ Matrix factorization

The diagram shows the factorization of a 4x6 term-document matrix into a 4x2 matrix of base vectors and a 2x6 matrix of document vectors. Handwritten annotations include blue and red circles around specific elements and arrows pointing from the base vectors to the document vectors.

	D ₁	D ₂	D ₃	D ₄	D ₅	D ₆
1	1	1	1	1	0	0
2	1	1	1	1	0	0
3	1	1	1	2	1	1
4	0	0	0	1	1	1
	1	2	3	4	5	6

	B ₁	B ₂
1	1	0
2	1	0
3	1	1
4	0	1

=

	D' ₁	D' ₂	D' ₃	D' ₄	D' ₅	D' ₆
	1	1	1	1	0	0
	0	0	0	1	1	1
	1	2	3	4	5	6

Handwritten labels above the second matrix: B₁, B₁, B₁, B₁+B₂, B₂, B₂ with arrows pointing to the columns of the document vector matrix.

Equivalently: the 4 x 6 term-document matrix can be written as a product of a 4 x 2 matrix with a 2 x 6 matrix

The base vectors B_1 and B_2 are the underlying "**concepts**"

The vectors D'_1, \dots, D'_6 are the representation of the documents in the (lower-dimensional) "**concept space**"

■ The goal of LSI

- Given an $m \times n$ term-document matrix A and $k < \text{rank}(A)$
- Then find a matrix A' of (column) rank k such that the difference between A' and A is **as small as possible**
- Formally: $A' = \text{argmin}_{A' \text{ } m \times n \text{ with rank } k} \|A - A'\|$
- For the $\| \dots \|$ we take the so-called **Frobenius-norm**

This is defined as $\|D\| := \sqrt{\sum_{ij} D_{ij}^2}$

The reason for using this norm is purely technical: that way, the math on the next slides works out nicely

Latent Semantic Indexing 7/8

■ How to compute such an approximation A'

- We first compute the so-called **singular value decomposition (SVD)** of the given matrix A :

- Theorem: for any $m \times n$ matrix A of rank r , there exist U, S, V such that $A = U \cdot S \cdot V$, and where

U is an $m \times r$ matrix with $U^T \cdot U = I_r$ I_r = the $r \times r$ identity matrix

S is an $r \times r$ matrix with non-zero entries only on its diag.

V is an $r \times n$ matrix with $V \cdot V^T = I_r$

- The decomposition is **unique** up to simultaneous permutation of the rows/columns of U, S , and V

Standard form: diagonal entries of S positive and sorted

*the columns of U
are normalized to 1
and orthogonal
to each other*

*the rows of V
are L2-normalized
and orthogonal
to each other*

Latent Semantic Indexing 8/8

■ Using the SVD, our task becomes easy

- Let $A = U \cdot S \cdot V$ be the SVD of A
- For a given $k < \text{rank}(A)$ let

U_k = the first k columns of U now an $m \times k$ matrix

S_k = the upper $k \times k$ part of S now a $k \times k$ matrix

V_k = the first k rows of V now a $k \times n$ matrix

- Note that U_k is column-orthonormal just like U

That is, $U_k^T \cdot U_k = I_k \dots$ and similarly, $V_k \cdot V_k^T = I_k$

- Let $A_k = U_k \cdot S_k \cdot V_k$

Then A_k is a matrix of rank k that minimizes $\|A - A_k\|$

Computing the SVD 1/4

■ Overview

- Can be computed from the **Eigenvector decomposition**
See next slides for some of the mathematics behind
- In practice, the more direct **Lanczos** method is used
- This method has complexity **$O(k \cdot \text{nnz})$** , where **k** is the rank and **nnz** is the number of non-zero values in the matrix

Note that for term-document matrices $\text{nnz} \ll n \cdot m$

- When implementing this for large matrices in Python, one should use the method **`scipy.sparse.linalg.svds`**

Using **`numpy.linalg.svd`** on large matrices takes forever

iterative algorithm

Computing the SVD 2/4

$$\overbrace{\begin{pmatrix} 4 & 3 \\ 3 & 4 \end{pmatrix}}^B = \overbrace{\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}}^U \overbrace{\begin{pmatrix} 7 & 0 \\ 0 & 1 \end{pmatrix}}^{\text{diag}(\lambda_1, \lambda_2)} \overbrace{\frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & -1 \end{pmatrix}}^{U^T}$$

■ The Eigenvector Decomposition (EVD) $\begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \xrightarrow{\text{norm.}} \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$

- A real symmetric $n \times n$ matrix B has n pairwise orthogonal unit eigenvectors u_1, \dots, u_n (with eigenvalues $\lambda_1, \dots, \lambda_n$)

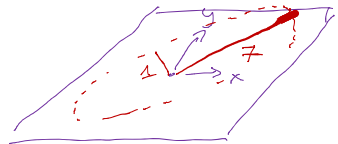
That is, $B \cdot u_i = \lambda_i \cdot u_i$ and $u_i \cdot u_j = 0$, for $i \neq j$, and $|u_i| = 1$

Equivalently, $B = U \cdot \text{diag}(\lambda_1, \dots, \lambda_n) \cdot U^T \dots$ and $U \cdot U^T = I$

- To compute all E-values λ , simply solve **$\det(B - \lambda \cdot I) = 0$**

- Example:

$$B = \begin{pmatrix} 4 & 3 \\ 3 & 4 \end{pmatrix}$$



One Eigenvector: $x_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ Eigenvalue $\lambda_1 = 7$

$$B \cdot x_1 = \begin{pmatrix} 4 & 3 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 4+3 \\ 3+4 \end{pmatrix} = \begin{pmatrix} 7 \\ 7 \end{pmatrix} = 7 \cdot x_1$$

Other Eigenvector: $x_2 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$ Eigenvalue $\lambda_2 = 1$

$$B \cdot x_2 = \begin{pmatrix} 4 & 3 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 1 \\ -1 \end{pmatrix} = \begin{pmatrix} 4-3 \\ 3-4 \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \end{pmatrix} = 1 \cdot x_2$$

λ is Eigenvalue $\Leftrightarrow \exists x \neq 0 \quad B \cdot x = \lambda \cdot x$ $n \times n$ identity matrix

$$\Leftrightarrow \exists x \neq 0 \quad (B - \lambda \cdot I) \cdot x = 0$$

$$\Leftrightarrow \det(B - \lambda \cdot I) = 0$$

Computing the SVD 3/4

■ From EVD to SVD, part 1

– Let \mathbf{A} be an arbitrary $\mathbf{m} \times \mathbf{n}$ matrix

– Then the matrix $\mathbf{A} \cdot \mathbf{A}^T$ is square and symmetric

$$\underbrace{(\mathbf{A} \cdot \mathbf{A}^T)^T}_{m \times m} = \underbrace{\mathbf{A}^T^T}_{m \times n} \cdot \underbrace{\mathbf{A}^T}_{n \times m} = \mathbf{A} \cdot \mathbf{A}^T$$

– Hence there exists an orthonormal matrix \mathbf{U} and a diagonal matrix \mathbf{S}_1 such that

$$\mathbf{A} \cdot \mathbf{A}^T = \mathbf{U} \cdot \mathbf{S}_1 \cdot \mathbf{U}^T$$

– Analogously, the matrix $\mathbf{A}^T \cdot \mathbf{A}$ is square and symmetric

– Hence there also exists an orthonormal matrix \mathbf{V} and a diagonal \mathbf{S}_2 such that

$$\mathbf{A}^T \cdot \mathbf{A} = \mathbf{V} \cdot \mathbf{S}_2 \cdot \mathbf{V}^T$$

Computing the SVD 4/4

■ From EVD to SVD, part 2

- Let us assume that a decomposition $\overset{m \times n}{\mathbf{A}} = \overset{m \times r}{\mathbf{U}} \cdot \overset{r \times r}{\mathbf{S}} \cdot \overset{r \times n}{\mathbf{V}}$ exists
- Then $\mathbf{A} \cdot \mathbf{A}^T = \mathbf{U} \cdot \mathbf{S}^2 \cdot \mathbf{U}^T$ and $\mathbf{A}^T \cdot \mathbf{A} = \mathbf{V}^T \cdot \mathbf{S}^2 \cdot \mathbf{V}$

That is: the **columns of \mathbf{U}** are the Eigenvectors of $\mathbf{A} \cdot \mathbf{A}^T$
and the **rows of \mathbf{V}** are the Eigenvectors of $\mathbf{A}^T \cdot \mathbf{A}$

And the singular values (diagonal values from \mathbf{S}) are just the square roots of the eigenvalues (which are the same for $\mathbf{A} \cdot \mathbf{A}^T$ and for $\mathbf{A}^T \cdot \mathbf{A}$)

$$\begin{aligned} \mathbf{A} \cdot \mathbf{A}^T &= (\mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V}) \cdot (\mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V})^T = \mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V} \cdot \mathbf{V}^T \cdot \overset{= \mathbf{I}_{r \times r}}{\mathbf{S}^T} \cdot \mathbf{U}^T = \mathbf{U} \cdot \mathbf{S}^2 \cdot \mathbf{U}^T \\ \mathbf{A}^T \cdot \mathbf{A} &= (\mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V})^T \cdot (\mathbf{U} \cdot \mathbf{S} \cdot \mathbf{V}) = \mathbf{V}^T \cdot \mathbf{S}^T \cdot \underbrace{\mathbf{U}^T \cdot \mathbf{U}}_{= \mathbf{I}_{r \times r}} \cdot \mathbf{S} \cdot \mathbf{V} = \mathbf{V}^T \cdot \mathbf{S}^2 \cdot \mathbf{V} \end{aligned}$$

Using LSI for better Retrieval 1/8

- **Variant 1:** work with A_k instead of A

	D_1	D_2	D_3	D_4	D_5	D_6
internet	1	1	0	1	0	0
web	1	0	1	1	0	0
surfing	1	1	1	2	1	1
beach	0	0	0	1	1	1

D'_1	D'_2	D'_3	D'_4	D'_5	D'_6
0.9	0.6	0.6	1.0	0.0	0.0
0.9	0.6	0.6	1.0	0.0	0.0
1.1	0.9	0.9	2.1	1.0	1.0
-0.1	0.1	0.1	0.9	1.0	1.0

Our example A from the beginning

best rank-2 approximation A_2

Using LSI for better Retrieval 2/8

■ **Variant 1:** work with \mathbf{A}_k instead of \mathbf{A}

- Problem: \mathbf{A}_k is a dense matrix, that is, most / all of it's $m \cdot n$ entries will be non-zero

Typically, both m and n will be very large, and then already storing this matrix is infeasible

Example: $m = 1000$ and $n = 10M \rightarrow m \cdot n = \mathbf{10\ G}$

Using LSI for better Retrieval 3/8

■ **Variant 2:** work with \mathbf{V}_k instead of with \mathbf{A}

- Recall: \mathbf{V}_k gives the representation of the documents in the k -dimensional concept space

	\mathbf{D}_1	\mathbf{D}_2	\mathbf{D}_3	\mathbf{D}_4	\mathbf{D}_5	\mathbf{D}_6
internet	1	1	0	1	0	0
web	1	0	1	1	0	0
surfing	1	1	1	2	1	1
beach	0	0	0	1	1	1

Our example \mathbf{A} from the beginning

\mathbf{D}'_1	\mathbf{D}'_2	\mathbf{D}'_3	\mathbf{D}'_4	\mathbf{D}'_5	\mathbf{D}'_6
0.4	0.3	0.3	0.7	0.3	0.3
0.5	0.2	0.2	0.0	-0.6	-0.6

\mathbf{V}_2 from the SVD of \mathbf{A}

- **Variant 2:** work with \mathbf{V}_k instead of with \mathbf{A}
 - Observation: V_k is a dense matrix, that is, most or all of its $k \cdot n$ entries are non-zero
 - Note: the original matrix A has $m' \cdot n$ non-zero entries, where m' is the average number of words in a document
 - So storing V_k instead of A is ok if $k \approx m'$ or smaller
 - No need for a sparse representation (inverted index)
 - This is a good variant to use in practice

Using LSI for better Retrieval 5/8

■ **Variant 2:** work with \mathbf{V}_k instead of with \mathbf{A}

- Problem 2: we need to map the query to concept space

The dot-product similarity of query \mathbf{q} with all documents is

$$\mathbf{q}^T \cdot \mathbf{A}_k = \mathbf{q}^T \cdot (\mathbf{U}_k \cdot \mathbf{S}_k \cdot \mathbf{V}_k) = \underbrace{(\mathbf{q}^T \cdot \mathbf{U}_k \cdot \mathbf{S}_k)}_{1 \times n} \cdot \underbrace{\mathbf{V}_k}_{n \times m}$$

$\mathbf{q}_k^T := \mathbf{q}^T \cdot \mathbf{U}_k \cdot \mathbf{S}_k$ is the query mapped to "concept space"

- The dot product $\mathbf{q}_k^T \cdot \mathbf{V}_k$ requires time $\sim \mathbf{n} \cdot \mathbf{k} \dots$ since both \mathbf{q}_k and \mathbf{V}_k are dense
- In comparison: computing the similarities of \mathbf{q} with the original documents requires time **$\mathbf{O}(\mathbf{n} \cdot \text{\#query words})$**

Using LSI for better Retrieval 6/8

■ **Variant 3:** expand the original documents

- In Variant 2, we have transformed both the query and the documents to concept space
- LSI can also be viewed as doing **document expansion** in the original space + no change in the query
- Namely, let $T_k = U_k \cdot U_k^T$ ^{$m \times q$ $q \times m$} this is an $m \times m$ matrix
- Then one can easily prove that $A_k = T_k \cdot A$ ^{$m \times m$ $m \times m$ $m \times m$}

do it as an (optional) exercise !

Using LSI for better Retrieval 7/8

■ Variant 3: expand the original documents

- Here is some intuition for \mathbf{T}_k , assuming **0** or **1** entries

In practice, we can get 0-1 entries by setting all entries in \mathbf{T} above a certain threshold to 1, and all others to 0

This 1 says: if there is "internet" in D_i then add "web" to D_i'

	internet	web	surfing	beach		D_i		D'_i	
internet	1	1	0	0	1	1	1	1	1
web	1	1	0	0	2	0	2	1	2
surfing	0	0	1	0	3	1	3	1	3
beach	0	0	0	1	4	0	4	0	4

\bullet

■ Linear combination with original scores

- Experience: LSI adds some useful information to the term-document matrix, but also a lot of **noise**
- In practice, one therefore uses a linear combination of the original scores and the LSI scores:

Variant 1: $\text{scores} = \lambda \cdot q^T \cdot A + (1 - \lambda) \cdot q^T \cdot A_k$

Variant 2: $\text{scores} = \lambda \cdot q^T \cdot A + (1 - \lambda) \cdot q_k^T \cdot V_k$

Variant 3: $\text{scores} = \lambda \cdot q^T \cdot A + (1 - \lambda) \cdot q^T \cdot T_k \cdot A$

References

■ Further reading

- Textbook Chapter 18: Matrix decompositions & LSI
<http://nlp.stanford.edu/IR-book/pdf/18lsi.pdf>
- Deerwester, Dumais, Landauer, Furnas, Harshman
[Indexing by Latent Semantic Analysis](#), JASIS 41(6), 1990

■ Web resources

- http://en.wikipedia.org/wiki/Latent_semantic_indexing
- http://en.wikipedia.org/wiki/Singular_value_decomposition
- <http://www.numpy.org/>
- <http://www.scipy.org/>