

# Latent Semantic Analysis

An overview of singular value decomposition  
in order to tackle semantics

Tim Van de Cruys

March 22, 2006

# Outline

- 1 Introduction
- 2 Preliminaries
  - Term-document matrix
  - Vector Space Measures
- 3 LSA: technique and example
  - Singular Value Decomposition
  - Example
- 4 Applications
  - Clustering
  - Varia

# Latent Semantic Analysis: Introduction 1/2

- Application of a mathematical/statistical technique to simulate how humans learn the semantics of a word
- LSA finds 'latent semantic dimensions' according to which words and documents can be identified
- Words (and passages) are represented as high-dimensional vectors in this semantic space
- Goal: counter data sparseness (poverty of the stimulus) and get rid of noise

# Latent Semantic Analysis: Introduction 1/2

- Application of a mathematical/statistical technique to simulate how humans learn the semantics of a word
- LSA finds 'latent semantic dimensions' according to which words and documents can be identified
- Words (and passages) are represented as high-dimensional vectors in this semantic space
- Goal: counter data sparseness (poverty of the stimulus) and get rid of noise

# Latent Semantic Analysis: Introduction 1/2

- Application of a mathematical/statistical technique to simulate how humans learn the semantics of a word
- LSA finds 'latent semantic dimensions' according to which words and documents can be identified
- Words (and passages) are represented as high-dimensional vectors in this semantic space
- Goal: counter data sparseness (poverty of the stimulus) and get rid of noise

# Latent Semantic Analysis: Introduction 1/2

- Application of a mathematical/statistical technique to simulate how humans learn the semantics of a word
- LSA finds 'latent semantic dimensions' according to which words and documents can be identified
- Words (and passages) are represented as high-dimensional vectors in this semantic space
- Goal: counter data sparseness (poverty of the stimulus) and get rid of noise

# Latent Semantic Analysis: Introduction 2/2

What is Latent Semantic Analysis technically speaking?

- The application a **singular value decomposition**
- to a **term-document matrix**
- to improve **vector space measures**

## Latent Semantic Analysis: Introduction 2/2

What is Latent Semantic Analysis technically speaking?

- The application a **singular value decomposition**
- to a **term-document matrix**
- to improve **vector space measures**



# Latent Semantic Analysis: Introduction 2/2

What is Latent Semantic Analysis technically speaking?

- The application a **singular value decomposition**
- to a **term-document matrix**
- to improve **vector space measures**

# Latent Semantic Analysis: Introduction 2/2

What is Latent Semantic Analysis technically speaking?

- The application a **singular value decomposition**
- to a **term-document matrix**
- to improve **vector space measures**

# Bag-of-word semantics

- LSA represents 'bag-of-word' semantics
- Idea that meaning of a passage equals the sum of the meaning of its words
- Meaning = an unordered set of word tokens, syntax is not taken into account
- Done by representing several passages in a **term-document matrix**

# Bag-of-word semantics

- LSA represents 'bag-of-word' semantics
- Idea that meaning of a passage equals the sum of the meaning of its words
- Meaning = an unordered set of word tokens, syntax is not taken into account
- Done by representing several passages in a **term-document matrix**

# Bag-of-word semantics

- LSA represents 'bag-of-word' semantics
- Idea that meaning of a passage equals the sum of the meaning of its words
- Meaning = an unordered set of word tokens, syntax is not taken into account
- Done by representing several passages in a **term-document matrix**

# Bag-of-word semantics

- LSA represents 'bag-of-word' semantics
- Idea that meaning of a passage equals the sum of the meaning of its words
- Meaning = an unordered set of word tokens, syntax is not taken into account
- Done by representing several passages in a **term-document matrix**

# Term-document matrix 1/2

Consider two documents:

- België is een koninkrijk in het midden van Europa, met als hoofdstad **Brussel**. **Brussel** heeft een Nederlandstalige en een Franstalige universiteit, maar de grootste studentenstad is **Leuven**. **Leuven** telt 27.000 studenten.
- Nederland is een West-Europees land aan de Noordzee. De hoofdstad van Nederland is **Amsterdam**. **Amsterdam** telt twee universiteiten. **Groningen** is een belangrijke studentenstad. In **Groningen** studeren 37.000 studenten.

# Term-document matrix 1/2

Consider two documents:

- België is een koninkrijk in het midden van Europa, met als hoofdstad **Brussel**. **Brussel** heeft een Nederlandstalige en een Franstalige universiteit, maar de grootste studentenstad is **Leuven**. **Leuven** telt 27.000 studenten.
- Nederland is een West-Europees land aan de Noordzee. De hoofdstad van Nederland is **Amsterdam**. **Amsterdam** telt twee universiteiten. **Groningen** is een belangrijke studentenstad. In **Groningen** studeren 37.000 studenten.



# Term-document matrix 1/2

Consider two documents:

- België is een koninkrijk in het midden van Europa, met als hoofdstad **Brussel**. **Brussel** heeft een Nederlandstalige en een Franstalige universiteit, maar de grootste studentenstad is **Leuven**. **Leuven** telt 27.000 studenten.
- Nederland is een West-Europees land aan de Noordzee. De hoofdstad van Nederland is **Amsterdam**. **Amsterdam** telt twee universiteiten. **Groningen** is een belangrijke studentenstad. In **Groningen** studeren 37.000 studenten.

# Term-document matrix 2/2

$$\begin{bmatrix} & B & NL \\ \textit{Groningen} & 0 & 2 \\ \textit{Leuven} & 2 & 0 \\ \textit{Amsterdam} & 0 & 2 \\ \textit{Brussel} & 2 & 0 \end{bmatrix}$$

Note: Values for words are transformed to values that represent their **importance** in the passage (entropy, mutual information)

## Matrix

	<i>B</i>	<i>NL</i>
<i>Groningen</i>	0	2
<i>Leuven</i>	2	0
<i>Amsterdam</i>	0	2
<i>Brussel</i>	2	0

## Apply cosine similarity measure

- $$\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{\|\vec{x}\| \|\vec{y}\|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$$

- Examples:

$$\cos(\text{Leuven}, \text{Amsterdam}) = \frac{0 \cdot 0 + 2 \cdot 2}{\sqrt{2^2} \sqrt{2^2}} = \frac{4}{2 \cdot 2} = 1$$

$$\cos(\text{Groningen}, \text{Brussel}) = \frac{0 \cdot 2 + 2 \cdot 0}{\sqrt{0^2} \sqrt{2^2}} = \frac{0}{0 \cdot 2} = 0$$

## Matrix

	<i>B</i>	<i>NL</i>
<i>Groningen</i>	0	2
<i>Leuven</i>	2	0
<i>Amsterdam</i>	0	2
<i>Brussel</i>	2	0

## Apply cosine similarity measure

- $\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$

- Examples:

- $\cos(\textit{Groningen}, \textit{Amsterdam}) = \frac{4}{\sqrt{16}} = 1$

- $\cos(\textit{Groningen}, \textit{Brussel}) = \frac{0}{\sqrt{16}} = 0$

## Matrix

	<i>B</i>	<i>NL</i>
<i>Groningen</i>	0	2
<i>Leuven</i>	2	0
<i>Amsterdam</i>	0	2
<i>Brussel</i>	2	0

## Apply cosine similarity measure

- $\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$
- Examples:
  - $\cos(\text{Groningen}, \text{Amsterdam}) = \frac{4}{\sqrt{16}} = 1$
  - $\cos(\text{Groningen}, \text{Brussel}) = \frac{0}{\sqrt{16}} = 0$

## Matrix

	<i>B</i>	<i>NL</i>
<i>Groningen</i>	0	2
<i>Leuven</i>	2	0
<i>Amsterdam</i>	0	2
<i>Brussel</i>	2	0

## Apply cosine similarity measure

- $\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$
- Examples:
  - $\cos(\text{Groningen}, \text{Amsterdam}) = \frac{4}{\sqrt{16}} = 1$
  - $\cos(\text{Groningen}, \text{Brussel}) = \frac{0}{\sqrt{16}} = 0$

## Matrix

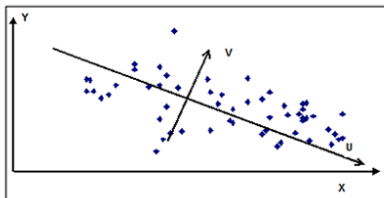
	<i>B</i>	<i>NL</i>
<i>Groningen</i>	0	2
<i>Leuven</i>	2	0
<i>Amsterdam</i>	0	2
<i>Brussel</i>	2	0

## Apply cosine similarity measure

- $\cos(\vec{x}, \vec{y}) = \frac{\vec{x} \cdot \vec{y}}{|\vec{x}| |\vec{y}|} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}}$
- Examples:
  - $\cos(\text{Groningen}, \text{Amsterdam}) = \frac{4}{\sqrt{16}} = 1$
  - $\cos(\text{Groningen}, \text{Brussel}) = \frac{0}{\sqrt{16}} = 0$

# Singular Value Decomposition

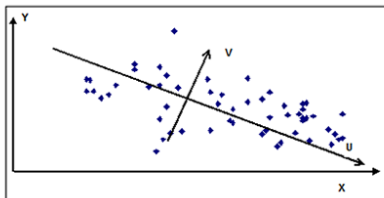
- Mathematical/statistical technique closely related to principal component analysis, factor analysis
- Find the dimensions that explain most variability by finding **eigenvectors** of matrix
- Only keep the  $n$  most important dimensions ( $n=50-1000$ )





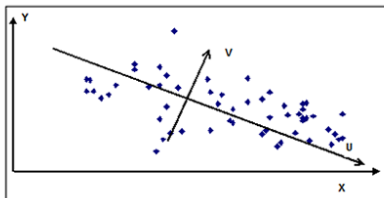
# Singular Value Decomposition

- Mathematical/statistical technique closely related to principal component analysis, factor analysis
- Find the dimensions that explain most variability by finding **eigenvectors** of matrix
- Only keep the  $n$  most important dimensions ( $n=50-1000$ )



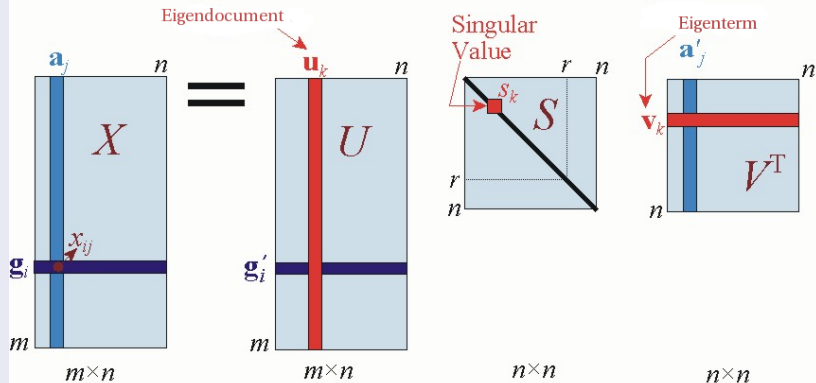
# Singular Value Decomposition

- Mathematical/statistical technique closely related to principal component analysis, factor analysis
- Find the dimensions that explain most variability by finding **eigenvectors** of matrix
- Only keep the  $n$  most important dimensions ( $n=50-1000$ )



## SVD: three matrices

$$X = USV^T$$



## Example 1a

$$A \begin{bmatrix} & B & NL \\ \textit{Groningen} & 0 & 2 \\ \textit{Leuven} & 2 & 0 \\ \textit{Amsterdam} & 0 & 2 \\ \textit{Brussel} & 2 & 0 \end{bmatrix} = U \begin{bmatrix} 0.00 & 0.71 \\ -0.71 & 0.00 \\ 0.00 & 0.71 \\ -0.71 & 0.00 \end{bmatrix} * S \begin{bmatrix} 2.83 & 0 \\ 0 & 2.83 \end{bmatrix} * V^T \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$$

## Example 1b

$$\begin{bmatrix} & B & NL & B \\ \text{Groningen} & 0 & 2 & 0 \\ \text{Leuven} & 2 & 0 & 0 \\ \text{Amsterdam} & 0 & 2 & 0 \\ \text{Brussel} & 2 & 0 & 1 \end{bmatrix} = U \begin{bmatrix} 0.00 & -0.71 & \cancel{0.00} \\ -0.66 & 0.00 & \cancel{0.75} \\ 0.00 & -0.71 & \cancel{0.00} \\ -0.75 & 0.00 & \cancel{\cancel{0.66}} \end{bmatrix} * S \begin{bmatrix} 2.92 & 0 & \emptyset \\ 0 & 2.83 & \emptyset \\ \emptyset & \emptyset & \cancel{0.68} \end{bmatrix} * V^T \begin{bmatrix} -0.97 & 0.00 & 0.26 \\ 0 & -1.00 & 0.00 \\ \cancel{\cancel{0.26}} & \cancel{0.00} & \cancel{\cancel{0.97}} \end{bmatrix} \cong A' \begin{bmatrix} 0.0 & 2.0 & 0.0 \\ 1.9 & 0.0 & 0.5 \\ 0.0 & 2.0 & 0.0 \\ 2.1 & 0 & 0.6 \end{bmatrix}$$

## Methodological remarks

- LSA in CLEF-corpus: 4 years of Dutch newspaper texts (Algemeen Dagblad, NRC Handelsblad)
- terms = nouns  
documents = articles
- 20.000 terms \* 100.000 documents matrix
- reduced to 300 latent semantic dimensions

## Methodological remarks

- LSA in CLEF-corpus: 4 years of Dutch newspaper texts (Algemeen Dagblad, NRC Handelsblad)
- terms = nouns  
documents = articles
- 20.000 terms \* 100.000 documents matrix
- reduced to 300 latent semantic dimensions

## Methodological remarks

- LSA in CLEF-corpus: 4 years of Dutch newspaper texts (Algemeen Dagblad, NRC Handelsblad)
- terms = nouns  
documents = articles
- 20.000 terms \* 100.000 documents matrix
- reduced to 300 latent semantic dimensions



## Methodological remarks

- LSA in CLEF-corpus: 4 years of Dutch newspaper texts (Algemeen Dagblad, NRC Handelsblad)
- terms = nouns  
documents = articles
- 20.000 terms \* 100.000 documents matrix
- reduced to 300 latent semantic dimensions

## Dimension 37

contingent_periode	0.55974	verslagweek	0.51471
NBC's	0.321241	courant	0.31842
week_staats	0.226316	dag_geld_rente	0.192096
kas_reserve	0.151812	belening	0.131637
belasting_afdracht	0.127845	geldmarkt_tarief	0.109959
verkrapping	0.10905	beleningen	0.0904901
geldmarkt	0.0747359	hoofde	0.0622487
DNB	0.0586182	storting	0.0567692
contingent	0.0542929	voorschot	0.0505641
omloop	0.0499995	mutatie	0.0472168
benutting	0.0324582	procentpunt	0.0323419
voorschot_rente	0.030187	schatkist	0.029555
ambtenaren_salaris	0.0289534	beroep_contingent	0.0286949
bankbiljet	0.0264982	basispunt	0.0243368
bankwezen	0.0237	liquiditeiten	0.0233299

## Dimension 64

23u	0.105716	Sneak	0.0904175
16u15	0.0882233	preview	0.0798345
On deadly ground	0.0636528	Cool runnings	0.0633079
22u30	0.0602807	17u30	0.0547169
Trois couleurs	0.0408388	Sister	0.0407604
ardilla	0.0395078	roja	0.0395078
The snapper	0.0364226	Mrs. Doubtfire	0.0335517
21u15	0.0328192	Monk	0.0325375
Intersection	0.024875	14u45	0.02288
spirits	0.022065	euro	0.021874
The piano	0.0212932	Aladdin	0.016921
Desmet	0.0166517	Ace Ventura	0.0162278
Mr. Jones	0.0149521	The three musketeers	0.0138651
La	0.0125177	rocker	0.01078
Philadelphia	0.0104788	15u30	0.0104311

## Clustering with LSA 1/2

- azijn bieslook bleekselderij blokje bosuitje boter bouillon champignon citroen citroensap crème deciliter deeg deksel dressing eetlepel folie garnaal gehakt gram ham hoofdgerecht keukenpapier knoflook koekepan komkommer kook lepel mossel mosterd olijf olijfolie oven pan pannetje paprika peper peterselie plak plakje prei reepje room salade saus slagroom spinazie takje theelepel Tip tomaat ui vocht voor\_gerecht vrucht vlees vulling warmtebron zeef zout
- artillerie Banja Luka Belgrado Bihac directrice enclave ex-Joegoslavië her\_opening Knin Krajina Kroaat Kroatië luchtmachtbasis massagraf Milosevic Montenegro Oost-Slavonië Radovan Karadzic Servië Serviër Tudjman UNPROFOR VN-soldaat vredesplan Zagreb

## Clustering with LSA 1/2

- azijn bieslook bleekselderij blokje bosuitje boter bouillon champignon citroen citroensap crème deciliter deeg deksel dressing eetlepel folie garnaal gehakt gram ham hoofdgerecht keukenpapier knoflook koekepan komkommer kook lepel mossel mosterd olijf olijfolie oven pan pannetje paprika peper peterselie plak plakje prei reepje room salade saus slagroom spinazie takje theelepel Tip tomaat ui vocht voor\_gerecht vrucht vlees vulling warmtebron zeef zout
- artillerie Banja Luka Belgrado Bihac directrice enclave ex-Joegoslavië her\_opening Knin Krajina Kroaat Kroatië luchtmachtbasis massagraf Milosevic Montenegro Oost-Slavonië Radovan Karadzic Servië Serviër Tudjman UNPROFOR VN-soldaat vredesplan Zagreb

## Clustering with LSA 2/2

- aantasting afbraak broeikaseffect CO dikte halfmond  
katalysator KNMI kooldioxyde kristal meting Montreal Nature  
ocean ozon\_laag stabilisatie straling verbranding vulkaan  
waarneming zonlicht zuurstof

## Clustering with syntactic relations (distributional similarity)

- **bieslook bosuitje citroensap deciliter eetlepel eierdooier  
gember gram kaneel knoflook koriander mosterd peper  
peterselie theelepel tijm**
- bak container doos kist koffer kom pan pot schaal zak
- Albanië Armenië Georgië Kazachstan Kroatië Macedonië  
Oekraïne Oezbekistan Wit-Rusland
- Hutu jood Kroaat moslim Palestijn Serviër

## Clustering with syntactic relations (distributional similarity)

- bieslook bosuitje citroensap deciliter eetlepel eierdooier  
gember gram kaneel knoflook koriander mosterd peper  
peterselie theelepel tijm
- bak container doos kist koffer kom pan pot schaal zak
- Albanië Armenië Georgië Kazachstan Kroatië Macedonië  
Oekraïne Oezbekistan Wit-Rusland
- Hutu jood Kroaat moslim Palestijn Serviër



## Clustering with syntactic relations (distributional similarity)

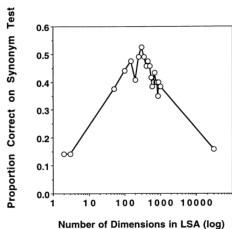
- bieslook bosuitje citroensap deciliter eetlepel eierdooier  
gember gram kaneel knoflook koriander mosterd peper  
peterselie theelepel tijm
- bak container doos kist koffer kom pan pot schaal zak
- Albanië Armenië Georgië Kazachstan Kroatië Macedonië  
Oekraïne Oezbekistan Wit-Rusland
- Hutu jood Kroaat moslim Palestijn Serviër

## Clustering with syntactic relations (distributional similarity)

- bieslook bosuitje citroensap deciliter eetlepel eierdooier  
gember gram kaneel knoflook koriander mosterd peper  
peterselie theelepel tijm
- bak container doos kist koffer kom pan pot schaal zak
- Albanië Armenië Georgië Kazachstan Kroatië Macedonië  
Oekraïne Oezbekistan Wit-Rusland
- Hutu jood Kroaat moslim Palestijn Serviër

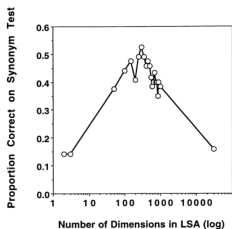
# LSA & synonym tests

- LSA trained on Grolier Encyclopedia, and given synonym test (TOEFL test).
- LSA scores 65%, identical to the average score of a large sample of students applying for college entrance in the United States from non-English speaking countries.



# LSA & synonym tests

- LSA trained on Grolier Encyclopedia, and given synonym test (TOEFL test).
- LSA scores 65%, identical to the average score of a large sample of students applying for college entrance in the United States from non-English speaking countries.



# LSA & essay grading

- LSA was used to assign scores to essay questions
- Trained on instructional text read by students, and a few example essays (or pre-graded essays)
- LSA scores found to be as good as expert assigned scores
- Moreover: LSA correlated significantly better with individual expert graders than one expert correlated with another
- LSA grades essays as good as expert **and** more objective

## LSA & essay grading

- LSA was used to assign scores to essay questions
- Trained on instructional text read by students, and a few example essays (or pre-graded essays)
- LSA scores found to be as good as expert assigned scores
- Moreover: LSA correlated significantly better with individual expert graders than one expert correlated with another
- LSA grades essays as good as expert **and** more objective

# LSA & essay grading

- LSA was used to assign scores to essay questions
- Trained on instructional text read by students, and a few example essays (or pre-graded essays)
- LSA scores found to be as good as expert assigned scores
- Moreover: LSA correlated significantly better with individual expert graders than one expert correlated with another
- LSA grades essays as good as expert **and** more objective

# LSA & essay grading

- LSA was used to assign scores to essay questions
- Trained on instructional text read by students, and a few example essays (or pre-graded essays)
- LSA scores found to be as good as expert assigned scores
- Moreover: LSA correlated significantly better with individual expert graders than one expert correlated with another
- LSA grades essays as good as expert **and** more objective



# LSA & essay grading

- LSA was used to assign scores to essay questions
- Trained on instructional text read by students, and a few example essays (or pre-graded essays)
- LSA scores found to be as good as expert assigned scores
- Moreover: LSA correlated significantly better with individual expert graders than one expert correlated with another
- LSA grades essays as good as expert **and** more objective