

Dimensionality Reduction Models for Language Data

Overview and Applications

Tim Van de Cruys

University of Groningen

Statistics Seminar
February 25th, 2009



Matrices

- Capturing co-occurrence frequencies of two entities



Matrices

- Capturing co-occurrence frequencies of two entities

	doc1	doc2	doc3	doc4
word1				
word2				
word3				
word4				



Matrices

- Capturing co-occurrence frequencies of two entities

	word1	word2	word3	word4
word1				
word2				
word3				
word4				



Matrices

- Capturing co-occurrence frequencies of two entities

	dep1	dep2	dep3	dep4
word1				
word2				
word3				
word4				



Matrices

- Capturing co-occurrence frequencies of two entities

	rood	lekker	snel	tweedehands
appel	2	1	0	0
wijn	2	2	0	0
auto	1	0	1	2
vrachtwagen	1	0	1	1



Matrices

- Capturing co-occurrence frequencies of two entities

	rood	lekker	snel	tweedehands
appel	7	9	0	0
wijn	12	6	0	0
auto	7	0	8	4
vrachtwagen	2	0	3	4



Matrices

- Capturing co-occurrence frequencies of two entities

	rood	lekker	snel	tweedehands
appel	56	98	0	0
wijn	44	34	0	0
auto	23	0	31	39
vrachtwagen	4	0	18	29



Matrices

- Capturing co-occurrence frequencies of two entities

	rood	lekker	snel	tweedehands
appel	728	592	1	0
wijn	1035	437	0	2
auto	392	0	487	370
vrachtwagen	104	0	393	293



Introduction

Two reasons for performing dimensionality reduction:

- Intractable computations
 - When number of elements and number of features is too large, similarity computations may become intractable
 - reduction of the number of features makes computation tractable again
- Generalization capacity
 - the dimensionality reduction is able to describe the data better, or is able to capture intrinsic semantic features
 - dimensionality reduction is able to improve the results (counter data sparseness and noise)



Latent Semantic Analysis: Introduction 1/2

- Application of a mathematical/statistical technique to simulate how humans learn the semantics of words
- 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 of **singular value decomposition**
- to a **term-document matrix**
- to improve **similarity calculations**



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 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

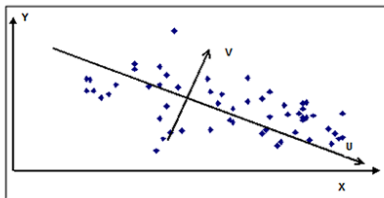
Apply cosine similarity measure

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



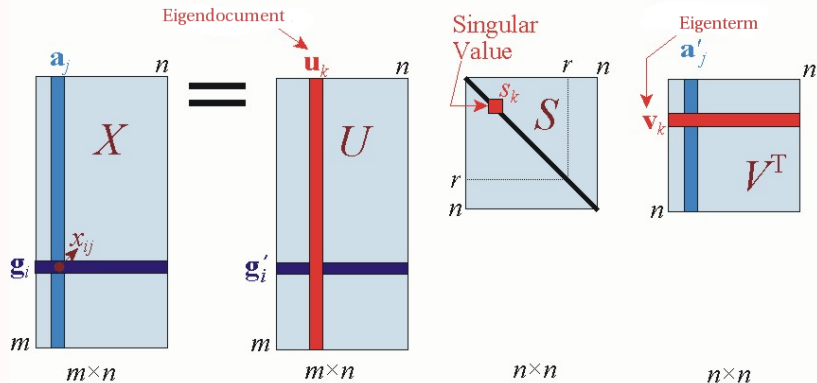
Singular Value Decomposition

- Mathematical/statistical technique closely related to principal component analysis
- SVD: generalization of PCA, computations on the actual data
- rectangular matrix instead of square covariance matrix
- Find the dimensions that explain most variance by solving number of **eigenvector** problems
- Only keep the n most important dimensions ($n = 50 - 300$)



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

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



Methodological remarks

- LSA in (part of) Twente Nieuws Corpus: 10 years of Dutch newspaper texts (AD, NRC, TR, VK, PAR)
- terms = nouns
documents = paragraphs
- 20,000 terms * 2,000,000 documents matrix
- reduced to 300 dimensions



Dimension 4

politiek	0.23	Amerikaans	0.08
land	0.22	kabinet	0.08
minister	0.20	Europa	0.08
partij	0.17	verkiezing	0.08
president	0.15	Nederlands	0.08
regering	0.14	internationaal	0.07
Europees	0.12	steun	0.07
Nederland	0.12	leider	0.07
premier	0.11	economisch	0.07
militair	0.09	VVD	0.07



Dimension 38

maand	0.22	cd	0.09
muziek	0.19	wereld	0.09
politie	0.15	geld	0.08
publiek	0.15	programma	0.08
school	0.14	groei	0.08
hoor	0.14	zing	0.08
volg	0.14	lied	0.08
minister	0.12	miljoen	0.08
lijk	0.11	win	0.08
nummer	0.11	val	0.08



Clustering with LSA 1/2

- azijn bak bestrijk bestrooi boter bouillon bout deeg dek_af
deksel eierdooier folie gaar garde giet goudbruin grill hak
keukenpapier knapperig knoflook koekepan koel_af koelkast
kook korrel lepel meng mengsel oven pan part pel pit rijst roer
roer_door roerend room rooster royaal saus schep schep_om
schik schil serveer smelt smoor snijd snipper strooi sudder
vergiet verhit verwarm voeg_toe voorverwarmd vork
vrucht vlees vuur wok wrijf_in zacht zeef zout



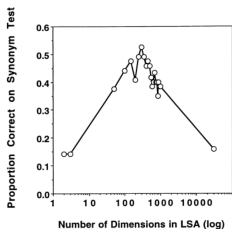
Clustering with LSA 2/2

- Breuker cellist cello Coleman compositie Concert concertzaal conservatorium contra_bas dj Duke Ensemble Es Gardiner house improvisatie instrument interval jazz Jazz klank klarinet luister Marsalis Monk musiceer musicus muziek_liefhebber muzikaal orgel pianist piano Productions recital Reinbert Rieu saxofonist saxofoon slagwerk slagwerker solist strijker symfonische trompet trompettist twintigste_eeuws violist viool Zappa zweep_op



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



Introduction

- LSA criticized for a number of reasons:
 - dimensionality reduction is best fit in least-square sense, but this is only valid for normally distributed data; language data is not normally distributed
 - Dimensions may contain negative values; it is not clear what negativity on a semantic scale should designate
- Shortcomings are remedied by subsequent techniques (PLSA, LDA, NMF, ...)



Technique 1/2

- Given a non-negative matrix V , find non-negative matrix factors W and H such that:

$$V_{n \times m} \approx W_{n \times r} H_{r \times m} \quad (1)$$

- Choosing $r \ll n, m$ reduces data
- Constraint on factorization: all values in three matrices need to be *non-negative values* (≥ 0)
- Constraint brings about a *parts-based* representation: only additive, no subtractive relations are allowed



Technique 2/2

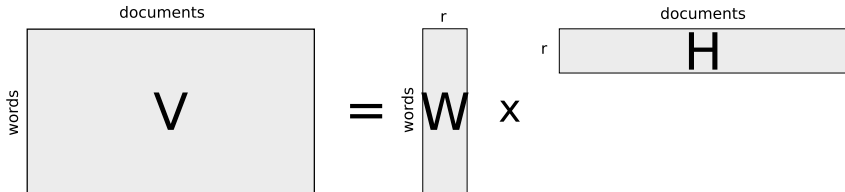
- Different kinds of NMF's that minimize different cost functions:
 - Square of Euclidean distance
 - Kullback-Leibler divergence
 - ⇒ better suited for language phenomena
- To find NMF is to minimize $D(V||WH)$ with respect to W and H , subject to the constraints $W, H \geq 0$
- This can be done with *update rules*

$$H_{a\mu} \leftarrow H_{a\mu} \frac{\sum_i W_{ia} \frac{V_{i\mu}}{(WH)_{i\mu}}}{\sum_k W_{ka}} \quad W_{ia} \leftarrow W_{ia} \frac{\sum_{\mu} H_{a\mu} \frac{V_{i\mu}}{(WH)_{i\mu}}}{\sum_v H_{av}} \quad (2)$$

- these update rules find a *local minimum* in the minimization of KL divergence



Graphical Representation



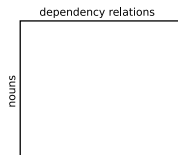
Results

- Context vectors (5k nouns \times 2k co-occurring nouns) extracted from CLEF corpus
- NMF is able to capture 'semantic' dimensions
- Examples:
 - *bus* 'bus', *taxi* 'taxi', *trein* 'train', *halte* 'stop', *reiziger* 'traveler', *perron* 'platform', *tram* 'tram', *station* 'station', *chauffeur* 'driver', *passagier* 'passenger'
 - *bouillon* 'broth', *slagroom* 'cream', *ui* 'onion', *eierdooier* 'egg yolk', *laurierblad* 'bay leaf', *zout* 'salt', *deciliter* 'decilitre', *boter* 'butter', *bleekselderij* 'celery', *saus* 'sauce'



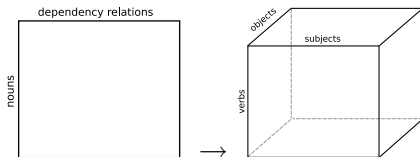
Two-way vs. three-way

- all methods use two way co-occurrence frequencies \longrightarrow matrix
- suitable for two-way problems
 - words \times documents
 - nouns \times dependency relations
- not suitable for n -way problems
 - words \times documents \times authors
 - verbs \times subjects \times direct objects



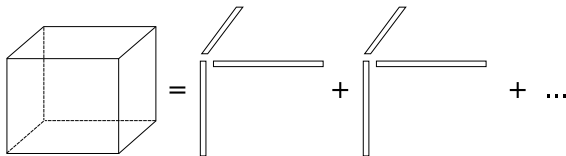
Two-way vs. three-way

- all methods use two way co-occurrence frequencies \longrightarrow matrix
- suitable for two-way problems
 - words \times documents
 - nouns \times dependency relations
- not suitable for n -way problems \longrightarrow tensor
 - words \times documents \times authors
 - verbs \times subjects \times direct objects

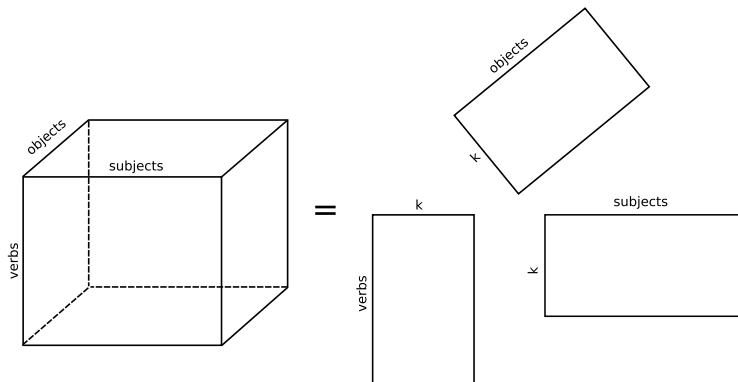


Technique

- Idea similar to non-negative matrix factorization
- Calculations are different
- $\min_{x_i \in \mathbb{R}_{\geq 0}^{D_1}, y_i \in \mathbb{R}_{\geq 0}^{D_2}, z_i \in \mathbb{R}_{\geq 0}^{D_3}} \| T - \sum_{i=1}^k x_i \circ y_i \circ z_i \|_F^2$



Graphical representation



Methodology

- Three-way extraction of selectional preferences
- Approach applied to Dutch, using TWENTE NIEUWS CORPUS (500M words)
- parsed with Alpino
- three-way co-occurrence of verbs with subjects and direct objects extracted
- adapted with extension of pointwise mutual information
- Resulting tensor 1K verbs \times 10K subjects \times 10K direct objects
- reduction to k dimensions ($k = 50, 100, 300$)



Example dimension: police action

subjects	su_s	verbs	v_s	objects	obj_s
<i>politie</i> 'police'	.99	<i>houd_aan</i> 'arrest'	.64	<i>verdachte</i> 'suspect'	.16
<i>agent</i> 'policeman'	.07	<i>arresteer</i> 'arrest'	.63	<i>man</i> 'man'	.16
<i>autoriteit</i> 'authority'	.05	<i>pak_op</i> 'run in'	.41	<i>betoger</i> 'demonstrator'	.14
<i>Justitie</i> 'Justice'	.05	<i>schiet_dood</i> 'shoot'	.08	<i>relschopper</i> 'rioter'	.13
<i>recherche</i> 'detective force'	.04	<i>verdenk</i> 'suspect'	.07	<i>raddraaier</i> 'instigator'	.13
<i>marechaussee</i> 'military police'	.04	<i>tref_aan</i> 'find'	.06	<i>overvaller</i> 'raider'	.13
<i>justitie</i> 'justice'	.04	<i>achterhaal</i> 'overtake'	.05	<i>Roemeen</i> 'Romanian'	.13
<i>arrestatieteam</i> 'special squad'	.03	<i>verwijder</i> 'remove'	.05	<i>actievoerder</i> 'campaigner'	.13
<i>leger</i> 'army'	.03	<i>zoek</i> 'search'	.04	<i>hooligan</i> 'hooligan'	.13
<i>douane</i> 'customs'	.02	<i>spoor_op</i> 'track'	.03	<i>Algerijn</i> 'Algerian'	.13



Example dimension: legislation

subjects	su_s	verbs	v_s	objects	obj_s
<i>meerderheid</i> 'majority'	.33	<i>steun</i> 'support'	.83	<i>motie</i> 'motion'	.63
<i>VVD</i>	.28	<i>dien_in</i> 'submit'	.44	<i>voorstel</i> 'proposal'	.53
<i>D66</i>	.25	<i>neem_aan</i> 'pass'	.23	<i>plan</i> 'plan'	.28
<i>Kamermeerderheid</i>	.25	<i>wijs_af</i> 'reject'	.17	<i>wetsvoorstel</i> 'bill'	.19
<i>fractie</i> 'party'	.24	<i>verwerp</i> 'reject'	.14	<i>hem</i> 'him'	.18
<i>PvdA</i>	.23	<i>vind</i> 'think'	.08	<i>kabinet</i> 'cabinet'	.16
<i>CDA</i>	.23	<i>aanvaard</i> 'accepts'	.05	<i>minister</i> 'minister'	.16
<i>Tweede Kamer</i>	.21	<i>behandel</i> 'treat'	.05	<i>beleid</i> 'policy'	.13
<i>partij</i> 'party'	.20	<i>doe</i> 'do'	.04	<i>kandidatuur</i> 'candidature'	.11
<i>Kamer</i> 'Chamber'	.20	<i>keur_goed</i> 'pass'	.03	<i>amendement</i> 'amendment'	.09



Example dimension: exhibition

subjects	su_s	verbs	v_s	objects	obj_s
<i>tentoonstelling</i> 'exhibition'	.50	<i>toon</i> 'display'	.72	<i>schilderij</i> 'painting'	.47
<i>expositie</i> 'exposition'	.49	<i>omvat</i> 'cover'	.63	<i>werk</i> 'work'	.46
<i>galerie</i> 'gallery'	.36	<i>bevat</i> 'contain'	.18	<i>tekening</i> 'drawing'	.36
<i>collectie</i> 'collection'	.29	<i>presenteer</i> 'present'	.17	<i>foto</i> 'picture'	.33
<i>museum</i> 'museum'	.27	<i>laat</i> 'let'	.07	<i>sculptuur</i> 'sculpture'	.25
<i>oeuvre</i> 'oeuvre'	.22	<i>koop</i> 'buy'	.07	<i>aquarel</i> 'aquarelle'	.20
<i>Kunsthall</i>	.19	<i>bezit</i> 'own'	.06	<i>object</i> 'object'	.19
<i>kunstenaar</i> 'artist'	.15	<i>zie</i> 'see'	.05	<i>beeld</i> 'statue'	.12
<i>dat</i> 'that'	.12	<i>koop_aan</i> 'acquire'	.05	<i>overzicht</i> 'overview'	.12
<i>hij</i> 'he'	.10	<i>in huis heb</i> 'own'	.04	<i>portret</i> 'portrait'	.11

