

Explorations in 3D: Semantic Tensor Spaces

Tim Van de Cruys

INRIA, ROCQUENCOURT

Distributional Semantics workshop
June 23, 2010
Groningen

Distributional similarity

Distributional similarity models are able to infer (lexical) semantics from text:

- Semantically similar words (syntactic context, small window)
 - **Parijs** 'Paris': *Londen* 'London', *Berlijn* 'Berlin', *New York*, *Milaan* 'Milan', *Rome*, *Wenen* 'Vienna', *Moskou* 'Moscow', *Stockholm*, *Frankfurt*
 - **liefde** 'love': *vriendschap* 'friendship', *passie* 'passion', *verlangen* 'desire', *angst* 'fear', *verdriet* 'sadness', *emotie* 'emotion', *schoonheid* 'beauty', *geloof* 'faith'

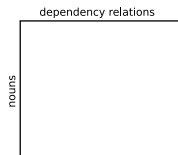
Distributional similarity

Distributional similarity models are able to infer (lexical) semantics from text:

- Semantic dimensions (large window)
 - (**transport**): *bus* 'bus', *taxi* 'taxi', *trein* 'train', *halte* 'stop', *reiziger* 'traveler', *perron* 'platform', *tram* 'tram', *station* 'station', *chauffeur* 'driver', *passagier* 'passenger'
 - (**food**): *bouillon* 'broth', *slagroom* 'cream', *ui* 'onion', *eierdooier* 'egg yolk', *laurierblad* 'bay leaf', *zout* 'salt', *deciliter* 'decilitre', *boter* 'butter', *bleekselderij* 'celery', *saus* 'sauce'

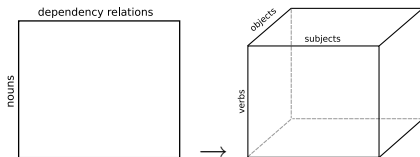
Two-way vs. three-way

- 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



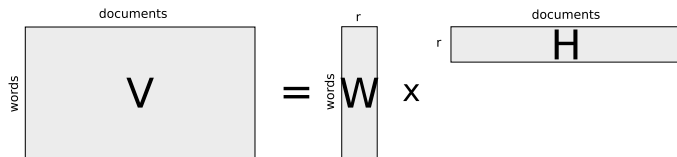
Technique

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

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

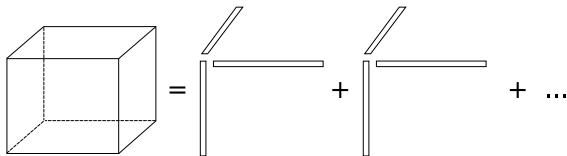
- Choosing $r \ll n, m$ reduces data
- Constraint on factorization: all values in three matrices need to be *non-negative values* (≥ 0)
- Constraint brings about a *parts-based* representation: only additive, no subtractive relations are allowed

Graphical Representation

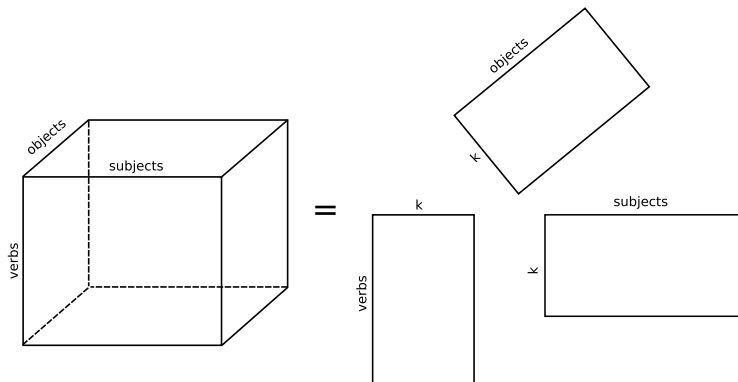


Technique

- 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



Introduction 1/2

- Standard selectional preference models: two-way co-occurrences
- Keeping track of single relationships
- But: two-way selectional preference models are not sufficiently rich
- Compare:
 - *The skyscraper is playing coffee.*
 - *The turntable is playing the piano.*

Introduction 2/2

- *The skyscraper is playing coffee.*
 - *(play, su, scyscraper)* ↓
 - *(play, obj, coffee)* ↓
- *The turntable is playing the piano.*
 - *(play, su, turntable)* ↑
 - *(play, obj, piano)* ↑
 - *(play, su, turntable, obj, piano)* ↓

Methodology

- Three-way extraction of selectional preferences
- Approach applied to Dutch, using TWENTE NIEUWS CORPUS (500M words of newspaper texts)
- parsed with Dutch dependency parser 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 ($k = 50, 100, 300$)

Examples: 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

Examples: 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

Examples: 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

Examples: quality count

- 44 dimensions contain similar, framelike semantics
- 43 dimensions contain less clear-cut semantics
 - single verbs account for one dimension
 - verb senses are mixed up
- 13 dimensions based on syntax rather than semantics
 - fixed expressions
 - pronomina

Evaluation: methodology

- pseudo-disambiguation task to test generalization capacity (standard automatic evaluation for selectional preferences)

<i>s</i>	<i>v</i>	<i>o</i>	<i>s'</i>	<i>o'</i>
<i>jongere</i>	<i>drink</i>	<i>bier</i>	<i>coalitie</i>	<i>aandeel</i>
'youngster'	'drink'	'beer'	'coalition'	'share'
<i>werkgever</i>	<i>riskeer</i>	<i>boete</i>	<i>doel</i>	<i>kopzorg</i>
'employer'	'risk'	'fine'	'goal'	'worry'
<i>directeur</i>	<i>zwaai</i>	<i>scepter</i>	<i>informatieur</i>	<i>vodka</i>
'manager'	'sway'	'sceptre'	'informer'	'wodka'

- 10-fold cross validation ($\pm 300,000$ co-occurrences)

Evaluation: models

- Evaluation of 4 different models
- 2 matrix models
 - $1\text{K verbs} \times (10\text{K subjects} + 10\text{K direct objects})$
 - singular value decomposition (\mathbb{R})
 - non-negative matrix factorization ($\mathbb{R}_{\geq 0}$)
- 2 tensor models
 - $1\text{K verbs} \times 10\text{K subjects} \times 10\text{K direct objects}$
 - parallel factor analysis (\mathbb{R})
 - non-negative tensor factorization ($\mathbb{R}_{\geq 0}$)

Evaluation: results

	dimensions		
	50 (%)	100 (%)	300 (%)
SVD	69.60 \pm 0.41	62.84 \pm 1.30	45.22 \pm 1.01
NMF	81.79 \pm 0.15	78.83 \pm 0.40	75.74 \pm 0.63
PARAFAC	85.57 \pm 0.25	83.58 \pm 0.59	80.12 \pm 0.76
NTF	89.52 \pm 0.18	90.43 \pm 0.14	90.89 \pm 0.16

Introduction

- **Problem:** ambiguity
- WAAL

Introduction

- **Problem:** ambiguity
- WAAL



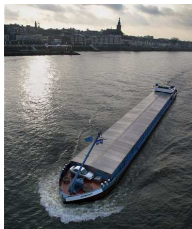
Introduction

- **Problem:** ambiguity
- WAAL



Introduction

- **Problem:** ambiguity
- WAAL

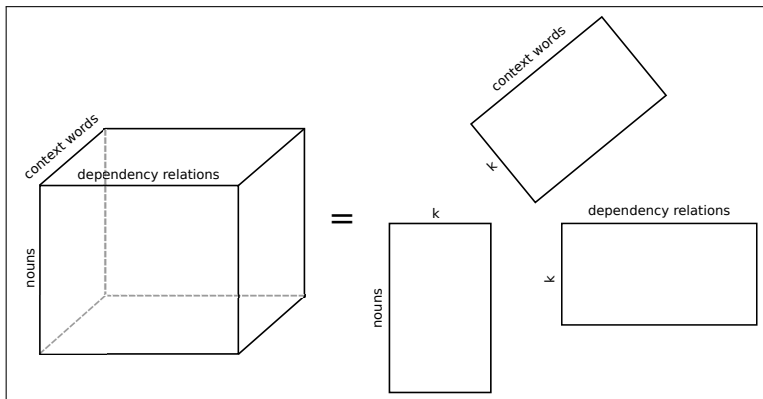


- Main research question: can 'bag of words' context and syntactic context be combined to differentiate between various senses of a word?

Methodology

- Goal: classification of nouns according to both 'bag of words' context and syntactic context
- \Rightarrow extract three-way co-occurrences of nouns with dependency relations with bag of words context words
- \Rightarrow three-way tensor of $5k$ nouns \times $80k$ dependency relations \times $1k$ context words
- Factorized with non-negative tensor factorization to 50 dimensions
- (For Dutch, TWNC, parsed with Alpino, paragraph as context)

Graphical Representation



Sense subtraction

- 'switch off' one dimension of an ambiguous word to reveal other possible senses
- From matrix W , we know which dimensions are the most important for a certain word
- Matrix H gives the importance of each dependency relation given a dimension
- 'subtract' dependency relations that are responsible for a given dimension from the original noun vector
 - $\vec{v}_{new} = \vec{v}_{orig}(\vec{1} - \vec{h}_{dim})$
 - each dependency relation is multiplied by a scaling factor, according to the load of the feature on the subtracted dimensions

Example dimension: food

nouns	dependency rels	context words
<i>knoflook</i> 'garlic'	<i>voeg_toe</i> _{OBJ} 'add'	<i>voeg_toe</i> 'add'
<i>peper</i> 'pepper'	<i>roer</i> _{OBJ} 'stir'	<i>minuut</i> 'minute'
<i>zout</i> 'salt'	<i>zout</i> _{COO} 'salt'	<i>vuur</i> 'fire'
<i>ui</i> 'oignon'	<i>peper</i> _{COO} 'pepper'	<i>water</i> 'water'
<i>peterselie</i> 'persil'	<i>knoflook</i> _{COO} 'garlic'	<i>stuk</i> 'piece'
<i>tomaat</i> 'tomato'	<i>meng</i> _{OBJ} 'mix'	<i>leg</i> 'lay'
<i>gember</i> 'ginger'	<i>ui</i> _{COO} 'oignon'	<i>rood</i> 'red'
<i>suiker</i> 'sugar'	<i>bak</i> _{OBJ} 'cook'	<i>vis</i> 'fish'
<i>citroen_sap</i> 'lemon juice'	<i>strooi</i> _{OBJ} 'sprinkle'	<i>dik</i> 'thick'
<i>koriander</i> 'coriander'	<i>laurier</i> _{COO} 'bay leaf'	<i>warm</i> 'warm'

Example dimension: music

nouns	dependency rels	context words
<i>muziek</i> 'music'	<i>klink</i> _{SUB} 'sound'	<i>muziek</i> 'music'
<i>dans</i> 'dance'	<i>klassiek</i> _{ADJ} 'classical'	<i>dans</i> 'dance'
<i>jazz</i>	<i>componeer</i> _{OBJ} 'compose'	<i>klassiek</i> 'classical'
<i>liedje</i> 'song'	<i>dans</i> _{COO} 'dance'	<i>cd</i>
<i>nummer</i> 'song'	<i>zing</i> _{OBJ} 'sing'	<i>lied</i> 'song'
<i>lied</i> 'song'	<i>hoor</i> _{OBJ} 'hear'	<i>klink</i> 'sound'
<i>cd</i>	<i>draai</i> _{OBJ} 'play'	<i>zing</i> 'sing'
<i>melodie</i> 'melody'	<i>theater</i> _{COO} 'theatre'	<i>theater</i> 'theatre'
<i>opera</i> 'opera'	<i>film</i> _{COO} 'film'	<i>nummer</i> 'song'
<i>rock</i>	<i>literatuur</i> _{COO} 'literature'	<i>tekst</i> 'text'

Conclusion

- Tensor space: novel model to investigate three-way (n -way) co-occurrence data
- Possible to generalize over co-occurrence data with appropriate factorization models
- Model applicable to and beneficial for:
 - Three-way selectional preference induction
 - Word sense discrimination