

Exploring Topic Models for Word Sense Discrimination

Tim Van de Cruys

University of Groningen

CLIN

December 7, 2007

Nijmegen



Semantic similarity

- Most work on semantic similarity relies on the DISTRIBUTIONAL HYPOTHESIS (Harris 1954)
- Take a word and its contexts:
 - tasty *smoutebol*
 - greasy *smoutebol*
 - a portion of *smoutebol*
 - *smoutebol* with loads of powdered sugar
- By looking at a word's context, one can infer its meaning

Semantic similarity

- Most work on semantic similarity relies on the DISTRIBUTIONAL HYPOTHESIS (Harris 1954)
- Take a word and its contexts:

- tasty *smoutebol*
- greasy *smoutebol*
- a portion of *smoutebol*
- *smoutebol* with loads of powdered sugar

⇒ **FOOD**

- By looking at a word's context, one can infer its meaning

Semantic similarity

- Most work on semantic similarity relies on the DISTRIBUTIONAL HYPOTHESIS (Harris 1954)
- Take a word and its contexts:

- tasty *smoutebol*
- greasy *smoutebol*
- a portion of *smoutebol*
- *smoutebol* with loads of powdered sugar



- By looking at a word's context, one can infer its meaning

Two kinds of context

- 1 'Bag of words' context
 - a window around the word is used as context
 - e.g. a fixed numbers of words, the paragraph in which a word appears, . . .
 - often used with some form of dimensionality reduction
- 2 Syntactic context
 - a corpus is parsed, dependency triples are extracted
 - e.g. <apple, obj, eat>, <apple, adj, red>
 - typically does not use any form of dimensionality reduction

Ambiguity

- **Problem:** ambiguity
 - Compare:
 - a trendy bar*
 - ↔ *an iron bar*
 - ↔ *today's air pressure: 1.013 bar*
 - Different meanings, but they are considered the same entity by a naive algorithm
- Main research question: can 'bag of words' context and syntactic context be combined to differentiate between various senses of a word?

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



Results

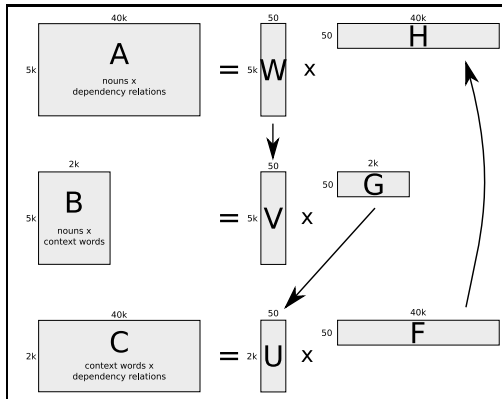
- Context vectors (5k nouns \times 2k co-occurring nouns) extracted from CLEF corpus
- NMF is able to capture 'semantic' dimensions)
- Examples:
 - *bus* 'bus', *taxi* 'taxi', *trein* 'train', *halte* 'stop', *reiziger* 'traveler', *perron* 'platform', *tram* 'tram', *station* 'station', *chauffeur* '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'



Methodology

- Goal: classification of nouns according to both 'bag of words' context and syntactic context
- \Rightarrow Construct three matrices capturing co-occurrence frequencies for each mode
 - nouns cross-classified by dependency relations
 - nouns cross-classified by (bag of words) context words
 - dependency relations cross-classified by context words
- \Rightarrow Apply NMF to matrices, but interleave the process
- Result of former factorization is used to initialize factorization of the next one

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

Combination with clustering

- A simple clustering algorithm (e.g. K-means) assigns ambiguous nouns to its predominant sense
- Centroid of the cluster is fold into topic model
- The dimensions that define the centroid are subtracted from the ambiguous noun vector
- Adapted noun vector is fed to the clustering algorithm again

Experimental Design

- Approach applied to Dutch, using CLEF corpus (Dutch newspaper texts '94-'95)
- Corpus parsed with Dutch dependency parser ALPINO
- three matrices constructed with:
 - 5k nouns \times 40k dependency relations
 - 5k nouns \times 2k context words
 - 40k dependency relations \times 2k context words
- Factorization to 100 dimensions

Example dimension: transport

- nouns:** *auto* 'car', *wagen* 'car', *tram* 'tram', *motor* 'motorbike', *bus* 'bus', *metro* 'subway', *automobilist* 'driver', *trein* 'train', *stuur* 'steering wheel', *chauffeur* 'driver'
- context words:** *auto* 'car', *trein* 'train', *motor* 'motorbike', *bus* 'bus', *rij* 'drive', *chauffeur* 'driver', *fiets* 'bike', *reiziger* 'reiziger', *passagier* 'passenger', *vervoer* 'transport'
- dependency relations:** *viertraps_{adj}* 'four pedal', *verplaats_{met_obj}* 'move with', *toeter_{adj}* 'honk', *tank_{in_houd_obj}* [parsing error], *tank_{subj}* 'refuel', *tank_{obj}* 'refuel', *rij_{voorbij_subj}* 'pass by', *rij_{voorbij_adj}* 'pass by', *rij_{af_subj}* 'drive off', *peperduur_{adj}* 'very expensive'



Pop: most similar words

pop music ↔ *doll*

- 1 *pop, rock, jazz, meubilair* 'furniture', *popmuziek* 'pop music', *heks* 'witch', *speelgoed* 'toy', *kast* 'cupboard', *servies* '[tea] service', *vraagteken* 'question mark'
- 2 *pop, meubilair* 'furniture', *speelgoed* 'toy', *kast* 'cupboard', *servies* '[tea] service', *heks* 'witch', *vraagteken* 'question mark', *sieraad* 'jewel', *sculptuur* 'sculpture', *schoen* 'shoe'
- 3 *pop, rock, jazz, popmuziek* 'pop music', *heks* 'witch', *danseres* 'dancer', *servies* '[tea] service', *kopje* 'cup', *house* 'house music', *aap* 'monkey'



Pop: clusters

- 1 *house, jazz, pop, rock*
- 2 *elektronica* 'electronics' , *horloge* 'watch' , *juweel* 'jewel' ,
kleding 'clothes' , *parfum* 'perfume' , *pop* 'doll' , *schoen* 'shoe' ,
sieraad 'jewel' , *speelgoed* 'toy' , *textiel* 'textile'



Barcelona: most similar words

Spanish city ↔ Spanish football club

- 1 *Barcelona, Arsenal, Inter, Juventus, Vitesse, Milaan 'Milan', Madrid, Parijs 'Paris', Wenen 'Vienna', München 'Munich'*
- 2 *Barcelona, Milaan 'Milan', München 'Munich', Wenen 'Vienna', Madrid, Parijs 'Paris', Bonn, Praag 'Prague', Berlijn 'Berlin', Londen 'London'*
- 3 *Barcelona, Arsenal, Inter, Juventus, Vitesse, Parma, Anderlecht, PSV, Feyenoord, Ajax*

Barcelona: clusters

- 1 *Anderlecht, Arsenal, Barcelona, Bayern, Inter, Juventus, Milan, Parma*
- 2 *Antwerpen 'Antwerp', Barcelona, Berlijn 'Berlin', Bonn, Bremen, Brussel 'Brussels', Frankfurt, Hamburg, Londen 'London', Milaan 'Milan', München 'Munich', Parijs 'Paris', Rome, Wenen 'Vienna'*



Conclusion

- Combining bag of words data and syntactic data is useful
 - bag of words data (factorized with NMF) puts its finger on topical dimensions
 - syntactic data is particularly good at finding similar words
 - a clustering approach allows one to determine which topical dimension(s) are responsible for a certain sense
 - and adapt the (syntactic) feature vector of the noun accordingly
 - subtracting the more dominant sense to discover less dominant senses



Future Work

- Proper evaluation, comparison with other WSD algorithms
- Extend the clustering part
 - Fully automatic clustering
 - Use of 'tight' clusters instead of 'broad' K-means results
- Scale the method to large data sets (20K nouns \times 100K dependency relations, larger corpus)

