

Lexical Association Measures

Collocation Extraction

Pavel Pecina

pecina@ufal.mff.cuni.cz

Institute of Formal and Applied Linguistics
Charles University, Prague



Doctoral thesis defense
Prague, September 24, 2008

Presentation outline

1/25

1. Introduction
2. Collocation extraction
3. Lexical association measures
4. Reference data
5. Empirical evaluation
6. Combining association measures
7. Conclusions

Lexical association

2/25

Collocational association

- ▶ restricts combination of words into phrases (beyond grammar)
- ▶ *collocations*
crystal clear, cosmetic surgery, weapons of mass destruction

Semantic association

- ▶ reflects semantic relationship between words
- ▶ *synonymy, antonymy, hypo/hyperonymy, meronymy, etc.*
sick – ill, baby – infant, dog – cat

Cross-language association

- ▶ corresponds to potential translations of words between languages
- ▶ *translation equivalents*
maison – house, baum – tree, květina – flower

Lexical association

2/25

Collocational association

- ▶ restricts combination of words into phrases (beyond grammar)
- ▶ *collocations*
crystal clear, cosmetic surgery, weapons of mass destruction

→ *lexicon*

Semantic association

- ▶ reflects semantic relationship between words
- ▶ *synonymy, antonymy, hypo/hyperonymy, meronymy, etc.*
sick – ill, baby – infant, dog – cat

→ *thesaurus*

Cross-language association

- ▶ corresponds to potential translations of words between languages
- ▶ *translation equivalents*
maison – house, baum – tree, květina – flower

→ *dictionary*

Measuring lexical association

3/25

Motivation

- ▶ **automatic** acquisition of associated words (*into a lexicon/thesarus/dictionary*)

Tool: Lexical association measures

- ▶ mathematical formulas determining **strength of association** between two (or more) words based on their occurrences and cooccurrences in a **corpus**.

Applications

- ▶ lexicography, natural language generation, word sense disambiguation
- ▶ bilingual word alignment, identification of translation equivalents
- ▶ information retrieval, cross-lingual information retrieval
- ▶ keyword extraction, named entity recognition
- ▶ syntactic constituent boundary detection, dependency parsing
- ▶ **collocation extraction**

Goals, objectives, and limitations

4/25

Goal

- ▶ application of lexical association measures to **collocation extraction**

Objectives

1. to compile a comprehensive **inventory** of lexical association measures
2. to build several **reference data** sets for collocation extraction
3. to **evaluate** the lexical association measures on these data sets
4. to explore the possibility of **combining** these measures into more complex models and **improve performance** in collocation extraction

Limitations

- ✗ focus on bigram (*two-word*) collocations

Collocational association

5/25

Collocability

- ▶ the ability of words to combine with other words in text
- ▶ governed by a system of rules and constraints: *syntactic*, *semantic*, *pragmatic*
- ▶ must be adhered to in order to produce correct, meaningful, fluent utterances
- ▶ specified *intensionally* or *extensionally*

Collocations

- ▶ **extensionally** restricted word combinations
- ▶ should be listed in a **lexicon** and learned in the same way as single words

Types of collocations

1. idioms (*to kick the bucket*, *to hear st. through the grapevine*)
2. proper names (*New York*, *Old Town*)
3. technical terms (*car oil*, *stock owl*)
4. phrasal verbs (*to switch off*, *to look after*)
5. light verb compounds (*to take a nap*, *to do homework*)
6. lexically restricted expressions (*strong tea*, *broad daylight*)

Collocation extraction

6/25

Task

- ▶ to extract a list of collocations (*types*) from a text corpus (**collocation lexicon**)
- ▶ no need to identify particular occurrences (*instances*) of collocations

Methods

- ▶ based on **extraction principles** verifying characteristic collocation properties
- ▶ i.e. **hypotheses** about word occurrences and cooccurrences in the corpus
- ▶ formulated as **lexical association measures**
- ▶ compute **association score** for each collocation candidate from the corpus
- ▶ the scores indicate **a chance** of a candidate **to be a collocation**

Extraction principles

1. *“Collocation components occur together more often than by chance”*
2. *“Collocations occur as units in information-theoretically noisy environment”*
3. *“Collocations occur in different contexts to their components”*

Extraction pipeline

7/25

1. linguistic preprocessing (*morphological and syntactic level*)
2. identification of **collocation candidates** (*dependency/surface/distance bigrams*)
3. extraction of occurrence and cooccurrence statistics (*frequency, contexts*)
4. **filtering** the candidates to improve precision (*POS patterns, frequency*)
5. application of a chosen lexical association measure
6. **ranking/classification** of collocation candidates according to their scores

Ranking

<i>červený kříž</i>	15.66
<i>řádková čárka</i>	14.01
<i>aritmetická operace</i>	10.52
<i>podavač papíru</i>	10.17
<i>systém typu</i>	3.54
<i>na další</i>	0.54
<i>program v</i>	0.35
<i>úroveň být</i>	0.25

Classification

<i>červený kříž</i>	1
<i>řádková čárka</i>	1
<i>aritmetická operace</i>	1
<i>podavač papíru</i>	1
<hr/>	
<i>systém typu</i>	0
<i>na další</i>	0
<i>program v</i>	0
<i>úroveň být</i>	0

Lexical association measures

8/25

#	Name	Formula
1.	Joint probability	$P(xy)$
2.	Conditional probability	$P(y x)$
3.	Reverse conditional probability	$P(x y)$
4.	Pointwise mutual information	$\log \frac{P(xy)}{P(x)P(y)}$
5.	Mutual dependency (MD)	$\log \frac{P(xy)^2}{P(x)P(y)}$
6.	Log frequency biased MD	$\log \frac{P(xy)}{P(x)P(y)} + \log P(xy)$
7.	Normalized expectation	$\frac{f(xy)}{f(x)}$
8.	Mutual expectation	$\frac{f(xy)(1-f(xy))}{f(x)(1-f(x))} - P(xy)$
9.	Salience	$\log \frac{P(xy)}{P(x)P(y)} - \log f(xy)$
10.	Pearson's χ^2 test	$\sum_{i,j} \frac{(f_{ij}-f_{i.}f_{.j})^2}{f_{i.}f_{.j}}$
11.	Fisher's exact test	$\frac{f(x,y)(1-f(x,y))f(y)}{f(x)(1-f(x))f(y)}$
12.	t test	$\frac{f(xy)-f(x)f(y)}{\sqrt{f(x)(1-f(x))f(y)(1-f(y))}}$
13.	z score	$\frac{f(xy)-f(x)f(y)}{\sqrt{f(x)(1-f(x))N}}$
14.	Poisson significance measure	$\frac{f(xy)-f(x)f(y)+\log f(xy)}{f(x)-f(x)\log f(x)+\log f(xy)}$
15.	Log likelihood ratio	$-2 \sum_{i,j} f_{ij} \log f_{ij} / f_{i.} / f_{.j}$
16.	Squared log likelihood ratio	$-2 \sum_{i,j} \log f_{ij}^2 / f_{i.} / f_{.j}$
17.	Russel-Rao	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
18.	Sokal-Michiner	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
19.	Rogers-Tanimoto	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j}) + \sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
20.	Hamann	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
21.	Third Sokal-Sneath	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
22.	Jaccard	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
23.	First Kulczynski	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
24.	Second Sokal-Sneath	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
25.	Second Kulczynski	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
26.	Fourth Sokal-Sneath	$\frac{\sum_{i,j} \min(f_{ij}, f_{i.}, f_{.j})}{\sum_{i,j} \max(f_{ij}, f_{i.}, f_{.j})}$
27.	Odds ratio	$\frac{\frac{f(xy)}{f(x)f(y)}}{\frac{f(x)f(y)}{f(xy)}}$
28.	Yulle's ω	$\frac{\log \frac{f(xy)}{f(x)f(y)}}{\log \frac{f(x)f(y)}{f(xy)}}$
29.	Yulle's Q	$\frac{\log \frac{f(xy)}{f(x)f(y)}}{\log \frac{f(x)f(y)}{f(xy)}}$
30.	Driver-Kroeber	$\frac{\log \frac{f(xy)}{f(x)f(y)}}{\log \frac{f(x)f(y)}{f(xy)}}$
31.	Fifth Sokal-Sneath	$\frac{\sqrt{(\alpha+\beta)(\alpha+\gamma)}}{\sqrt{(\alpha+\beta)(\alpha+\gamma)} + \sqrt{(\alpha+\beta)(\alpha+\gamma)}}$
32.	Pearson	$\frac{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)} + \sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}$
33.	Barni-Urbani	$\frac{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)} + \sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}$
34.	Braun-Blanquet	$\frac{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}{\sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)} + \sqrt{(\alpha+\beta)(\alpha+\gamma)(\alpha+\delta)(\beta+\gamma)}}$
35.	Simpson	$\frac{\min(\alpha+\beta, \alpha+\gamma)}{\max(\alpha+\beta, \alpha+\gamma)}$
36.	Michael	$\frac{\min(\alpha+\beta, \alpha+\gamma)}{\max(\alpha+\beta, \alpha+\gamma)}$
37.	Mountford	$\frac{\min(\alpha+\beta, \alpha+\gamma)}{\max(\alpha+\beta, \alpha+\gamma)}$
38.	Fager	$\frac{\sqrt{(\alpha+\beta)(\alpha+\gamma)}}{\sqrt{(\alpha+\beta)(\alpha+\gamma)}} - \frac{1}{2} \max(h, c)$
39.	Unigram subtuples	$\log \frac{f_{ij}}{f_{i.}f_{.j}} - 3.29 \sqrt{\frac{1}{2} + \frac{1}{2} + \frac{1}{2}}$
40.	U cost	$\log(1 + \frac{\min(\alpha, \beta, \gamma, \delta)}{\max(\alpha, \beta, \gamma, \delta)})$
41.	S cost	$\log(1 + \frac{\min(\alpha, \beta, \gamma, \delta)}{\max(\alpha, \beta, \gamma, \delta)})$
42.	R cost	$\log(1 + \frac{\min(\alpha, \beta, \gamma, \delta)}{\max(\alpha, \beta, \gamma, \delta)})$
43.	T combined cost	$\sqrt{U \times S \times R}$
44.	Phi	$\frac{P(xy) - P(x)P(y)}{\sqrt{P(x)P(y)(1-P(x))(1-P(y))}}$
45.	Kappa	$\frac{P(xy) - P(x)P(y)}{1 - P(x)P(y)}$

#	Name	Formula
46.	J measure	$\max\{P(xy) \log \frac{P(xy)}{P(x)P(y)} + P(xy) \log \frac{P(xy)}{P(x)P(y)}, P(xy) \log \frac{P(xy)}{P(x)P(y)} + P(xy) \log \frac{P(xy)}{P(x)P(y)}\}$
47.	Gini index	$\max\{P(x)(P(y x)^2) + P(y)(P(x y)^2) - P(xy)^2, P(x)(P(y x)^2) + P(y)(P(x y)^2) - P(xy)^2, P(y)(P(x y)^2) + P(x)(P(y x)^2) - P(xy)^2\}$
48.	Confidence	$\max\{P(y x), P(x y)\}$
49.	Laplace	$\max\{\frac{N P(xy)+1}{N P(x)+1}, \frac{N P(y)+1}{N P(y)+1}\}$
50.	Conviction	$\max\{\frac{P(xy)P(x)}{P(y)}, \frac{P(xy)P(y)}{P(x)}\}$
51.	Platersky-Shapiro	$\max\{P(xy) - P(x)P(y), P(x y) - P(x)\}$
52.	Certainty factor	$\max\{\frac{P(xy)-P(x)P(y)}{P(x)}, \frac{P(xy)-P(x)P(y)}{P(y)}\}$
53.	Added value (AV)	$\max\{P(y x) - P(y), P(x y) - P(x)\}$
54.	Collective strength	$\frac{P(xy) + P(x)P(y)}{P(x)P(y) + P(x)P(y)} - \frac{1 - P(x)P(y) - P(x)P(y)}{1 - P(x) - P(y)}$
55.	Klogsen	$\sqrt{P(xy) \cdot AV}$
56.	Context entropy	$-\sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
57.	Left context entropy	$-\sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
58.	Right context entropy	$-\sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
59.	Left context divergence	$P(x) \log P(x) - \sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
60.	Right context divergence	$P(y) \log P(y) - \sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
61.	Cross entropy	$-\sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
62.	Reverse cross entropy	$-\sum_{i,j} P(w C_{ij}) \log P(w C_{ij})$
63.	Intersection measure	$\frac{P(xy)P(x)}{P(x)}$
64.	Euclidean norm	$\sqrt{\sum_{i,j} (P(w C_{ij}) - P(w C_{ij}))^2}$
65.	Cosine norm	$\frac{P(xy)P(x)P(y)}{\sqrt{P(x)P(y)P(xy)}}$
66.	L1 norm	$\sum_{i,j} P(w C_{ij}) - P(w C_{ij}) $
67.	Confusion probability	$\sum_{i,j} \frac{P(w C_{ij}) P(w C_{ij}) - P(w C_{ij}) }{P(w C_{ij})}$
68.	Reverse confusion probability	$\sum_{i,j} \frac{P(w C_{ij}) P(w C_{ij}) - P(w C_{ij}) }{P(w C_{ij})}$
69.	Jensen-Shannon divergence	$\frac{1}{2} D(p(w C_x) \frac{1}{2}(p(w C_x) + p(w C_y))) + D(p(w C_y) \frac{1}{2}(p(w C_x) + p(w C_y)))$
70.	Cosine of pointwise MI	$\sqrt{2 \times MI(w, x)} \cdot \sqrt{2 \times MI(w, y)}$
71.	KL divergence	$\sum_{i,j} P(w C_{ij}) \log \frac{P(w C_{ij})}{P(w C_{ij})}$
72.	Reverse KL divergence	$\sum_{i,j} P(w C_{ij}) \log \frac{P(w C_{ij})}{P(w C_{ij})}$
73.	Skew divergence	$D(P(w C_x) \alpha p(w C_x) + (1-\alpha)p(w C_y))$
74.	Reverse skew divergence	$D(P(w C_y) \alpha p(w C_x) + (1-\alpha)p(w C_y))$
75.	Phrase word cooccurrence	$\frac{1}{2} (\frac{P(xy)P(x)}{P(x)} + \frac{P(xy)P(y)}{P(y)})$
76.	Word association	$\frac{1}{2} (\frac{P(xy)P(x)}{P(x)} + \frac{P(xy)P(y)}{P(y)})$
Cosine context similarity:		
	$c_x = (x_i)$	$\cos(c_x, c_y) = \frac{P(xy)}{\sqrt{P(x)P(y)}}$
77.	in boolean vector space	$z_i = \delta f(w_i C_x)$
78.	in tf vector space	$z_i = f(w_i C_x)$
79.	in tf · idf vector space	$z_i = f(w_i C_x) \cdot \frac{1}{\sqrt{f(w_i)}}; df(w_i) = [x: w_i C_x]$
Dice context similarity:		
	$c_x = (x_i)$	$\text{dice}(c_x, c_y) = \frac{2P(xy)}{P(x) + P(y)}$
80.	in boolean vector space	$z_i = \delta f(w_i C_x)$
81.	in tf vector space	$z_i = f(w_i C_x)$
82.	in tf · idf vector space	$z_i = f(w_i C_x) \cdot \frac{1}{\sqrt{f(w_i)}}; df(w_i) = [x: w_i C_x]$

Table 1: Inventory of lexical association measures for collocation extraction.

Reference data set: *PDT-Dep*

9/25

Source corpus

- ▶ Prague Dependency Treebank 2.0, 1.5 mil. tokens
- ▶ manually annotated on *morphological* and *analytical* level

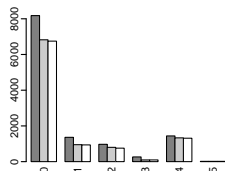
Collocation candidates

- ▶ **dependency bigrams**: direct dependency relation between components
- ▶ morphological normalization: *lemma proper + pos + gender + degree + negation*
- ▶ part-of-speech filter: *A:N, N:N, V:N, R:N, C:N, N:V, N:C, D:A, N:A, D:V, N:T, N:D, D:D*
- ▶ frequency filter: *minimal frequency required, $f > 5$*

Annotation

- ▶ three independent parallel annotations (*no context; full agreement required*)
- ▶ six categories, merged into two: **collocations** (1-5), **non-collocations** (0):

5. *idiomatic expressions*
4. *technical terms*
3. *support verb constructions*
2. *proper names*
1. *stock phrases, frequent unpredictable usages*
0. *non-collocations*



Additional data sets

10/25

PDT-Surf

- ▶ analogous to *PDT-Dep* (*corpus, filtering, annotation*)
- ▶ collocation candidates extracted as **surface bigrams**: pairs of adjacent words
- ▶ **assumption**: collocations cannot be modified by insertion of another word
- ▶ annotation consistent with *PDT-Dep*

CNC-Surf

- ▶ collocation candidates – instances of *PDT-Surf* in the *Czech National Corpus*
- ▶ SYN 2000 and 2005, 240 mil. tokens, morphologically tagged and lemmatized

PAR-Dist

- ▶ source corpus: **Swedish Parole**, 22 mil. tokens
- ▶ automatic morphological tagging and lemmatization
- ▶ **distance bigrams**: word pairs occurring within a distance of 1–3 words
- ▶ **annotation**: non-exhaustive manual extraction of **support verb constructions**
- ▶ no frequency filter applied

Reference data summary

11/25

<i>reference data set</i>	<i>PDT-Dep</i>	<i>PDT-Surf</i>	<i>CNC-Surf</i>	<i>PAR-Dist</i>
morphology	<i>manual</i>	<i>manual</i>	<i>auto</i>	<i>auto</i>
syntax	<i>manual</i>	<i>none</i>	<i>none</i>	<i>none</i>
bigram types	<i>dependency</i>	<i>surface</i>	<i>surface</i>	<i>distance</i>
tokens	1 504 847	1 504 847	242 272 798	22 883 361
bigram types	635 952	638 030	30 608 916	13 370 375
after frequency filtering	26 450	29 035	2 941 414	13 370 375
after part-of-speech filtering	12 232	10 021	1 503 072	898 324
collocation candidates	12 232	10 021	9 868	17 027
data sample size	100 %	100 %	0.66 %	1.90 %
true collocations	2 557	2 293	2 263	1 292
baseline precision (%)	21.02	22.88	22.66	7.59

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|}$$

$$\text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking

červený kříž	15.66
železná opona	15.23
řádková čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking

<i>červený kříž</i>	15.66
<i>železná opona</i>	15.23
<i>řádová čárka</i>	14.01
<i>kupónová knížka</i>	13.83
<i>autor knihy</i>	11.05
<i>aritmetická operace</i>	10.52
<i>podavač papíru</i>	10.17
<i>nová kniha</i>	10.09
<i>kulatý stůl</i>	7.03
<i>nová vlna</i>	6.59
<i>čerpací stanice</i>	6.04
<i>systém typu</i>	3.54
<i>centrum města</i>	1.54
<i>na další</i>	0.54
<i>program v</i>	0.35
<i>úroveň být</i>	0.25

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	0
aritmetická operace	0
podavač papíru	0
nová kniha	0
kulatý stůl	0
nová vlna	0
čerpací stanice	0
systém typu	0
centrum města	0
na další	0
program v	0
úroveň být	0

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	0
aritmetická operace	0
podavač papíru	0
nová kniha	0
kulatý stůl	0
nová vlna	0
čerpací stanice	0
systém typu	0
centrum města	0
na další	0
program v	0
úroveň být	0

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	0
aritmetická operace	0
podavač papíru	0
nová kniha	0
kulatý stůl	0
nová vlna	0
čerpací stanice	0
systém typu	0
centrum města	0
na další	0
program v	0
úroveň být	0

Precision	Recall
100 %	50 %

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
<hr/>	
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	1
<hr/>	
aritmetická operace	0
podavač papíru	0
nová kniha	0
kulatý stůl	0
nová vlna	0
čerpací stanice	0
systém typu	0
centrum města	0
na další	0
program v	0
úroveň být	0

Precision	Recall
100 %	50 %
80 %	50 %

Evaluation measures

12/25

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	1
aritmetická operace	1
podavač papíru	1
nová kniha	1
kulatý stůl	1
nová vlna	1
čerpací stanice	1
systém typu	1
centrum města	1
na další	1
program v	1
úroveň být	1

Precision	Recall
100 %	12 %
100 %	25 %
100 %	37 %
100 %	50 %
80 %	50 %
83 %	62 %
85 %	75 %
75 %	75 %
77 %	87 %
70 %	87 %
72 %	100 %
66 %	100 %
61 %	100 %
57 %	100 %
53 %	100 %
50 %	100 %

Evaluation measures

$$1) \text{ Precision} = \frac{|\text{correctly classified collocations}|}{|\text{total classified as collocations}|} \quad \text{Recall} = \frac{|\text{correctly classified collocations}|}{|\text{total collocations}|}$$

Ranking	
červený kříž	15.66
železná opona	15.23
řádová čárka	14.01
kupónová knížka	13.83
autor knihy	11.05
aritmetická operace	10.52
podavač papíru	10.17
nová kniha	10.09
kulatý stůl	7.03
nová vlna	6.59
čerpací stanice	6.04
systém typu	3.54
centrum města	1.54
na další	0.54
program v	0.35
úroveň být	0.25

Classification	
červený kříž	1
železná opona	1
řádová čárka	1
kupónová knížka	1
autor knihy	1
aritmetická operace	1
podavač papíru	1
nová kniha	1
kulatý stůl	1
nová vlna	1
čerpací stanice	1
systém typu	1
centrum města	1
na další	1
program v	1
úroveň být	1

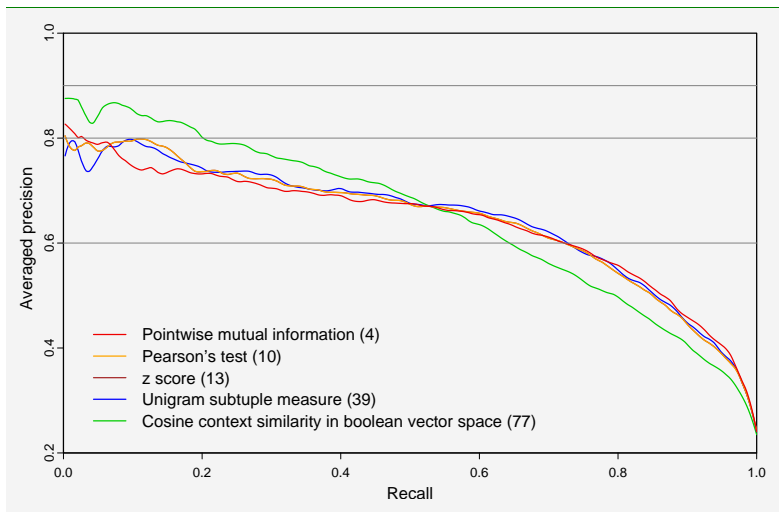
Precision	Recall
100 %	12 %
100 %	25 %
100 %	37 %
100 %	50 %
80 %	50 %
83 %	62 %
85 %	75 %
75 %	75 %
77 %	87 %
70 %	87 %
72 %	100 %
66 %	100 %
61 %	100 %
57 %	100 %
53 %	100 %
50 %	100 %

$$2) \text{ (Mean) Average Precision} = \frac{1}{r} \sum_{i=1}^r p_i$$

89.6 % = AP

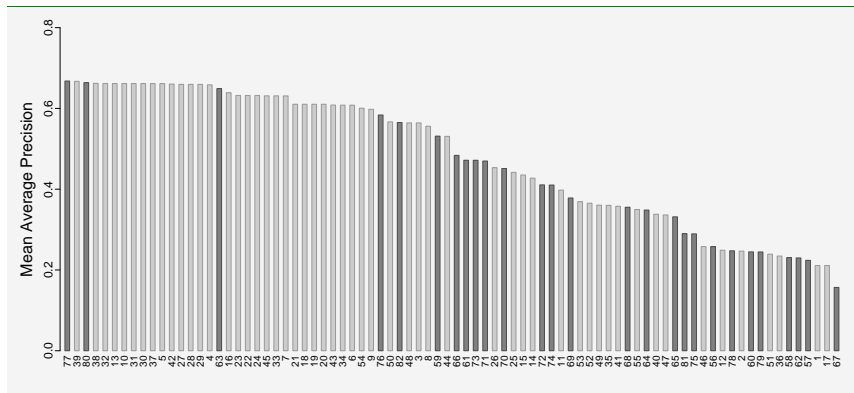
Results / Precision-Recall: *PDT-Dep*

13/25



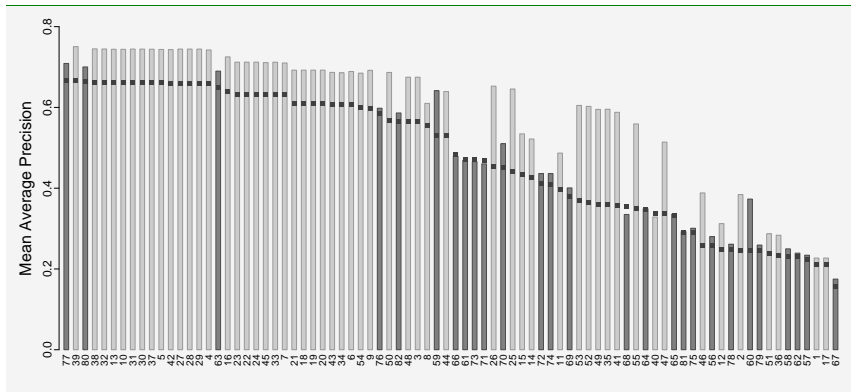
Results / Mean average precision: *PDT-Dep*

14/25



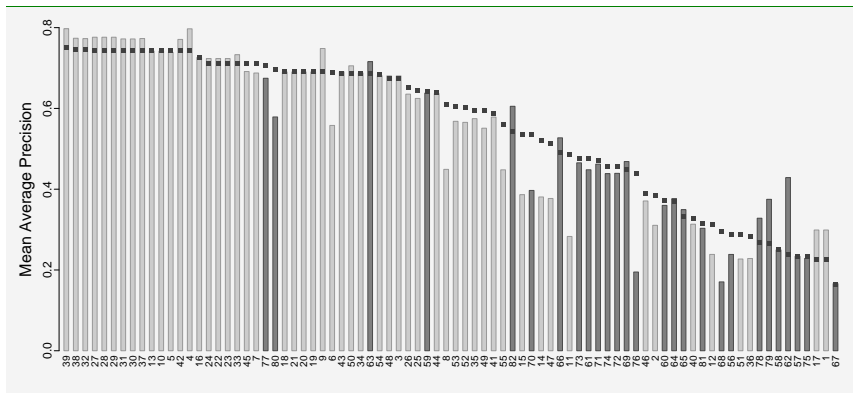
Results / Mean average precision: *PDT-Dep* vs. *PDT-Surf*

15/25



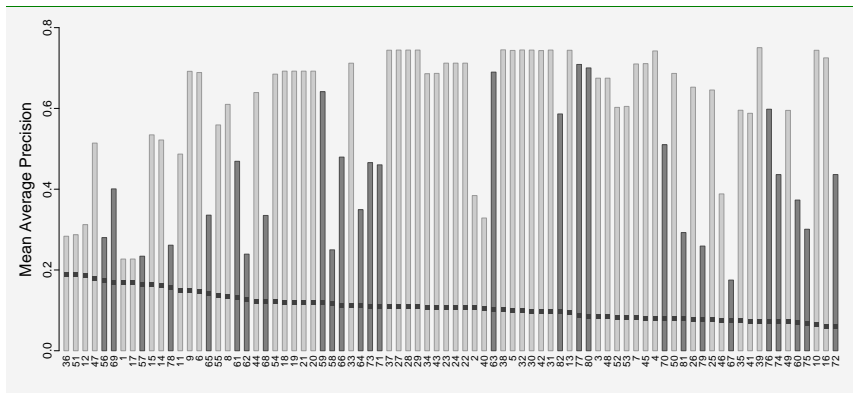
Results / Mean average precision: *PDT-Surf* vs. *CNC-Surf*

16/25



Results / Mean average precision: *PAR-Dist* vs. *PDT-Dep*

17/25

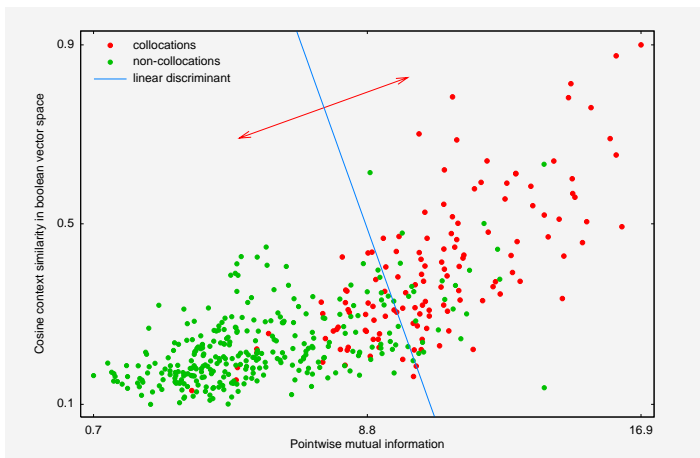


Combining association measures

18/25

Motivation

- ▶ different association measures discover different groups/types of collocations
- ▶ existence of uncorrelated association measures



Combination models

19/25

Framework

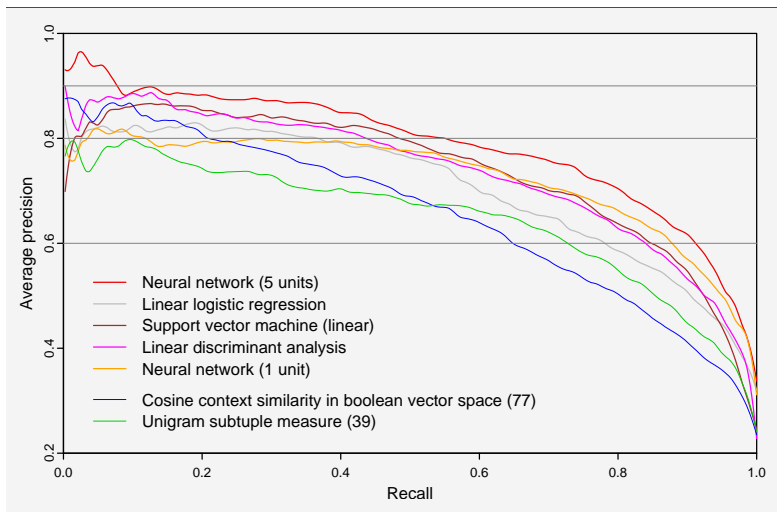
- ▶ each collocation candidate \mathbf{x}^i is described by the **feature vector** $\mathbf{x}^i = (x_1^i, \dots, x_{82}^i)^T$ consisting of scores of all association measures
- ▶ and assigned a **label** $y^i \in \{0, 1\}$ indicating whether the bigram is considered to be a true collocation ($y = 1$) or not ($y = 0$)
- ▶ we look for a **ranker function** $f(\mathbf{x}^i)$ determining the strength of lexical association between components of a candidate \mathbf{x}^i
- ▶ e.g. **linear combination** of association scores: $f(\mathbf{x}^i) = w_0 + w_1 x_1^i + \dots + w_{82} x_{82}^i$

Methods

1. *Linear logistic regression*
 2. *Linear discriminant analysis*
 3. *Support vector machines*
 4. *Neural networks*
- ▶ in the **training phase** used as regular classifiers on two-class data
 - ▶ in the **application phase** no classification threshold applies

Results / Precision-Recall: *PDT-Dep*

20/25



Results / Mean average precision: *PDT-Dep*

21/25

<i>method</i>	<i>precision at</i>				
	<i>R=20</i>	<i>R=50</i>	<i>R=80</i>	<i>MAP</i>	<i>+%</i>
Neural Network (5 units)	91.00	81.75	70.22	80.87	21.08
Linear Logistic Regression	86.96	79.74	64.63	77.36	15.82
Linear Discriminant Analysis	85.99	77.34	61.44	75.16	12.54
Neural Network (1 unit)	82.47	77.08	65.75	74.88	12.11
Support Vector Machine	81.33	76.08	61.49	73.03	9.35
Cosine similarity (77)	80.88	68.46	49.99	66.79	0.00
Unigram subtuples (39)	75.86	68.19	55.13	66.72	–

Adding linguistic features

22/25

Combined variables:

AM: association measures**AM+POS:** association measures + part-of-speech pattern**AM+POS+DEP:** association measures + part-of-speech + analytical function

<i>reference data set</i>	<i>PDT-Dep</i>	<i>PDT-Surf</i>	<i>CNC-Surf</i>	<i>PAR-Dist</i>
Baseline	21.02	22.88	22.66	7.59
Best AM	66.72	75.03	79.74	18.88
NNet.5 (AM)	80.87	84.35	86.30	35.78
NNet.5 (AM+POS)	82.79	86.48	88.22	–
NNet.5 (AM+POS+DEP)	84.53	–	–	–

Model reduction

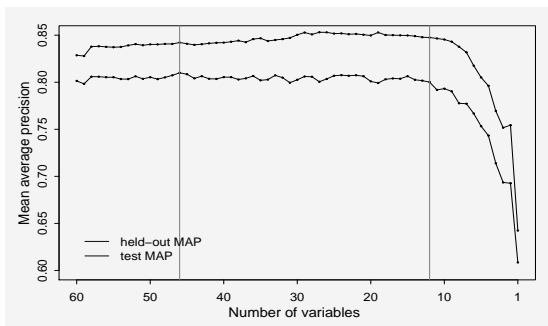
23/25

Motivation

- ▶ models combining all 82 association measures are too complex
- ▶ **redundant/improper** variables should be removed

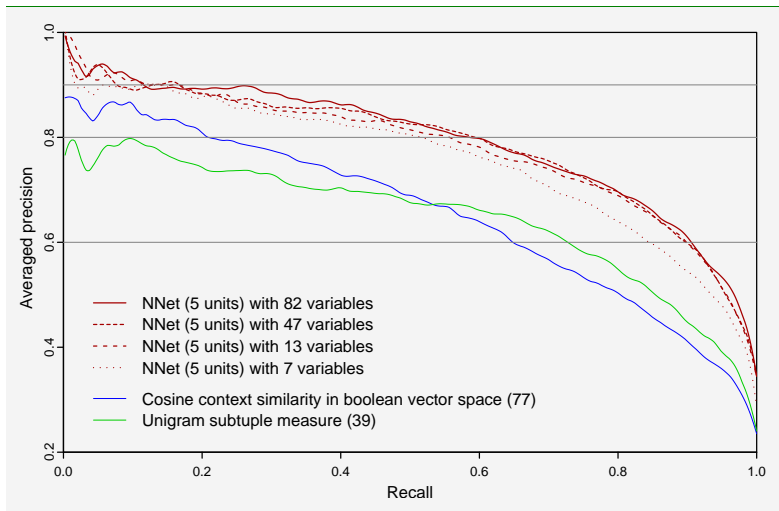
Algorithm

1. correlation based **clustering**; one representative selected from each cluster
2. **step-wise** procedure removing the variable causing minimal performance degradation in each iteration



Model reduction results / Precision-Recall: *PDT-Dep*

24/25



Conclusions

25/25

Objectives accomplished

1. to compile a comprehensive **inventory** of lexical association measures
 - *82 association measures based on 3 extraction principles, unified notation*
2. to build several **reference data** sets for collocation extraction
 - *4 reference data sets for evaluation in different settings*
3. to **evaluate** the lexical association measures on these data sets
 - *performance heavily depends on the task, data, language, etc.*
 - *no general “best-performing” measure selected*
4. to **combine** these measures and **improve performance** in collocation extraction
 - *combining is meaningful and greatly improves performance*
 - *combination models can be reduced to use only 13 variables*

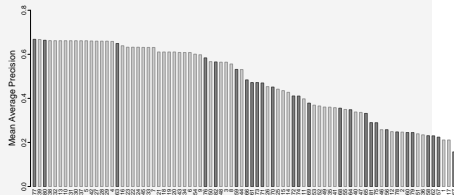
Other contribution

- ▶ overview of different notions of collocation (*definitions, typology, classification*)
- ▶ evaluation scheme (*average precision, crossvalidation, significance tests*)
- ▶ reference data sets used in MWE 2008 Shared Task

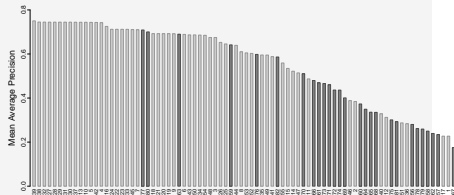
Context-based vs. statistical association measures

25/25

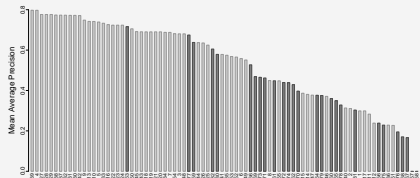
PDT-Dep



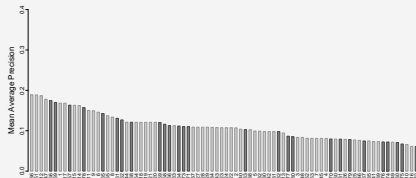
PDT-Surf



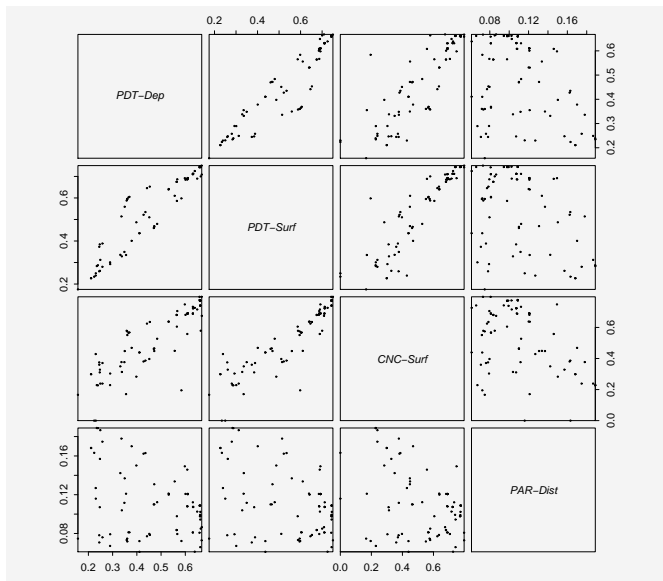
CNC-Surf



PAR-Dist



Comparison of AM evaluation results (MAP)



Extraction principle I

“Collocation components occur together more often than by chance”

- ▶ the corpus is interpreted as a sequence of **randomly generated words**
- ▶ word (*marginal*) probability ML estimations: $P(x) \approx \frac{f(x)}{N}$
- ▶ bigram (*joint*) probability ML estimations: $P(xy) \approx \frac{f(xy)}{N}$
- ▶ the **chance** \sim the **null hypothesis of independence**: $H_0: \hat{P}(xy) = P(x) \cdot P(y)$

Example

Data: $f(\text{životní prostředí}) = 220$	MLE: $p(\text{životní prostředí}) = 0.000146$
$f(\text{životní}) = 261$	$p(\text{životní}) = 0.000174$
$f(\text{prostředí}) = 439$	$p(\text{prostředí}) = 0.000292$

$H_0: \hat{p}(\text{životní prostředí}) = p(\text{životní}) \cdot p(\text{prostředí}) = 0.0000000005$
 $\hat{f}(\text{životní prostředí}) = 0.076$

AM: $PMI(\text{životní prostředí}) = \log \frac{f(\text{životní prostředí})}{\hat{f}(\text{životní prostředí})} = \log \frac{220}{0.076} = 11.53$

Extraction principle II

“Collocations occur as units in information-theoretically noisy environment”

- ▶ the corpus again interpreted as a sequence of **randomly generated words**
- ▶ at each point of the sequence we:
 1. estimate probability distribution of words occurring after (and before)
 2. measure uncertainty (*entropy*) what the next (or previous) word will be
- ▶ points with **high uncertainty** are likely to be **collocation boundaries**
- ▶ points with **low uncertainty** are likely to be **located within a collocation**

Example



Český *kapitálový trh* dnes ovlivnil pokles cen všech *cenných papírů* a zejména akcií.

Extraction principle III

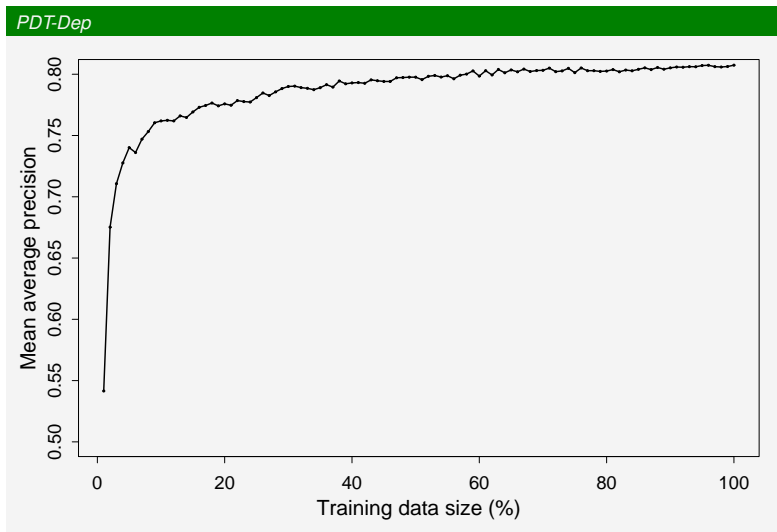
“Collocations occur in different contexts to their components”

- ▶ **non-compositionality**: meaning of a collocation must differ from the *union* of the meaning of its components
- ▶ modeling meanings by **empirical contexts**: a bag of words occurring within a specified context window of a word or an expression
- ▶ the **more different the contexts** of an expression to its components are, the higher the chance is that the expression is a collocation








Example

... není. **Maltské liry lze nakoupit pouze ve směnárnách, černý trh** s valutami neexistuje. Na Maltě je v porovnání s ...
 ... přestal. **V patách za krizí vstoupil do Bělehradu černý trh**, pašování a zvýšená kriminalita. Překupníci provázejí ...
 ... nebyli z toho obviněni. **Řídí gangy, které kontrolují černý trh** a okrádají cizince. Oba byli zbaveni funkcí a byl ...
 ... antidrogové hysterii. **Následkem toho neexistoval ani černý trh**, protože nebylo na čem vydělávat. V roce 1957 bylo ...
 ... doručeny k rychlému zpracování. **Naplnlo se již rozjíždí černý trh** se vstupenkami. Na závod na 5000 m v rychlobruslařů ...
 ... na čelném místě obchodu se zbraněmi. Zatímco **černý trh** se zbraněmi se pro celý svět stává čím dál tím větší. ...
 ... čtením v parlamentu. **Věřím, že brzy bude regulovat černý trh** s ohroženými druhy zvířat, míní. Promoravské strany ...
 ... jako malí čtyřletí a pětiletí kluci. **Byl to dobytčí trh** jako z minulého století. Se vším všudy prodávali ...
 ... přání než reálných možností. **Na rozdíl od dolaru se trh** amerických státních dluhopisů nezměnil. A novými ...
 ... opětnému nárůstu. **Podle Plan Econu si český kapitálový trh** bude v nejbližším roce počínat o něco lépe. Většina ...
 ... **To by mohlo vzhledem k propojení přes mezibankovní trh** depozit vést k řetězovým reakcím. Příliv kapitálu ...
 ... PVT, na ceně ztratil také indexový Tabák. **Volný trh** má však našťást i světlé stránky. K nim patří například ...
 ... spoluzakladatel. **Také v Maďarsku se uvolní mediální trh** již letos. Maďarsko jako první z postkomunistických ...
 **Mezi ně patří i OfficePorte Voice, který byl na trh** uveden pod heslem "více než modem". Obsahuje totiž ...

Learning curve



List of relevant publications

-  Pavel Pecina: **Lexical Association Measures and Collocation Extraction**, *Multiword expressions: Hard going or plain sailing? Special issue of the International Journal of Language Resources and Evaluation*, Springer, 2009 (accepted).
-  Pavel Pecina: **Machine Learning Approach to Multiword Expression Extraction**, *In Proceedings of the sixth International Conference on Language Resources and Evaluation (LREC) Workshop: Towards a Shared Task for Multiword Expressions*, Marrakech, Morocco, 2008.
-  Pavel Pecina: **Reference Data for Czech Collocation Extraction**, *In Proceedings of the sixth International Conference on Language Resources and Evaluation (LREC) Workshop: Towards a Shared Task for Multiword Expressions*, Marrakech, Morocco, 2008.
-  Pavel Pecina, Pavel Schlesinger: **Combining Association Measures for Collocation Extraction**, *In Proceedings of the 21th International Conference on Computational Linguistics and 44th Annual Meeting of the Association for Computational Linguistics (COLING/ACL)*, Sydney, Australia, 2006.
-  Silvie Cinková, Petr Podveský, Pavel Pecina, Pavel Schlesinger: **Semi-automatic Building of Swedish Collocation Lexicon**, *In Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC)*, Genova, Italy, 2006.
-  Pavel Pecina: **An Extensive Empirical Study of Collocation Extraction Methods**, *In Proceedings of the Association for Computational Linguistics Student Research Workshop (ACL)*, Ann Arbor, Michigan, USA, 2005.
-  Pavel Pecina, Martin Holub: **Semantically Significant Collocations**, *UFAL/CKL Technical Report TR-2002-13*, Faculty of Mathematics and Physics, Charles University, Prague, Czech Rep., 2002.