

An Extensive Empirical Study of Collocation Extraction Methods

Pavel Pecina

`pecina@ufal.mff.cuni.cz`

Institute of Formal and Applied Linguistics
Charles University, Prague



June 27, 2005

Outline

- 1 Introduction
 - Notion of Collocation
 - Motivation
 - The Task
- 2 Collocation Extraction
 - Methodology
 - Association Measures
 - Evaluation
- 3 Combining Association Measures
 - Classification and Ranking
 - Attribute Selection
- 4 Summary

Outline

- 1 Introduction
 - Notion of Collocation
 - Motivation
 - The Task
- 2 Collocation Extraction
 - Methodology
 - Association Measures
 - Evaluation
- 3 Combining Association Measures
 - Classification and Ranking
 - Attribute Selection
- 4 Summary

Definitions I

Firth (1951):

“Collocations of a given word are statements of the habitual or customary places of that word.”

Choueka (1988):

“A collocation is a sequence of two or more consecutive words, that has characteristics of a syntactic and semantic unit, and whose exact and unambiguous meaning or connotation cannot be derived directly from the meaning or connotation of its components.”

Čermák (1982):

“Individual words cannot be combined freely or randomly only by syntactic rules. The ability of a word to combine with other words (collocability) can be expressed:

- a) intensionally → valency*
- b) extensionally” → collocations*

Characteristic Properties

Non-compositionality

(kick the bucket, carriage return, white man)

- *The meaning of a collocation is not a straightforward composition of the meaning of its parts.*

Non-substitutability

(yellow wine, hit the bucket, make homework)

- *Components of collocation cannot be substituted with a related word or a synonym.*

Non-modifiability

(give a big hand, poor as church mice)

- *Collocations cannot be modified or syntactically transformed.*

Other properties

- *Collocations are not necessarily adjacent.* *(knock the door)*
- *Collocations cannot be directly translated.* *(ice cream)*
- *Collocations are domain-specific.* *(carriage return)*
- *Judging collocations is subjective.* *(new company)*

Types of Collocations

Collocations have both linguistic and lexicographic character and covers a wide range of lexical phenomena:

- light verb compounds – *verbs with little semantic content* (take, make, do)
- verb particle constructions, phrasal verbs (look up, take off, tell off)
- idioms – *fixed phrases* (kick the bucket)
- stock phrases (good morning)
- technological expressions – *concepts or objects in tech. dom.* (hard disk)
- proper names (Ann Arbor)

Motivation

Collocations can be used in a wide range of fields:

- Lexicography
- Machine translation
- Information retrieval, information extraction
- Word sense disambiguation
- Spell/grammar/style-checking
- Text classification and summarization
- Keyword extraction
- Language modeling
- Language generation

The Tasks

To build a collocation lexicon.

- 1 Creating manually annotated reference data
 - *of reasonable size.*
- 2 Evaluation of collocation extraction methods
 - *interval-wise by the means of precision-recall.*
- 3 Combining association measures for collocation extraction
 - *and achieve "better" results.*
- 4 Reduce number of combined measures
 - *and select the "best subset" of available association measures.*

Focus on bigram collocations

- 1 *Processing of longer expressions requires larger amounts of data.*
- 2 *Scalability of some methods to high order n-grams is limited.*

Outline

- 1 Introduction
 - Notion of Collocation
 - Motivation
 - The Task
- 2 Collocation Extraction**
 - Methodology
 - Association Measures
 - Evaluation
- 3 Combining Association Measures
 - Classification and Ranking
 - Attribute Selection
- 4 Summary

Collocation Extraction

- Most methods are based on **verification** of typical **collocation properties**.
- These properties are formally described by **mathematical formulas** that determine **degree of association** between words.
- Such formulas are called **association measures** and compute **association score** for each collocation candidate from a corpus.
- The scores indicate **a chance** of a candidate **to be a collocation**.
- The scores can be used for **ranking** or for **classification**:

Ranking

<i>red cross</i>	15.66
<i>decimal point</i>	14.01
<i>arithmetic operation</i>	10.52
<i>paper feeder</i>	10.17
<i>system type</i>	3.54
<i>and others</i>	0.54
<i>program in</i>	0.35
<i>level is</i>	0.25

Classification

<i>red cross</i>	1
<i>decimal point</i>	1
<i>arithmetic operation</i>	1
<i>paper feeder</i>	1
<i>system type</i>	0
<i>and others</i>	0
<i>program in</i>	0
<i>level is</i>	0

The Methodology

- 1 Identifying Word Base Forms:
 - Surface forms
 - Stems or lemmas
 - Lemmas with additional morphosyntactic features
- 2 Extracting all possible collocation candidates:
 - Consequent word n-grams (*multi-word expressions*)
 - Sliding window
 - Syntactic structures (*dependency n-grams*)
- 3 Collecting cooccurrence statistics:
 - Frequency of word and n-gram occurrences
 - Immediate contexts
 - Global contexts
- 4 Computing association measures
- 5 Ranking or classification

Word Base Forms

Problem:

- Surface word forms too specific (*rich morphology, we work with Czech*)
- Lemmas too general (*loss of syntactic and semantic information*)

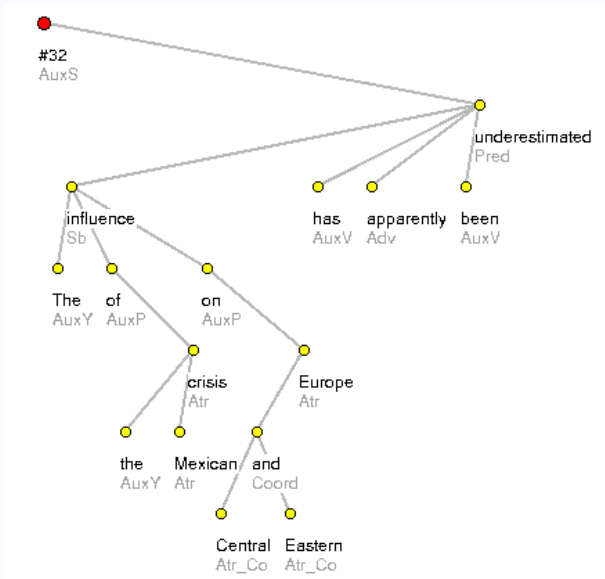
Solution:

- Lemmas with a subset of morphological tags

```

<f>nenahraditelná<l>nahraditelný_(*4)<t>AAFS1----1N----<r>8<g>7
      ↓           ↓ ↓           ↓↓
nahraditelný_(*4)  A F           1N
                ↓
                <f>nahraditelný_(*4)<t>A*F1N</f>
                ↓
                nenahraditelná
  
```

Dependency Bigrams



Cooccurrence Statistics

a) Contingency tables

$f(xy)$	$f(x\bar{y})$	$f(x*)$
$f(\bar{x}y)$	$f(\bar{x}\bar{y})$	$f(\bar{x}*)$
$f(*y)$	$f(*\bar{y})$	N

Example

	$X=black$	$X\neq black$	X
$Y=market$	<i>black market</i>	<i>new market</i>	<i>market</i>
$Y\neq market$	<i>black horse</i>	<i>new horse</i>	<i>horse</i>
Y	<i>black</i>	<i>new</i>	(all)

b) Contexts

C_w	global context of word w
C_{xy}	global context of bigram xy
C'_{xy}	left immediate context of xy
C''_{xy}	right immediate context of xy

Example

dobrá situace . Kapitálový **trh** je však stále nelikvidní
že to není samostatný **trh** a že je součástí širšího
bariérách v přístupu na **trh** , cenových rozdílech ,
banky . Americký akciový **trh** byl za silného obchodování
jít se svou kuží na **trh** . Pro vydání i mluvila

Context word probability distribution $P(w_j|x)$



Types of Association Measures

- 1 *“Collocations are very frequent word combinations.”*
 - ML estimations of joint and conditional probabilities
- 2 *“Collocation components occur together more often than by a chance.”*
 - Mutual information and derived measures
 - Statistical tests of independence
 - Likelihood measures
 - Other heuristic association measures and coefficients
- 3 *“Collocations occur as units in a (inf.-theoretically) noisy environment.”*
 - Immediate context measures
- 4 *“Collocations occur in different contexts than their components.”*
 - Information-theory measures
 - Information-retrieval similarity measures

Total: 84 association measures + 3 morphosyntactic features

Data

Source: Prague Dependency Treebank v 1.0

Sentences: 81,614

Word forms: 1,255,590

Dependency bigram types: 202,171

Reference bigram types ($f > 5$): 21,597

Reference collocation candidates (relevant POS): 8,904

- Data manually annotated according association strength.

4	<i>idioms and completely non-compositional expressions</i>	7
3	<i>partially non-compositional phrases, technical terms</i>	201
2	<i>names of persons, geographical places, and other entities</i>	2,698
1	<i>frequent compositional usages</i>	484
0	<i>non-collocations</i>	5,514

- All association measures computed for all bigrams.
- Comparison by precision-recall curves (no thresholds).

Data

Source: Prague Dependency Treebank v 1.0

Sentences: 81,614

Word forms: 1,255,590

Dependency bigram types: 202,171

Reference bigram types ($f > 5$): 21,597

Reference collocation candidates (relevant POS): 8,904

- Data manually annotated according association strength.

4	<i>idioms and completely non-compositional expressions</i>	2,906
3	<i>partially non-compositional phrases, technical terms</i>	
2	<i>names of persons, geographical places, and other entities</i>	
1	<i>frequent compositional usages</i>	5,998
0	<i>non-collocations</i>	

- All association measures computed for all bigrams.
- Comparison by precision-recall curves (no thresholds).

Data

Source: Prague Dependency Treebank v 1.0

Sentences: 81,614

Word forms: 1,255,590

Dependency bigram types: 202,171

Reference bigram types ($f > 5$): 21,597

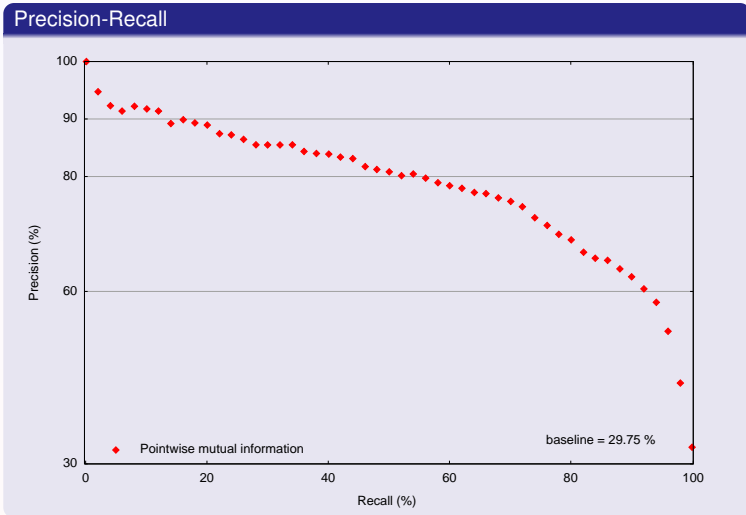
Reference collocation candidates (relevant POS): 8,904

- Data manually annotated according association strength.

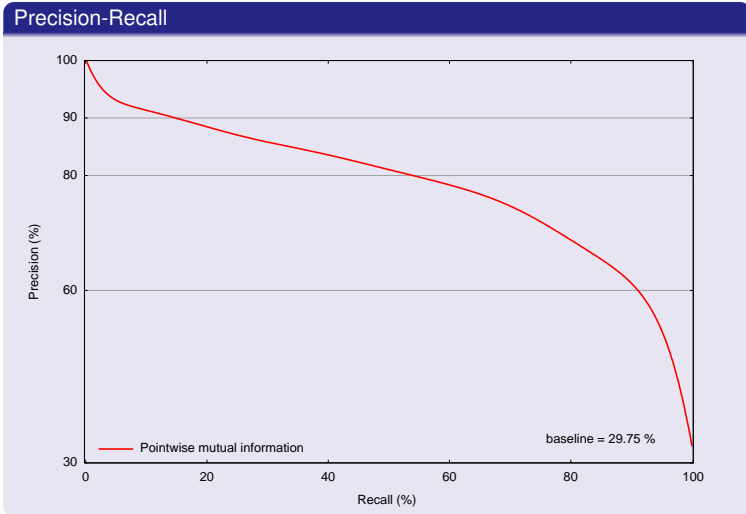
4	<i>idioms and completely non-compositional expressions</i>	29 %
3	<i>partially non-compositional phrases, technical terms</i>	
2	<i>names of persons, geographical places, and other entities</i>	
1	<i>frequent compositional usages</i>	71 %
0	<i>non-collocations</i>	

- All association measures computed for all bigrams.
- Comparison by precision-recall curves (no thresholds).

Precision-Recall Curves

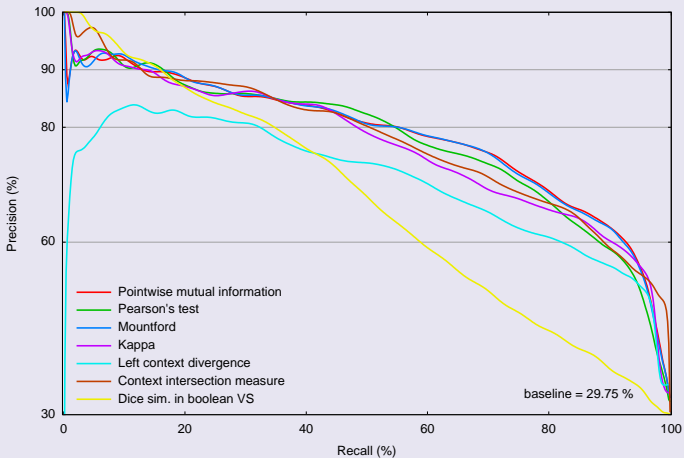


Precision-Recall Curves



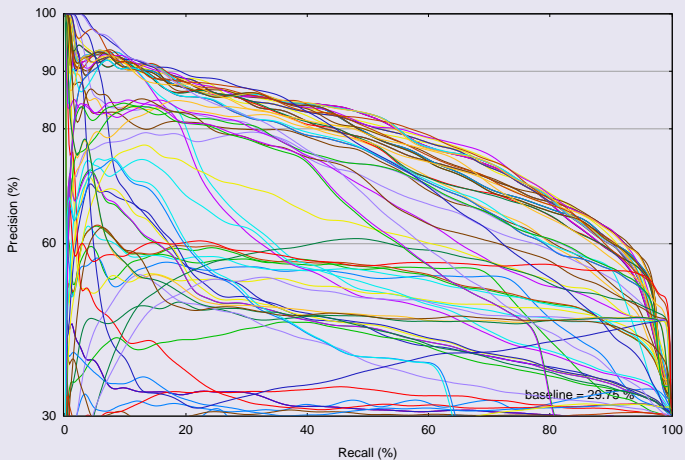
The Best Methods

Precision-Recall curves of the best association measures of each group



All Results

Precision-Recall curves of all association measures

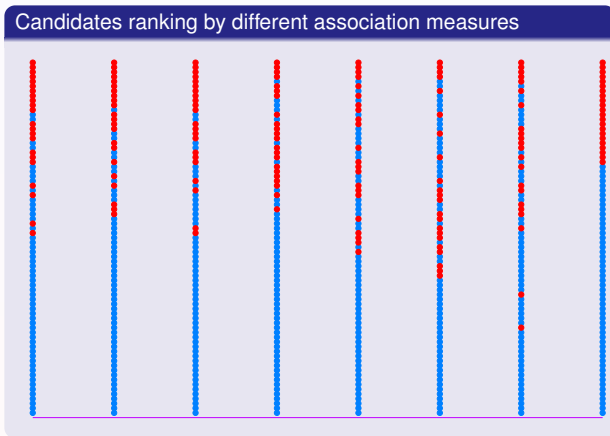


Outline

- 1 Introduction
 - Notion of Collocation
 - Motivation
 - The Task
- 2 Collocation Extraction
 - Methodology
 - Association Measures
 - Evaluation
- 3 Combining Association Measures**
 - Classification and Ranking**
 - Attribute Selection**
- 4 Summary

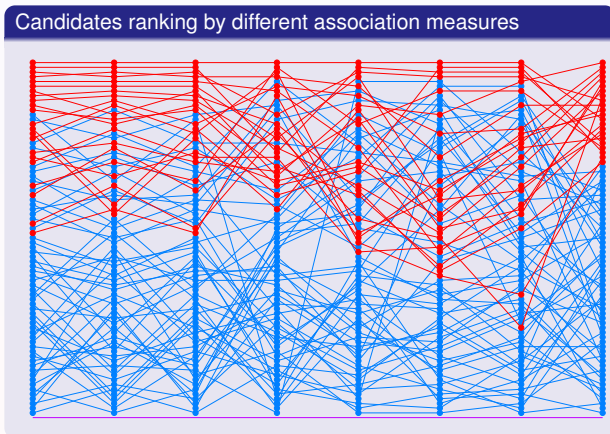
Motivation I

- Can we combine the association measures to get better results?



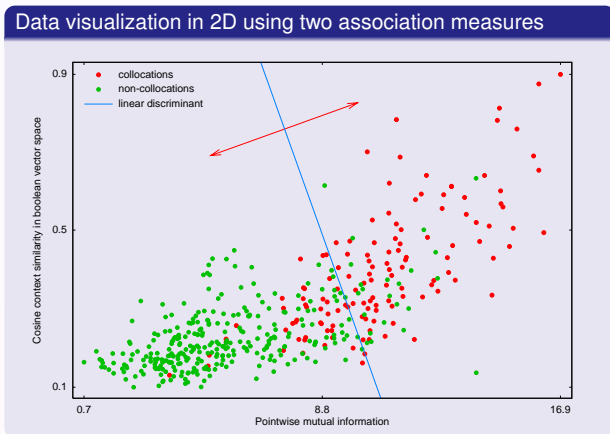
Motivation I

- Can we combine the association measures to get better results?



Motivation II

- Can we combine the association measures to get better results?



Combining Multiple Methods

Voting

- Each method votes whether the candidate is or is not a collocation.
- The final vote depends on the majority of the these votes.

$$\begin{array}{cccccc} x_1, & x_2, & x_3, & x_4 & \dots & x_n \\ \downarrow & \downarrow & \downarrow & \downarrow & & \downarrow \\ 1 & 0 & 1 & 1 & \dots & 0 \end{array} \Rightarrow y$$

Liner combination

- Each association score is weighted by its coefficient.
- The final score is defined as combination of these weighted scores.

$$\beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \dots + \beta_n x_n = y$$

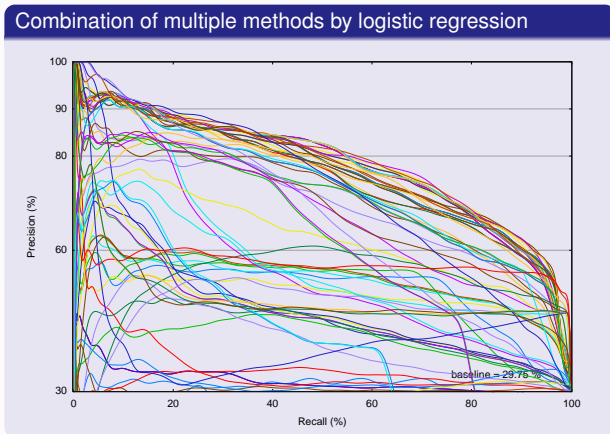
Logistic Regression

$$P(\mathbf{x} \text{ is collocation}) = \frac{\exp^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots + \beta_n x_n}}{1 + \exp^{\beta_0 + \beta_1 x_1 \dots + \beta_2 x_2 + \beta_n x_n}}$$

Logistic Regression

$$P(\mathbf{x} \text{ is collocation}) = \frac{\exp^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots + \beta_n x_n}}{1 + \exp^{\beta_0 + \beta_1 x_1 \dots + \beta_n x_n}}$$

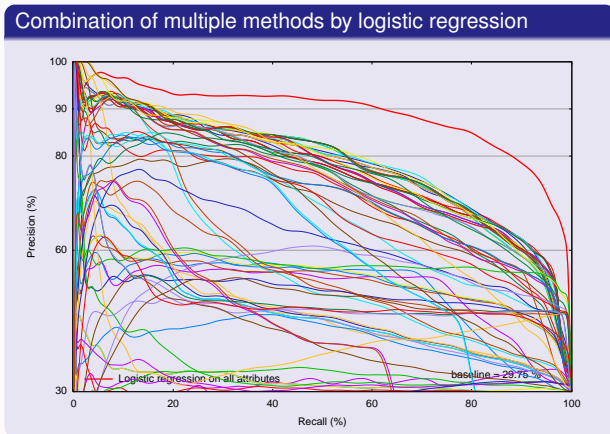
Combination of multiple methods by logistic regression



Logistic Regression

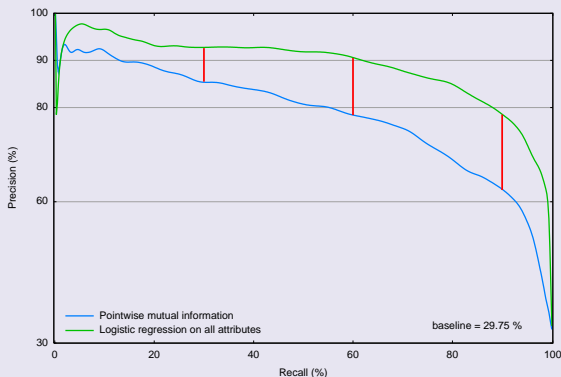
$$P(\mathbf{x} \text{ is collocation}) = \frac{\exp^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 \dots + \beta_n x_n}}{1 + \exp^{\beta_0 + \beta_1 x_1 \dots + \beta_n x_n}}$$

Combination of multiple methods by logistic regression



Logistic Regression: Results I

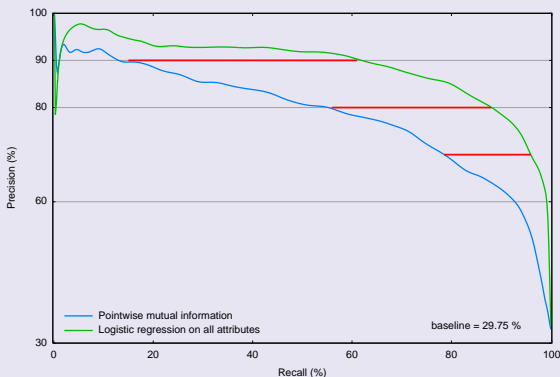
Precision improvement



<i>Recall</i>	30	60	90
P. mutual information	85.5	78.4	62.5
Logistic regression	92.6	89.5	84.5
Absolute improvement	7.1	11.1	22.0
Relative improvement	8.3	14.2	35.2

Logistic Regression: Results II

Recall improvement



<i>Precision</i>	90	80	70
P. mutual information	16.3	56.0	78.0
Logistic regression	55.8	86.7	96.7
Absolute improvement	39.2	30.7	17.7
Relative improvement	242.3	54.8	23.9

Attribute Selection

Can we reduce the number of combined association measures?

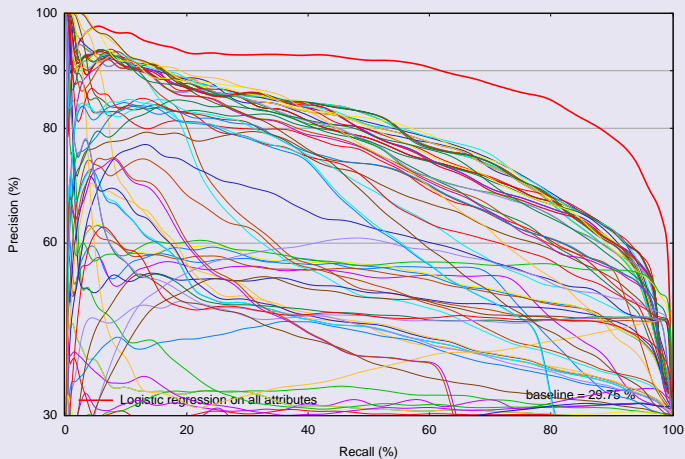
Greedy (stepwise) attribute selection:

- 1 Start with a full set of attributes.
- 2 Estimate parameters of the model.
- 3 Remove the attribute that minimally reduces the performance.
- 4 Repeat until the performance changes significantly.

Result: 87 reduced to 17

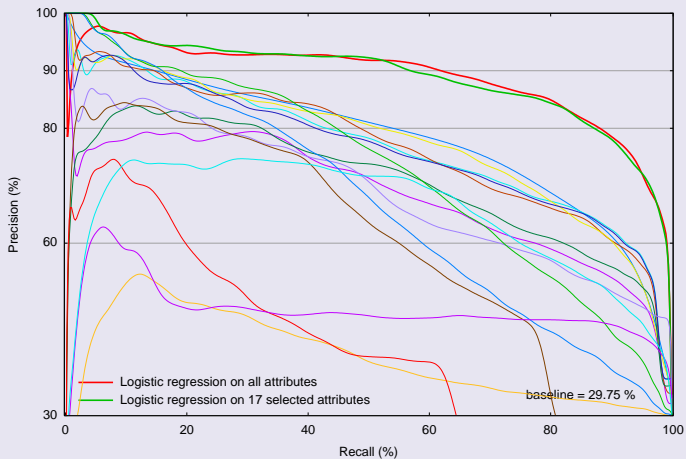
Attribute Selection: Beginning

Logistic regression on all attributes (84 +3 attributes)



Attribute Selection: End

Greedy attribute selection using logistic regression (17 attributes)



Outline

- 1 Introduction
 - Notion of Collocation
 - Motivation
 - The Task
- 2 Collocation Extraction
 - Methodology
 - Association Measures
 - Evaluation
- 3 Combining Association Measures
 - Classification and Ranking
 - Attribute Selection
- 4 Summary

Summary

Achieved results

- Empirical **evaluation** of 84 association measures.
Pointwise mutual information evaluated as one of the best measures.
- Statistical **combination** of multiple association measures.
Linear logistic regression gives significant performance improvement.
- **Selection** of the best subset of association measures.
Greedy algorithm reduced number of association measures to 17.

Outlook

- Multiple annotation of the reference data.
- Employing other classification (ranking) methods.

That's all folks ...

Thank you!

Association Measures I

1. Mean component offset	$\frac{1}{n} \sum_{i=1}^n d_i$
2. Variance component offset	$\frac{1}{n-1} \sum_{i=1}^n (d_i - \bar{d})^2$
3. Joint probability	$P(xy)$
4. Conditional probability	$P(y x)$
5. Reverse conditional prob.	$P(x y)$
6. Pointwise mutual inform.	$\log \frac{P(xy)}{P(x)P(*y)}$
7. Mutual dependency (MD)	$\log \frac{P(xy)^2}{P(x*)P(*y)}$
8. Log frequency biased MD	$\log \frac{P(xy)^2}{P(x*)P(*y)} + \log P(xy)$
9. Normalized expectation	$\frac{2f(xy)}{f(x*)+f(*y)}$
10. Mutual expectation	$\frac{2f(xy)}{f(x)+f(*y)} \cdot P(xy)$
11. Saliency	$\log \frac{P(xy)^2}{P(x*)P(*y)} \cdot \log f(xy)$
12. Pearson's χ^2 test	$\sum_{ij} \frac{(f_{ij} - \hat{f}_{ij})^2}{\hat{f}_{ij}}$
13. Fisher's exact test	$\frac{f(x*)!f(\bar{x}*)!f(*y)!f(*\bar{y})!}{N!f(xy)!f(x\bar{y})!f(\bar{x}y)!f(\bar{x}\bar{y})!}$
14. t test	$\frac{f(xy) - \hat{f}(xy)}{\sqrt{f(xy)(1 - (f(xy)/N))}}$
15. z score	$\frac{f(xy) - \hat{f}(xy)}{\sqrt{\hat{f}(xy)(1 - (\hat{f}(xy)/N))}}$
16. Poisson significance measure	$\frac{\hat{f}(xy) - f(xy) \log \hat{f}(xy) + \log f(xy)!}{\log N}$

Association Measures II

17. Log likelihood ratio	$-2 \sum_{ij} f_{ij} \log \frac{f_{ij}}{f_{i.} f_{.j}}$
18. Squared log likelihood ratio	$-2 \sum_{ij} \frac{\log^2 f_{ij}}{f_{ij}}$
Association coefficients:	
19. Russel-Rao	$\frac{a}{a+b+c+d}$
20. Sokal-Michiner	$\frac{a+d}{a+b+c+d}$
*21. Rogers-Tanimoto	$\frac{a+d}{a+2b+2c+d}$
22. Hamann	$\frac{(a+d)-(b+c)}{a+b+c+d}$
23. Third Sokal-Sneath	$\frac{b+c}{a+d}$
24. Jaccard	$\frac{a}{a+b+c}$
*25. First Kulczynsky	$\frac{a}{b+c}$
26. Second Sokal-Sneath	$\frac{a}{a+2(b+c)}$
27. Second Kulczynski	$\frac{1}{2} \left(\frac{a}{a+b} + \frac{a}{a+c} \right)$
28. Fourth Sokal-Sneath	$\frac{1}{4} \left(\frac{a}{a+b} + \frac{a}{a+c} + \frac{d}{d+b} + \frac{d}{d+c} \right)$
29. Odds ratio	$\frac{ad}{bc}$
30. Yulle's ω	$\frac{\sqrt{ad} - \sqrt{bc}}{\sqrt{ad} + \sqrt{bc}}$
*31. Yulle's Q	$\frac{ad - bc}{ad + bc}$
32. Driver-Kroeber	$\frac{a}{\sqrt{(a+b)(a+c)}}$

Association Measures III

33. Fifth Sokal-Sneath	$\frac{ad}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$
34. Pearson	$\frac{ad - bc}{\sqrt{(a+b)(a+c)(d+b)(d+c)}}$
35. Baroni-Urbani	$\frac{a + \sqrt{ad}}{a + b + c + \sqrt{ad}}$
36. Braun-Blanquet	$\frac{a}{\max(a+b, a+c)}$
37. Simpson	$\frac{a}{\min(a+b, a+c)}$
38. Michael	$\frac{4(ad - bc)}{(a+d)^2 + (b+c)^2}$
39. Mountford	$\frac{2a}{2bc + ab + ac}$
40. Fager	$\frac{a}{\sqrt{(a+b)(a+c)}} - \frac{1}{2} \max(b, c)$
41. Unigram subtuples	$\log \frac{ad}{bc} - 3.29 \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$
42. U cost	$\log \left(1 + \frac{\min(b, c) + a}{\max(b, c) + a} \right)$
43. S cost	$\log \left(1 + \frac{\min(b, c)}{a+1} \right) - \frac{1}{2}$
44. R cost	$\log \left(1 + \frac{a}{a+b} \right) \cdot \log \left(1 + \frac{a}{a+c} \right)$
45. T combined cost	$\sqrt{U \times S \times \bar{R}}$
46. Phi	$\frac{P(xy) - P(x^*)P(*y)}{\sqrt{P(x^*)P(*y)(1 - P(x^*)) (1 - P(*y))}}$
47. Kappa	$\frac{P(xy) + P(\bar{x}\bar{y}) - P(x^*)P(*y) - P(\bar{x}^*)P(*\bar{y})}{1 - P(x^*)P(*y) - P(\bar{x}^*)P(*\bar{y})}$
48. J measure	$\max \left[P(xy) \log \frac{P(y x)}{P(*y)} + P(x\bar{y}) \log \frac{P(\bar{y} x)}{P(*\bar{y})}, \right. \\ \left. P(xy) \log \frac{P(x y)}{P(x^*)} + P(\bar{x}y) \log \frac{P(\bar{x} y)}{P(\bar{x}^*)} \right]$

Association Measures IV

49. Gini index	$\begin{aligned} & \max[P(x^*)(P(y x)^2 + P(\bar{y} x)^2) - P(*y)^2 \\ & \quad + P(x^*)(P(y \bar{x})^2 + P(\bar{y} \bar{x})^2) - P(*\bar{y})^2, \\ & \quad P(*y)(P(x y)^2 + P(\bar{x} y)^2) - P(x^*)^2 \\ & \quad + P(*\bar{y})(P(x \bar{y})^2 + P(\bar{x} \bar{y})^2) - P(\bar{x}^*)^2] \end{aligned}$
50. Confidence	$\max[P(y x), P(x y)]$
51. Laplace	$\max\left[\frac{NP(xy)+1}{NP(x^*)+2}, \frac{NP(xy)+1}{NP(*y)+2}\right]$
52. Conviction	$\max\left[\frac{P(x^*)P(*y)}{P(\bar{x}\bar{y})}, \frac{P(\bar{x}^*)P(*y)}{P(\bar{x}\bar{y})}\right]$
53. Piatersky-Shapiro	$P(xy) - P(x^*)P(*y)$
54. Certainty factor	$\max\left[\frac{P(y x) - P(*y)}{1 - P(*y)}, \frac{P(x y) - P(x^*)}{1 - P(x^*)}\right]$
55. Added value (AV)	$\max[P(y x) - P(*y), P(x y) - P(x^*)]$
56. Collective strength	$\frac{P(xy) + P(\bar{x}\bar{y})}{P(x^)P(y) + P(\bar{x}^*)P(*y)} \cdot \frac{1 - P(x^*)P(*y) - P(\bar{x}^*)P(*y)}{1 - P(xy) - P(\bar{x}\bar{y})}$
57. Klosgen	$\sqrt{P(xy)} \cdot AV$
Context measures:	
*58. Context entropy	$-\sum_w P(w C_{xy}) \log P(w C_{xy})$
59. Left context entropy	$-\sum_w P(w C_{xy}^l) \log P(w C_{xy}^l)$
60. Right context entropy	$-\sum_w P(w C_{xy}^r) \log P(w C_{xy}^r)$

Association Measures V

* 61. Left context divergence	$P(x^*) \log P(x^*) - \sum_w P(w C_{xy}^l) \log P(w C_{xy}^l)$
62. Right context divergence	$P(*y) \log P(*y) - \sum_w P(w C_{xy}^r) \log P(w C_{xy}^r)$
63. Cross entropy	$-\sum_w P(w C_x) \log P(w C_y)$
64. Reverse cross entropy	$-\sum_w P(w C_y) \log P(w C_x)$
65. Intersection measure	$\frac{2 C_x \cap C_y }{ C_x + C_y }$
66. Euclidean norm	$\sqrt{\sum_w (P(w C_x) - P(w C_y))^2}$
67. Cosine norm	$\frac{\sum_w P(w C_x) P(w C_y)}{\sqrt{\sum_w P(w C_x)^2 \cdot \sum_w P(w C_y)^2}}$
68. L1 norm	$\sum_w P(w C_x) - P(w C_y) $
69. Confusion probability	$\sum_w \frac{P(x C_w) P(y C_w) P(w)}{P(x^*)}$
70. Reverse confusion prob.	$\sum_w \frac{P(y C_w) P(x C_w) P(w)}{P(*y)}$
* 71. Jensen-Shannon diverg.	$\frac{1}{2} [D(p(w C_x) \frac{1}{2}(p(w C_x) + p(w C_y))) + D(p(w C_y) \frac{1}{2}(p(w C_x) + p(w C_y)))]$
72. Cosine of pointwise MI	$\frac{\sum_w M(w, x) M(w, y)}{\sqrt{\sum_w M(w, x)^2} \cdot \sqrt{\sum_w M(w, y)^2}}$
* 73. KL divergence	$\sum_w P(w C_x) \log \frac{P(w C_x)}{P(w C_y)}$
* 74. Reverse KL divergence	$\sum_w P(w C_y) \log \frac{P(w C_y)}{P(w C_x)}$

Association Measures VI

75. **Skew divergence**

$$D(p(w|C_X) || \alpha p(w|C_Y) + (1-\alpha)p(w|C_X))$$

76. **Reverse skew divergence**

$$D(p(w|C_Y) || \alpha p(w|C_X) + (1-\alpha)p(w|C_Y))$$

77. **Phrase word cocurrence**

$$\frac{1}{2} \left(\frac{f(x|C_{XY})}{f(xy)} + \frac{f(y|C_{XY})}{f(xy)} \right)$$

78. **Word association**

$$\frac{1}{2} \left(\frac{f(x|C_Y) - f(xy)}{f(xy)} + \frac{f(y|C_X) - f(xy)}{f(xy)} \right)$$

Cosine context similarity:

$$\frac{1}{2} (\cos(\mathbf{c}_X, \mathbf{c}_{XY}) + \cos(\mathbf{c}_Y, \mathbf{c}_{XY}))$$

$$\mathbf{c}_Z = (z_i); \cos(\mathbf{c}_X, \mathbf{c}_Y) = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2} \cdot \sqrt{\sum y_i^2}}$$

*79. **in boolean vector space**

$$z_i = \delta(f(w_j|C_Z))$$

80. **in *f* vector space**

$$z_i = f(w_j|C_Z)$$

81. **in *f*·*idf* vector space**

$$z_i = f(w_j|C_Z) \cdot \frac{N}{df(w_j)}; df(w_j) = |\{x : w_j \in C_X\}|$$

Dice context similarity:

$$\frac{1}{2} (\text{dice}(\mathbf{c}_X, \mathbf{c}_{XY}) + \text{dice}(\mathbf{c}_Y, \mathbf{c}_{XY}))$$

$$\mathbf{c}_Z = (z_i); \text{dice}(\mathbf{c}_X, \mathbf{c}_Y) = \frac{2 \sum x_i y_i}{\sum x_i^2 + \sum y_i^2}$$

*82. **in boolean vector space**

$$z_i = \delta(f(w_j|C_Z))$$

*83. **in *f* vector space**

$$z_i = f(w_j|C_Z)$$

*84. **in *f*·*idf* vector space**

$$z_i = f(w_j|C_Z) \cdot \frac{N}{df(w_j)}; df(w_j) = |\{x : w_j \in C_X\}|$$

*85. **Part of speech**

{Adjective:Noun, Noun:Noun, Noun:Verb, ...}

*86. **Dependency type**

{Attribute, Object, Subject, ...}

87. **Dependency structure**

{↗, ↖}