

Combining Association Measures for Collocation Extraction

Pavel Pecina & Pavel Schlesinger

Institute of Formal and Applied Linguistics
Charles University, Prague, Czech Republic

{pecina, schlesinger}@ufal.mff.cuni.cz

Outline

- ▶ Notion of collocation
 - ... with definitions, characteristic properties, and examples.
- ▶ Manual annotation of reference data
 - ... of reasonable size and quality.
- ▶ Collocation extraction and evaluation of basic methods
 - ... by means of precision-recall curves and mean average precision.
- ▶ Combining association measures for collocation extraction
 - ... to improve results of individual measures.
- ▶ Reduce number of combined measures
 - ... to remove redundant and ineffective measures.

Notion of Collocation

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.”

Evert (2004):

“A **collocation** is a word combination whose semantic and/or syntactic properties cannot be fully predicted from those of its components, and which therefore has to be listed in a lexicon.”

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 of Collocations

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)

Reference Data Annotation

Source: Prague Dependency Treebank 2.0 (new version of the Czech treebank)

Sentences: 87,980 (morphological and analytical annotation)

Word forms: 1,504,847

Dependency bigram types: 635,952

Reference bigram types (f>S): 26,450 (frequency filtering)

Reference collocation candidates: 12,232 (Part-of-Speech filtering)

Focus on bigram collocations:

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

Manual annotation:

- ▶ The list of collocation candidates processed by three linguists in parallel.
- ▶ Bigrams that all three annotators independently recognized as collocations (of any type) were considered true collocations (20.9%).

Collocation categories and types

- ▶ idiomatic expressions
 - studená válka (cold war), natáhnout bačkory (kick the bucket)
- ▶ technical terms
 - předseda vlády (prime minister), očítý svědek (eye witness)
- ▶ support verb constructions
 - mít pravdu (to be right), učinit rozhodnutí (make decision)
- ▶ names of persons, locations, and other entities
 - Pražský hrad (Prague Castle), Červený kříž (Red Cross)
- ▶ stock phrases
 - zásadní problém (major problem), konec roku (end of the year)

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 and can be used for **ranking** (highest to the top) or **classification** (by setting a threshold).

Ranking	Score
red cross	15.66
decimal point	14.01
arithmetic operation	10.52
paper feeder	10.17
system type	3.54
and others	0.54
program in level is	0.25

Classification	Score
red cross	1
decimal point	1
arithmetic operation	1
paper feeder	1
system type	0
and others	0
program in level is	0

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:

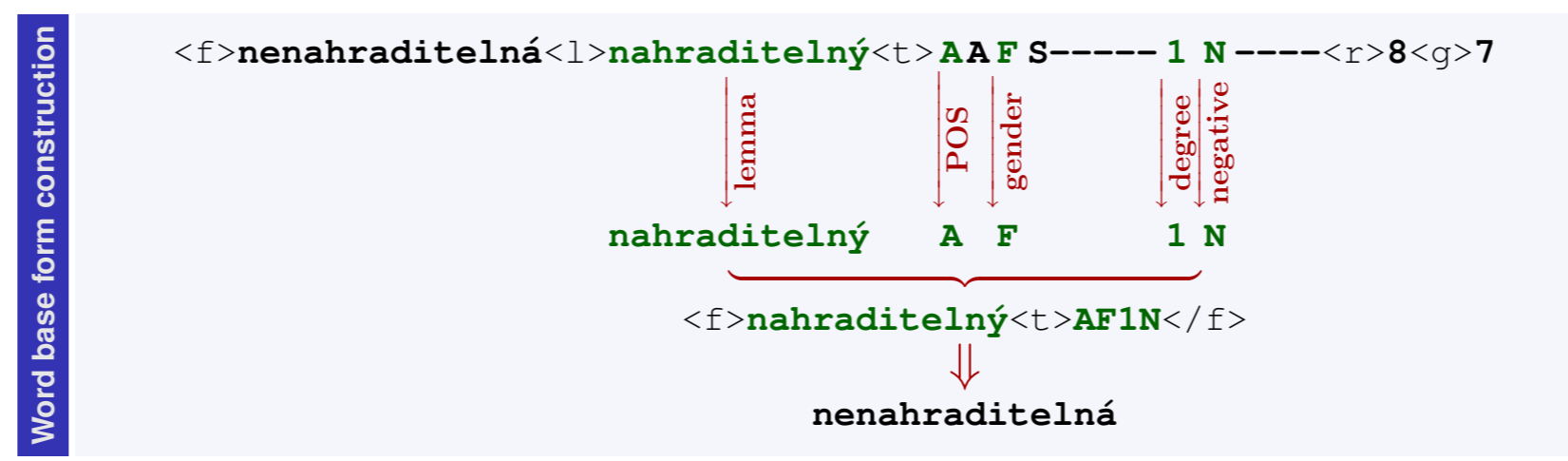
- ▶ frequency of word and n-gram occurrences
- ▶ immediate contexts
- ▶ empirical contexts

4. Computing association measures

5. Ranking or classification

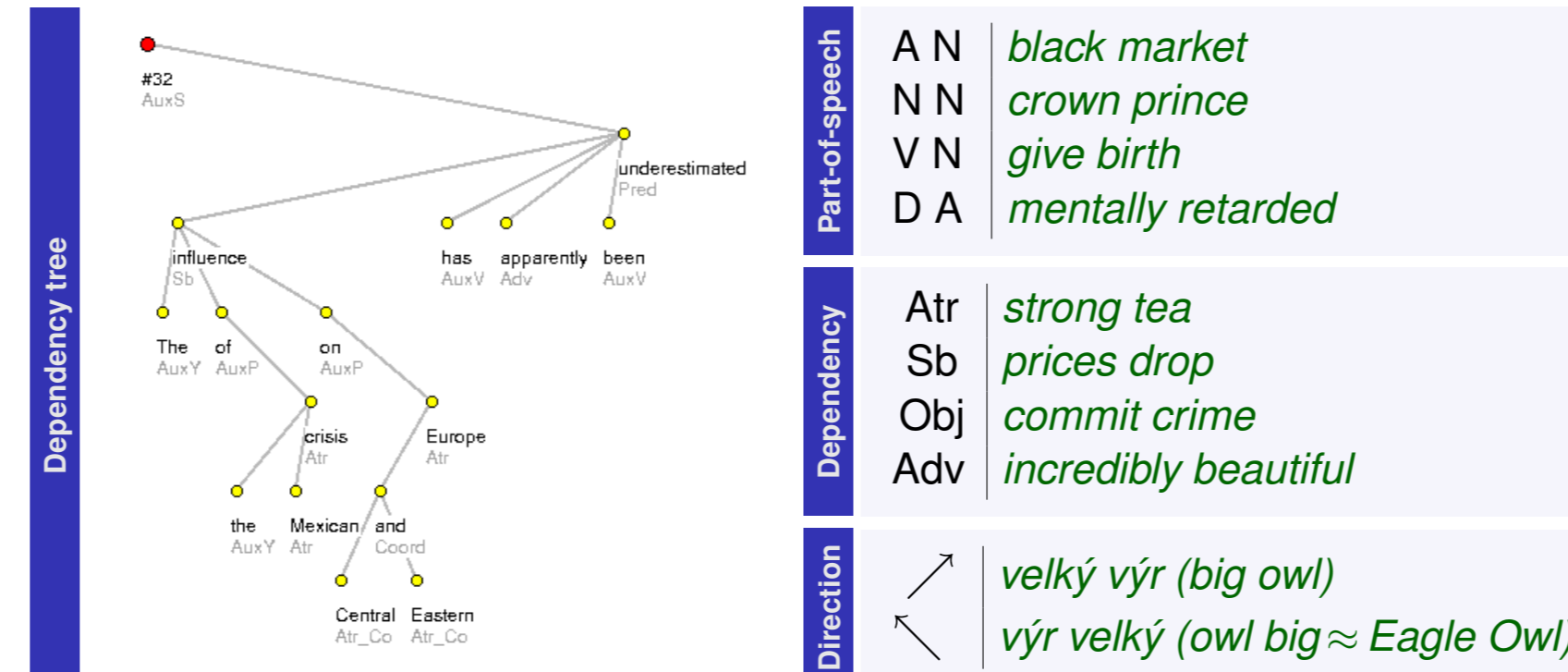
Word Base Forms

- ▶ Surface word forms too specific (rich Czech morphology)
- ▶ Pure lemmas too general (loss of syntactic and semantic information)
- ▶ Lemmas with a subset of morphological tag fairly optimal



Dependency Bigrams

- ▶ Dependency trees broken down to dependency bigrams consisting of word base forms, dependency type and direction:

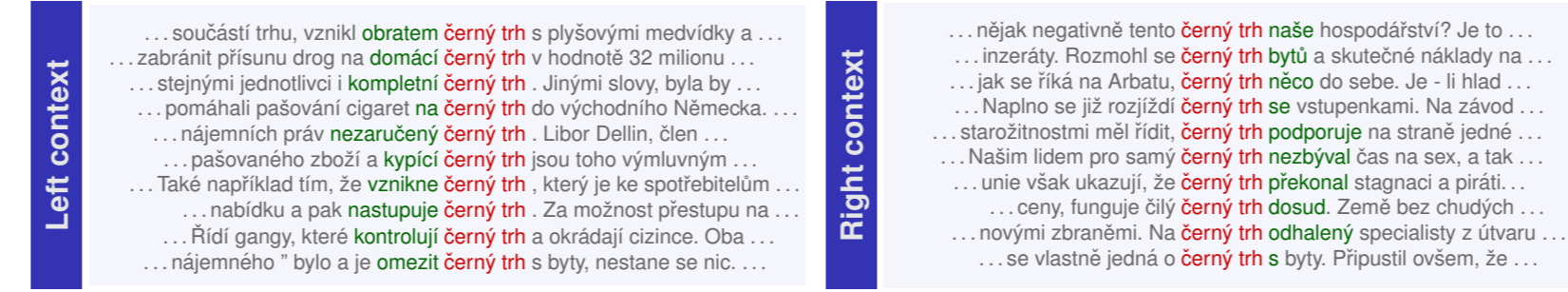


Cooccurrence and Context Statistics

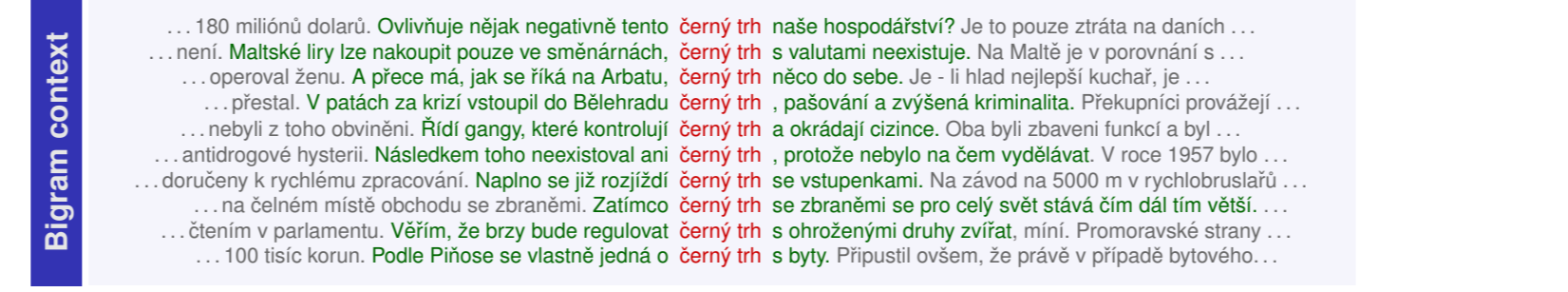
- ▶ Contingency tables – observed frequencies and marginal frequencies:

Contingency table		Example	
$f(x,y)$	$f(x,\bar{y})$	$f(x^*,y)$	$f(x^*,\bar{y})$
$f(\bar{x},y)$	$f(\bar{x},\bar{y})$	$f(\bar{x}^*,y)$	$f(\bar{x}^*,\bar{y})$
$f(*,y)$	$f(*,\bar{y})$	$f(*,y)$	$f(*,\bar{y})$
		N	N

- ▶ Immediate context – words immediately preceding or following the bigram:



- ▶ Empirical context – words occurring within a specified context window:



Collocation Hypotheses and Types of Association Measures

H₁: “Collocations are very frequent word combinations.”

- ▶ ML estimations of joint and conditional probabilities

H₂: “Collocation components occur together more often than by a chance.”

- ▶ Mutual information and derived measures
- ▶ Statistical tests of independence
- ▶ Likelihood measures

H₃: “Collocations occur as units in a (inf.-theoretically) noisy environment.”

- ▶ Information-theory measures of immediate context

H₄: “Collocations occur in different contexts than their components.”

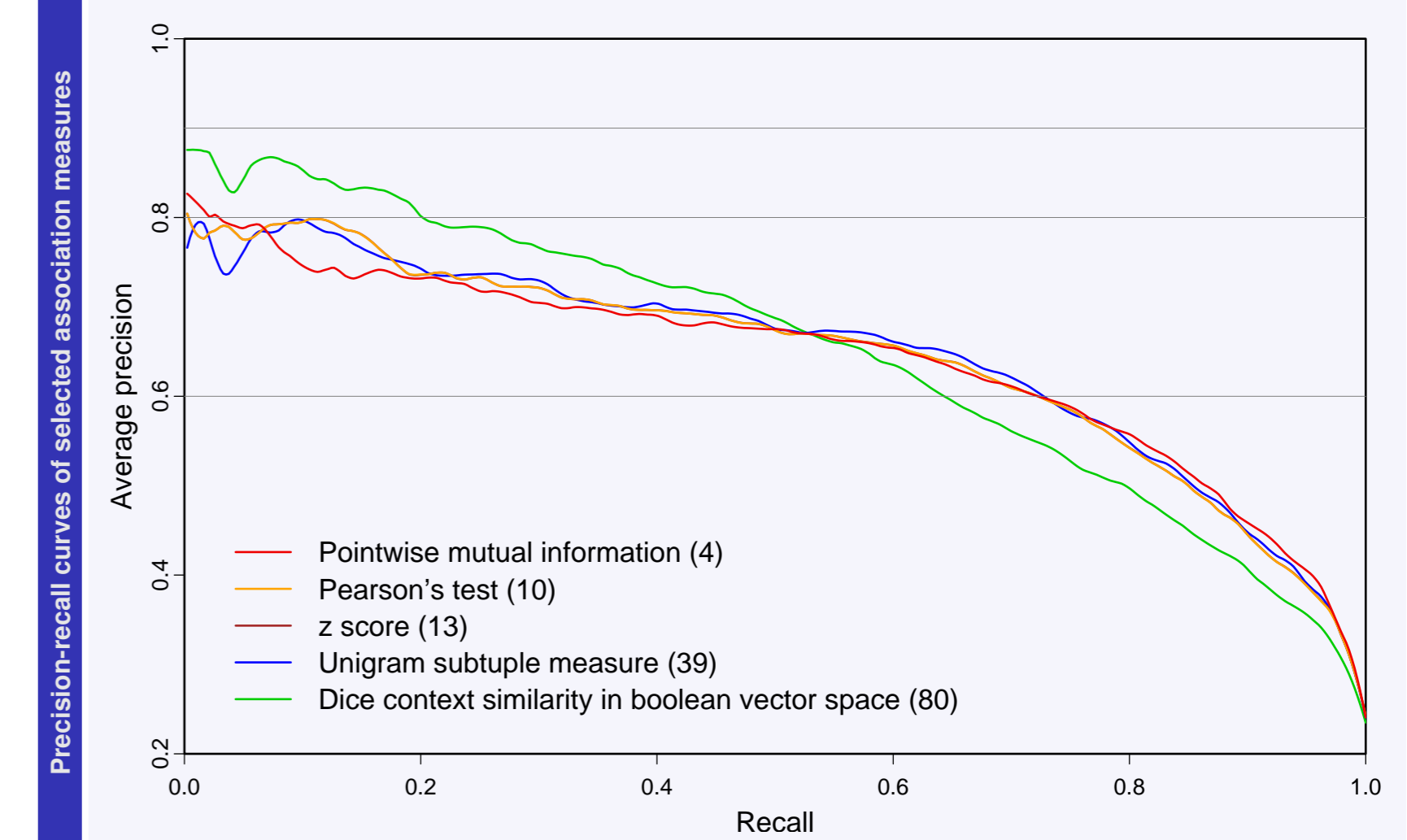
- ▶ Information-theory context measures
- ▶ Information-retrieval context similarity measures

Association Measures

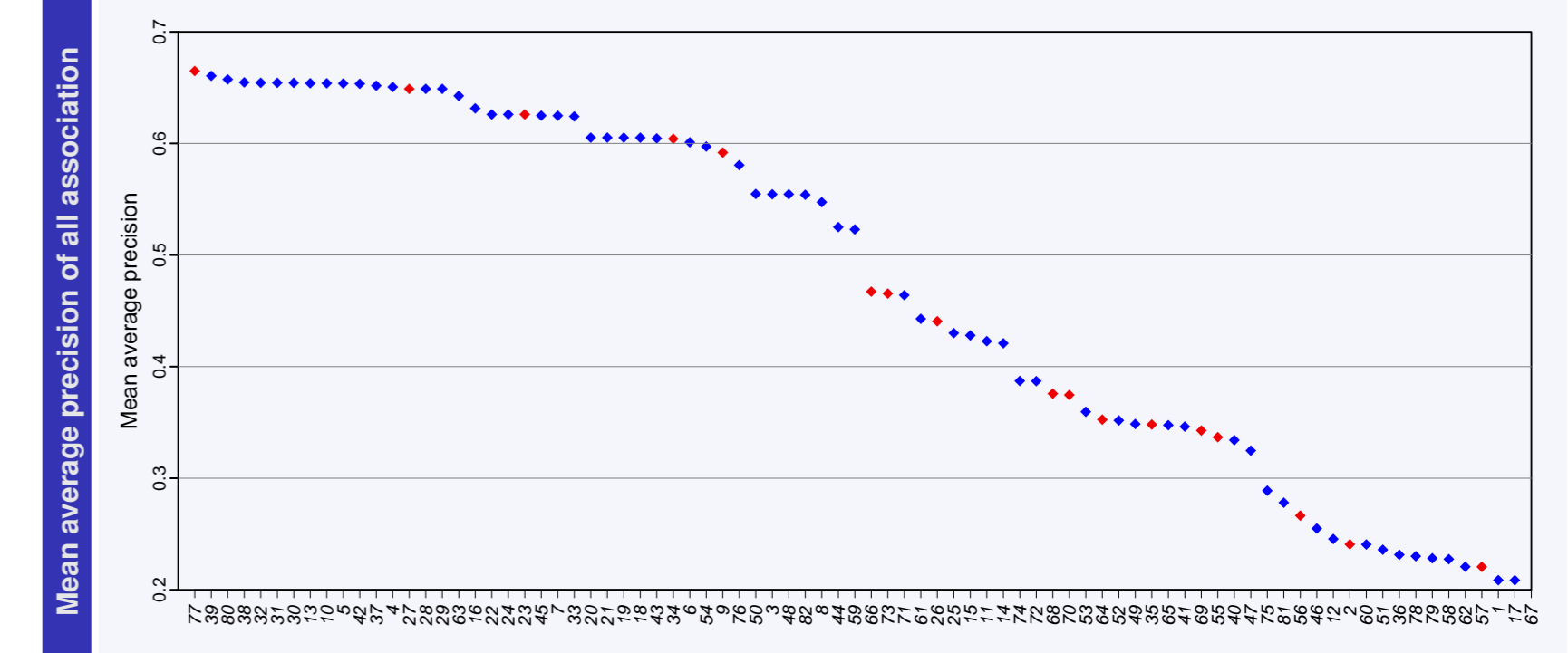
1. Joint probability	$P(x,y)$	47. Gini index	$\max\{P(x^+)(P(y^-)^2 + P(y^+)^2) - P(x^+)^2\}$
2. Conditional probability	$P(y x)$	48. Confidence	$\frac{P(x,y)}{P(x)}$
3. Reverse conditional prob.	$P(x y)$	49. Laplace	$\frac{P(x,y)}{P(y)}$
4. Pointwise mutual inform.	$-\log \frac{P(x,y)}{P(x)P(y)}$	50. Conviction	$\frac{P(x,y)}{P(x)P(y)}$
5. Mutual dependency (MD)	$-\log \frac{P(x,y)}{P(x)P(y)}$	51. Piatersky-Shapiro	$P(x,y) - P(x)P(y)$
6. Log frequency biased MD	$\log \frac{P(x,y)}{P(x)P(y)}$	52. Certainty factor	$\frac{P(x,y) - P(x)P(y)}{1 - P(x)P(y)}$
7. Normalized expectation	$\frac{P(x,y)}{P(x)}$	53. Added value (AV)	$\frac{P(x,y) - P(x)P(y)}{P(x)}$
8. Mutual expectation	$\frac{P(x,y)}{P(x)}$	54. Collective strength	$\frac{P(x,y) - P(x)P(y)}{P(x)}$
9. Salience	$\log \frac{P(x,y)}{P(x)P(y)}$	55. Klossen	$\frac{P(x,y) - P(x)P(y)}{P(x)}$
10. Pearson's χ^2 test	$\sum \frac{(f_{ij} - f_{ij}^e)^2}{f_{ij}^e}$	56. Context entropy	$-\sum P(w_i C_j) \log P(w_i C_j)$
11. Fisher's exact test	$\frac{f_{11}f_{22} - f_{12}f_{21}}{f_{1+}f_{2+}f_{+1}f_{+2}}$	57. Left context entropy	$-\sum P(C_j w_i) \log P(C_j w_i)$
12. t test	$\frac{f_{11} - f_{12}}{\sqrt{\frac{f_{1+}f_{2+}}{N}}}$	58. Right context entropy	$-\sum P(C_j w_i) \log P(C_j w_i)$
13. z score	$\frac{f_{11} - f_{12}}{\sqrt{\frac{f_{1+}f_{2+}}{N}}}$	59. Left context divergence	$-\sum P(w_i C_j) \log P(w_i C_j)$
14. Poisson significance meas.	$\frac{f_{11} - f_{12}}{\sqrt{f_{1+} + f_{2+}}}$	60. Right context divergence	$-\sum P(w_i C_j) \log P(w_i C_j)$
15. Log likelihood ratio	$-2 \sum \frac{f_{ij} \log \frac{f_{ij}}{f_{ij}^e}}{f_{ij}}$	61. Cross entropy	$-\sum P(w_i C_j) \log P(w_i C_j)$
16. Squared log likelihood rat.	$-2 \sum \frac{f_{ij} \log \frac{f_{ij}}{f_{ij}^e}}{f_{ij}}$	62. Reverse cross entropy	$-\sum P(w_i C_j) \log P(w_i C_j)$
Association coefficients:		63. Intersection measure	$\frac{P(x,y)}{P(x)P(y)}$
17. Russel-Rao	$\frac{1}{1 + \sqrt{1 - 2P(x,y)}}$	64. Euclidean norm	$\sqrt{\sum (P(w_i C_j) - P(w_i))^2}$
18. Sokal-Michener	$\frac{1 - 2P(x,y)}{1 + \sqrt{1 - 2P(x,y)}}$	65. Cosine norm	$\frac{P(x,y)}{\sqrt{P(x)P(y)}}$
19. Rogers-Tanimoto	$\frac{2P(x,y)}{1 + P(x,y)}$	66. L1 norm	$\sum P(w_i C_j) - P(w_i) $
20. Hamann	$\frac{(a+b) - 2c}{a+b+c}$	67. Confusion probability	$\sum \frac{P(x_i, C_j)P(x_i, C_k)}{P(x_i)}$
21. Third Sokal-Sneath	$\frac{2c}{a+b+c}$	68. Reverse confusion prob.	$\sum \frac{P(x_i, C_j)P(x_i, C_k)}{P(x_i)}$
22. Jaccard	$\frac{c}{a+b+c}$	69. Jensen-Shannon diverg.	$\frac{1}{2} (D(p_i(w_i C_j) p_i(w_i C_k)) + D(p_i(w_i C_k) p_i(w_i C_j)))$
23. First Kulczynsky	$\frac{c}{a+b}$	70. Cosine of pointwise MI	$\frac{P(x,y)}{P(x)P(y)}$
24. Second Sokal-Sneath	$\frac{2c}{a+b+c}$	71. KL divergence	$\sum P(w_i C_j) \log \frac{P(w_i C_j)}{P(w_i)}$
25. Second Kulczynsky	$\frac{2c}{a+b}$	72. Reverse KL divergence	$\sum P(w_i C_j) \log \frac{P(w_i C_j)}{P(w_i)}$
26. Fourth Sokal-Sneath	$\frac{2c}{a+b+c}$	73. Skew divergence	$D(p_i(w_i C_j) p_i(w_i C_k) + (1-\alpha)p_i(w_i C_j))$
27. Odds ratio	$\frac{ad}{bc}$	74. Reverse skew divergence	$D(p_i(w_i C_j) p_i(w_i C_k) + (1-\alpha)p_i(w_i C_j))$
28. Yulle's u	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	75. Phrase word cooccurrence	$\frac{P(x,y)}{P(x)P(y)}$
29. Yulle's Q	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	76. Word association	$\frac{1}{2} (\frac{f_{11}f_{22}}{f_{1+}f_{2+}} + \frac{f_{12}f_{21}}{f_{1+}f_{2+}})$
30. Fisher-Kroeber	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	77. In if vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) + \cos(\epsilon_{y,x}))$
31. Driver-Sneath	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	78. in if vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) - \cos(\epsilon_{y,x}))$
32. Pearson	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	79. in-if-if vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) - \cos(\epsilon_{y,x}))$
33. Baroni-Urbani	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	80. in boolean vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) - \cos(\epsilon_{y,x}))$
34. Braun-Blanquet	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	81. in if vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) - \cos(\epsilon_{y,x}))$
35. Simpson	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$	82. in-if-if vector space	$\frac{1}{2} (\cos(\epsilon_{x,y}) - \cos(\epsilon_{y,x}))$
36. Michael	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		
37. Mountford	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		
38. Fager	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		
39. Unigram subtuples	$\log \frac{(a+b)(c+d) + 2c}{(a+b)(c+d)}$		
40. U cost	$\log(1 + \frac{ad - bc}{(a+b)(c+d)})$		
41. S cost	$\log(1 + \frac{ad - bc}{(a+b)(c+d)})$		
42. R cost	$\log(1 + \frac{ad - bc}{(a+b)(c+d)})$		
43. T combined cost	$\sqrt{T \times S \times R}$		
44. Phi	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		
45. Kappa	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		
46. J measure	$\frac{ad - bc}{\sqrt{(a+b)(c+d)(a+d)(b+c)}}$		

Evaluation of Association Measures

- ▶ Precision: fraction of positive predictions correct (given a threshold)
- ▶ Recall: fraction of positives correctly predicted (given a threshold)
- ▶ Precision–recall curve: PR scores for all possible threshold values



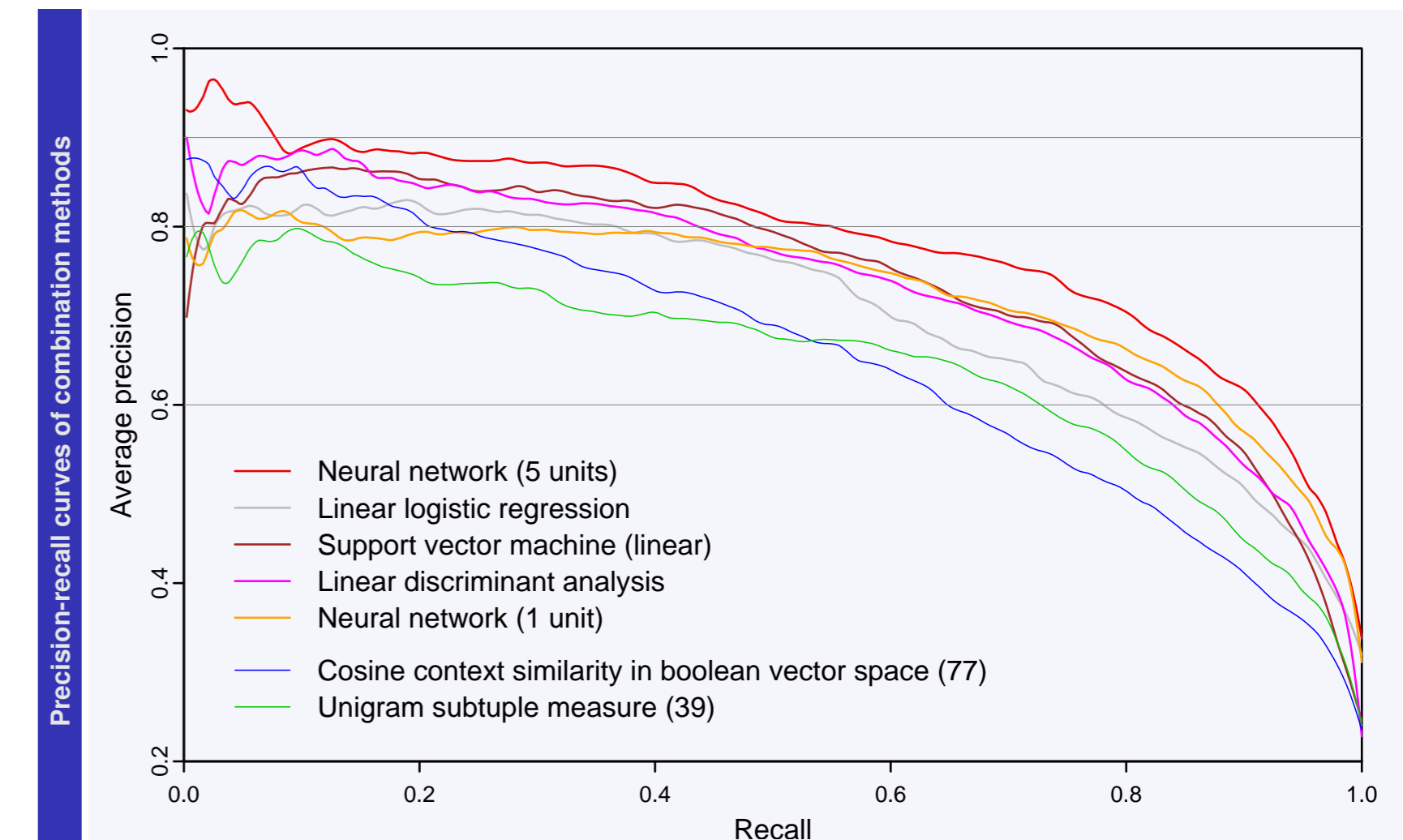
- ▶ Average precision: expected value of precision for all possible values of recall
- ▶ Mean average precision: mean of AP computed for crossvalidation data folds



Combining Association Measures

- ▶ Each collocation candidate x^i is described by the **feature vector** $x^i = (x_1^i, \dots, x_n^i)^T$ consisting of scores of all association measures.
- ▶ We look for a **ranker** function $f(x) \rightarrow \mathbb{R}$ that determines the strength of lexical association between components of bigram x :

- ▶ Linear logistic regression
- ▶ Linear discriminant analysis
- ▶ Support vector machines
- ▶ Neural networks



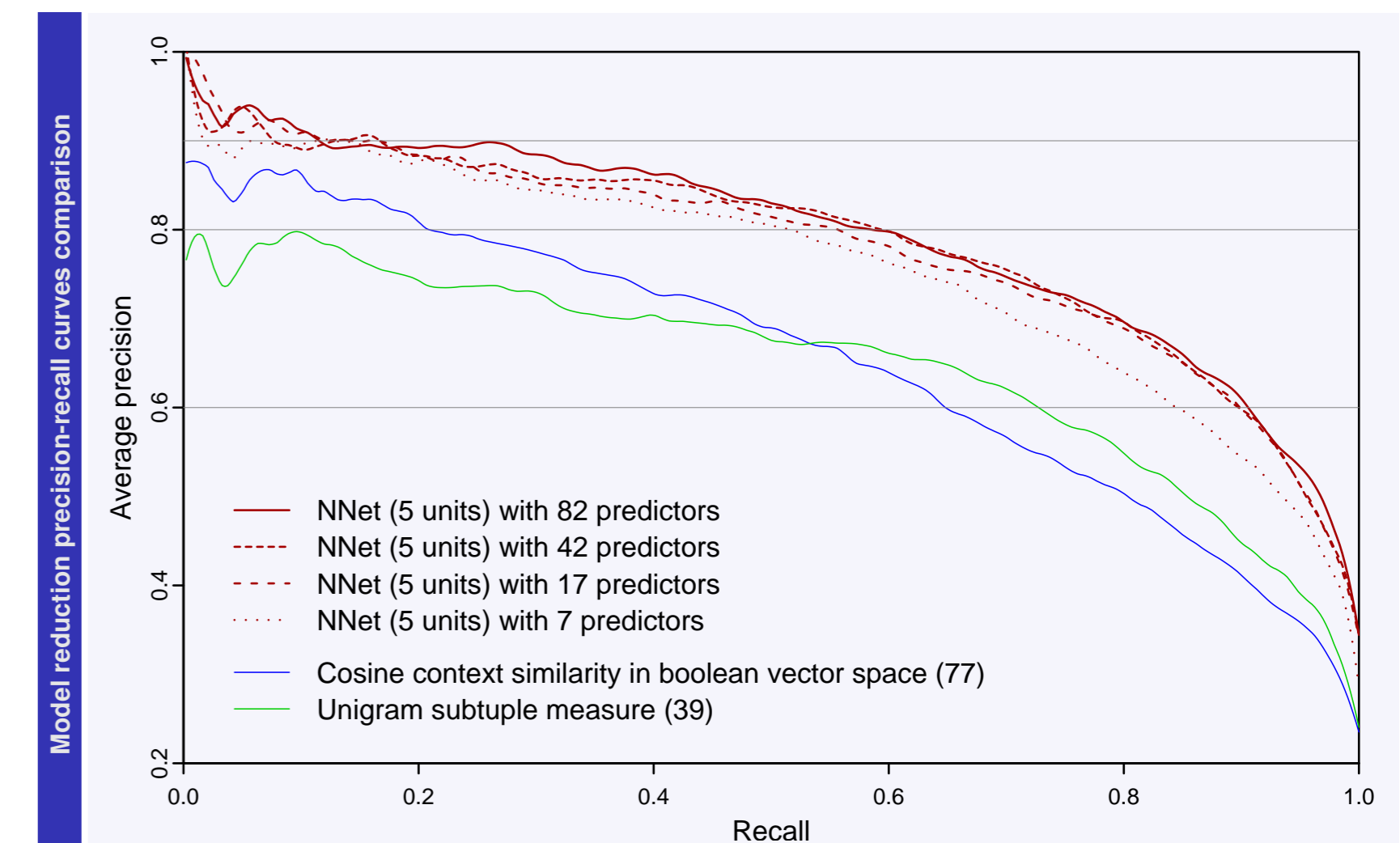
method	R=20	AP	R=50	R=80	MAP
Neural Network (5 units)	89.56	82.74	70.11	80.81	21.53
Neural Network (3 units)	89.41	81.99	69.64	79.71	19.88
Neural Network (2 units)	86.92	81.68	68.33	78.77	18.47
Support Vector Machine (linear)	85.72	79.49	63.86	75.66	13.79
Linear Discriminant Analysis	84.72	77.18	62.90	75.11	12.96
Support Vector Machine (quadratic)	84.29	79.54	64.24	74.53	12.09
Neural Network (1 unit)	77.98	76.83	66.75	73.25	10.17
Logistic Linear Regression	82.45	76.26	58.61	71.88	8.11
Cosine similarity (77)	80.94	68.90	50.54	66.49	0.00
Unigram subtuples (39)	74.55	67.49	55.16	65.74	-

Model Reduction

- ▶ Combination models are too complex in number of predictors used.
- ▶ Some association measures are very similar (analytically or empirically) and as predictors perhaps even redundant. Such measures have no use in the models, make their training problematic, and should be excluded.

The model reduction algorithm:

- Starts with **hierarchical clustering** of variables in order to group those with similar contribution to the model (by **Pearson's correlation coefficient**).
- After $82-d$ iteration, the variables are grouped into d clusters, one representative from each cluster is selected as a predictor into the **initial model**.
- In each next step, the algorithm removes one predictor causing minimal degradation of performance (measured by MAP on held-out data).



Results

- ▶ The best performing individual association measures are cosine context similarity in boolean vector space and unigram subtype measure.
- ▶ The best performing combination method is Neural Network with 5 units in the hidden layer with 21.53% relative improvement in terms of MAP.
- ▶ Number of variables in this model was reduced to 17 without significant degradation of its performance.