

# Evolving New Lexical Association Measures Using Genetic Programming

Or: How to make evolution do the down and dirty work?

Jan Šnajder & Bojana Dalbelo Bašić

University of Zagreb  
Faculty of Electrical Engineering and Computing  
Department of Electronics, Microelectronics, Computer and Intelligent Systems

**MEi:CogSci Conference**  
Dubrovnik, June 17, 2010

Knowledge  
TechnologiesLab

# Computer processing of language

- Natural language processing
- Computational linguistics
- Information retrieval (IR)
- Text mining (Text analytics)

# Outline

- 1 Collocations and multiword expression
- 2 Collocation extraction
- 3 Genetic programming of AMs
- 4 Conclusion

# Outline

1 Collocations and multiword expression

2 Collocation extraction

3 Genetic programming of AMs

4 Conclusion

# Multiword expressions

Working definition (Evert, Baldwin):

A **multiword expression (MWE)** is a combination of two or more words whose semantic, syntactic, or statistical properties are not entirely predictable from those of its components, and which therefore needs to be listed in a lexicon.

- E.g. *silver bullet, jump in, kick the bucket, traffic light, blond hair*

Three characteristic properties (Manning & Schütze):

- ① non-compositionality (semantically opaque)
- ② non-substitutability (lexically determined)
- ③ non-modifiability (syntactically rigid)

## Multiword expressions (2)

MWEs cover a wide range of lexicalized expressions:

- **Idioms** (*silver bullet, black sheep, kick the bucket*)
- **Proper names** (*Humpty Dumpty, Star Wars*)
- **Terminological expressions** (*operating system, visual cortex*)
- **Phrasal verbs** (*jump in, call back*)
- **Compound nouns** (*traffic light, personal computer, value added tax*)
- **Institutional phrases and clichés** (*by and large, down and dirty*)
- ...

Important in lexicography, language learning, ...

# Collocations

Our working definition:

A combination of two or more words that **have the tendency to co-occur** (because they correspond to a conventional way of saying things).

- Co-occurrence frequency as an epiphenomenon of lexicalization
- Collocation is an **empirical notion**, MWE is a **theoretical notion**
- Mostly overlapping
  - ▶ although some collocations are not MWE (*balloon boy*)
- Interaction between linguistics, statistics and computational linguistics

## Collocations – why we need them?

- IR indexing (Vechtomova et al.; Wacholder & Songi, 2003)
- IR query expansion (Mandala et al., 2000)
- Parsing (Baldwin, 2004)
- Information extraction (Lin, 1998)
- Machine translation (Gerber & Yang, 1997)
- Natural language generation (Smadja & McKeown, 1990)
- Word sense disambiguation (Wu & Chang, 2004),
- Terminology extraction (Goldman & Wehrli, 2001)
- Document classification (Scott & Matwin, 1999)
- ...

# Collocations and manual document indexing

The screenshot shows the WinAIDE application interface. On the left, a document titled 'Pravilnik za kontrolu salmonela i drugih određenih uzročnika zoonoza koje se prenose hransom' is displayed. The title is underlined and redacted. Below the title, the text reads: 'Na temelju članka 41. stavka 6. Zakona o veterinarstvu (»Narodne novine«, broj 70/97, 105/01 i 172/03), ministar poljoprivrede, šumarstva i vodnoga gospodarstva donosi'. The date '105 25.9.2006' and the number '2362' are also present. The right side of the interface contains two tables of collocation frequency and EuroVoc descriptor reliability.

Descriptor	Reliability
zoonosis : Id. 27737	100
contagious disease : Id. 1266	99
disease prevention : Id. 5216	98
health legislation : Id. 5821	97
veterinary inspection : Id. 5612	95
animal disease : Id. 792	94
food inspection : Id. 4688	94
public health : Id. 5259	71

eCADIS – Document indexing with EUROVOC descriptors (Kolar et al., 2005)

# Outline

1 Collocations and multiword expression

2 Collocation extraction

3 Genetic programming of AMs

4 Conclusion

# Collocation extraction from corpus

- **Idea:** use statistics to extract collocations from corpus
- **N-gram** – a sequence of two or more words (bigram, trigram, tetragram, . . . )
- The simplest approach: extract most frequent n-grams
- **Lexical association measures (AMs)** assign a value to an n-gram indicating how strong the words are associated
- AMs measure the affinity of one word towards the other:  
does the n-gram occur more often than expected by chance?

Example: Vjesnik (years 1999–2009) – 56MW

$$f(\text{traffic}) = 12364, f(\text{light}) = 5741, f(\text{traffic light}) = 922$$

$$f(\text{green}) = 8395, f(\text{dolphin}) = 1067, f(\text{green dolphin}) = 2$$





Endangered *Amazon river dolphin* (also known as the *pink dolphin*)

# Lexical association measures

## ① Mutual information

$$MI(x, y) = \log \frac{P(xy)}{P(x)P(y)}$$

## ② Dice coefficient

$$Dice(x, y) = \frac{2f(xy)}{f(x) + f(y)}$$

## ③ $\chi^2$ -test

$$\chi^2 = \sum_{i,j} \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

## ④ Log-likelihood

$$G^2 = 2 \sum_{i,j} O_{ij} \log \frac{O_{ij}}{E_{ij}}$$

## Extraction example

<b>Bigram</b>	<b>AM</b>
<i>prime ministar</i>	31.5
<i>last game</i>	6.2
<i>fish soup</i>	31.0
<i>adaptation period</i>	1.7
<i>pope John</i>	22.1
<i>new judge</i>	13.3
<i>least five</i>	12.4
<i>bread recipe</i>	6.2
<i>water quality</i>	11.0
<i>consumer basket</i>	26.2
<i>six euros</i>	8.2
<i>baloon boy</i>	43.3
<i>Andy Roddick</i>	50.2
<i>big opportunity</i>	6.9

## Extraction example

<b>Bigram</b>	<b>AM</b>
<i>Andy Roddick</i>	50.2
<i>baloon boy</i>	43.3
<i>prime ministar</i>	31.5
<i>fish soup</i>	31.0
<i>consumer basket</i>	26.2
<i>pope John</i>	22.1
<i>new judge</i>	13.3
<i>least five</i>	12.4
<i>water quality</i>	11.0
<i>six euros</i>	8.2
<i>big opportunity</i>	6.9
<i>last game</i>	6.2
<i>bread recipe</i>	6.2
<i>adaptation period</i>	1.7

# Extraction procedure

- ① Tokenization
- ② Stemming/lemmatization (inflectional normalization)
  - ▶ *prime ministers* → *prime minister*
- ③ POS tagging
  - ▶ *prime minister* → AN
- ④ Token counts
- ⑤ N-gram counts
- ⑥ POS n-gram filtering
  - ▶ AN, NN, ANN, AAN, NAN, NNN, ...
- ⑦ Computing AMs for the remaining n-grams
- ⑧ Applying cutoff threshold

# Heuristic AM

- Stopword-sensitive AM for trigrams (Petrović et al., 2006):

$$MI(xyz) = \begin{cases} 2 \log \frac{P(xyz)}{P(x)P(z)} & \text{if } stop(y) \\ MI(xyz) & \text{otherwise} \end{cases}$$

- Special treatment of n-grams with stopwords (propositions and conjunctions)
- E.g. *cure for cancer, ministry of truth, bed and breakfast*

## Extraction evaluation

- It is important to evaluate!
- It is important to evaluate!
- It is important to evaluate!
- **Precision ( $P$ )** – how many of the extracted n-grams are genuine collocations?
- **Recall ( $R$ )** – how many of all the genuine collocations are extracted?
- Combined measure:

$$F_1 = \frac{2PR}{P + R}$$

- Need a hand-annotated **sample of collocations**
- A highly subjective task (even, or especially, for linguists)

# Annotation sample

županije u iznosu  
odлуku Upravnog vijeća  
terora pasa latalica  
Hrvatski filmski savez  
veterani i kadeti  
lutkarstvo Umjetničke akademije  
sistem oružanih snaga  
zemlje Srednjeg istoka  
kriteriji za dobivanje  
neposrednoj blizini mjesta  
akcijskog plana vlade  
beljskog pogona Poljoprivreda  
drugoligaši i trećeligaši  
drugih izvora finansiranja  
zgrada nije etažirana  
vrata do vrata  
sredstva za opremanje  
dogradonačelnik Petar Mlinarić  
umirovljenika Ivana Matiševa  
tijelima lokalne uprave  
lokaciju i gradnju  
mjesto u sustavu  
operirano slijepo crijevo  
izbora u koaliciji  
groblju u Vinkovcima  
činjenje takvih djela  
kuna za otkup  
posao za državu  
klub od ispadanja  
obračuna potrošnje vode  
Gradsko kazalište Joza  
predsjednik Uprave HT-a  
stožer osječkog prvoligaša  
inozemni trgovачki lanci  
nadogradnja i prenamjena  
subote ili nedjelje

# Annotation sample

županije u iznosu  
odluku Upravnog vijeća  
terora pasa latalica  
**Hrvatski filmski savez**  
veterani i kadeti  
lutkarstvo Umjetničke akademije  
**sustav oružanih snaga**  
**zemlje Srednjeg istoka**  
kriteriji za dobivanje  
neposrednoj blizini mjesta  
**akcijskog plana vlade**  
beljskog pogona Poljoprivreda  
drugoligaši i trećeligaši  
drugi izvora financiranja  
zgrada nije etažirana  
**vrata do vrata**  
sredstva za opremanje  
**dogradonačelnik Petar Mlinarić**

umirovljenika Ivana Matiševa  
**tijelima lokalne uprave**  
lokaciju i gradnja  
mjesto u sustavu  
operirano slijepo crijevo  
izbora u koaliciji  
groblju u Vinkovcima  
činjenje takvih djela  
kuna za otkup  
posao za državu  
klub od ispadanja  
obračuna potrošnje vode  
Gradsko kazalište Joza  
**predsjednik Uprave HT-a**  
stožer osječkog prvoligaša  
**inozemni trgovački lanci**  
nadogradnja i prenamjena  
subote ili nedjelje

# Outline

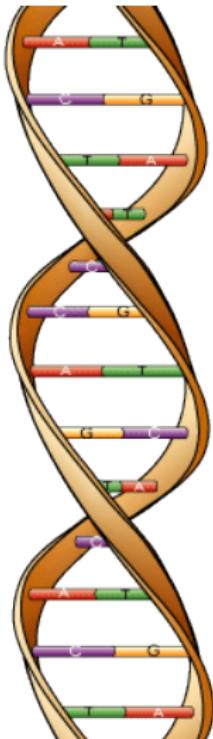
1 Collocations and multiword expression

2 Collocation extraction

3 Genetic programming of AMs

4 Conclusion

# Genetic programming



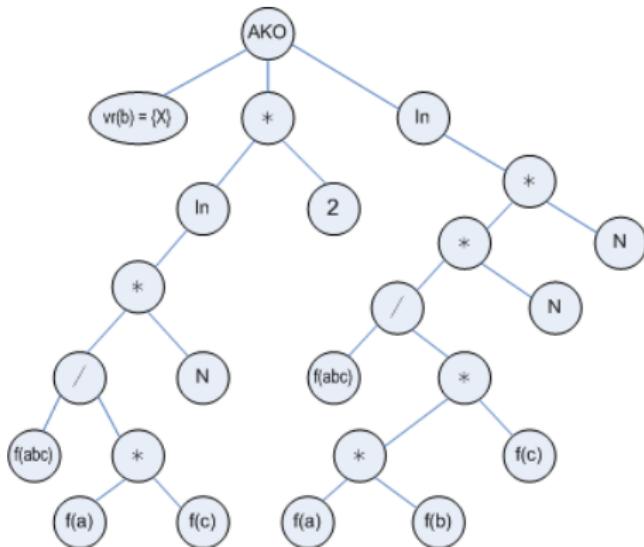
- **Genetic algorithm** – computational technique that mimics biological evolution for the purpose of solving complex optimization problems
  - ▶ population of chromosomes (solutions)
  - ▶ fitness function
  - ▶ selection
  - ▶ crossover
  - ▶ mutation
- Stochastic search through vast solution space
- **Genetic programming** – use of genetic algorithms to evolve computer programs
- Computer programs are **tree-like data structures**

# Evolving AMs

Idea: use genetic programming to **evolve new AMs** (Šnajder et al., 2008)

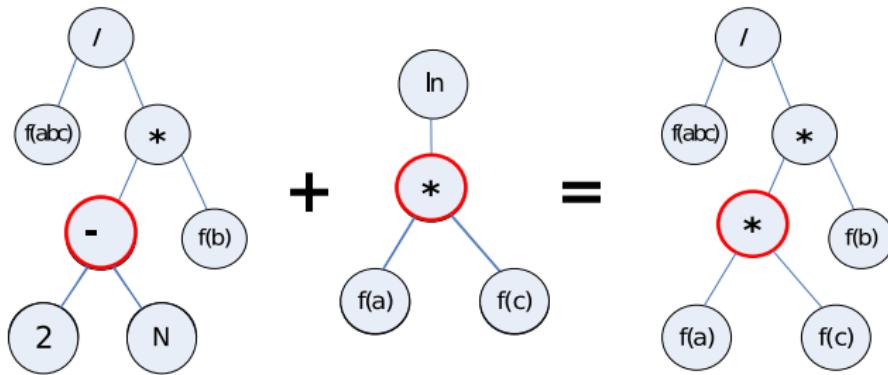
## Chromosomes = AMs

- Leaves:
  - ▶ constant
  - ▶ n-gram frequency
  - ▶ POS
- Inner nodes:
  - ▶ arithmetic operator:  
+, -, log
  - ▶ conditional branching:  
“if POS is T then...”



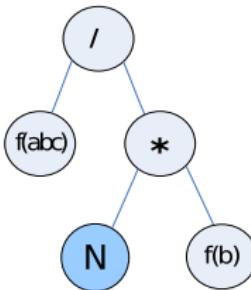
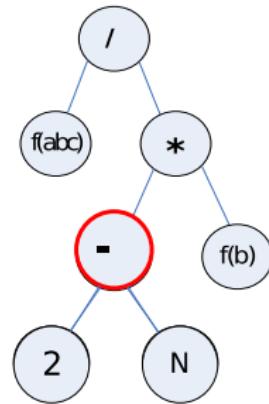
## AM crossover

- **Exchange** of genetic material (exploitation principle)
- Two AMs (the parents) are combined into a new AM (the child)
- **Swap** randomly chosen subtrees of parent solutions

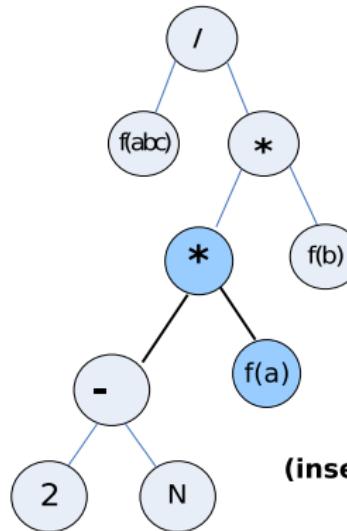


# AM mutation

- **Introduce** new genetic material (exploration principle)
- **Remove** a randomly selected node (25% chance)
- **Insert** a random node at a random position (75% chance)



(removal)



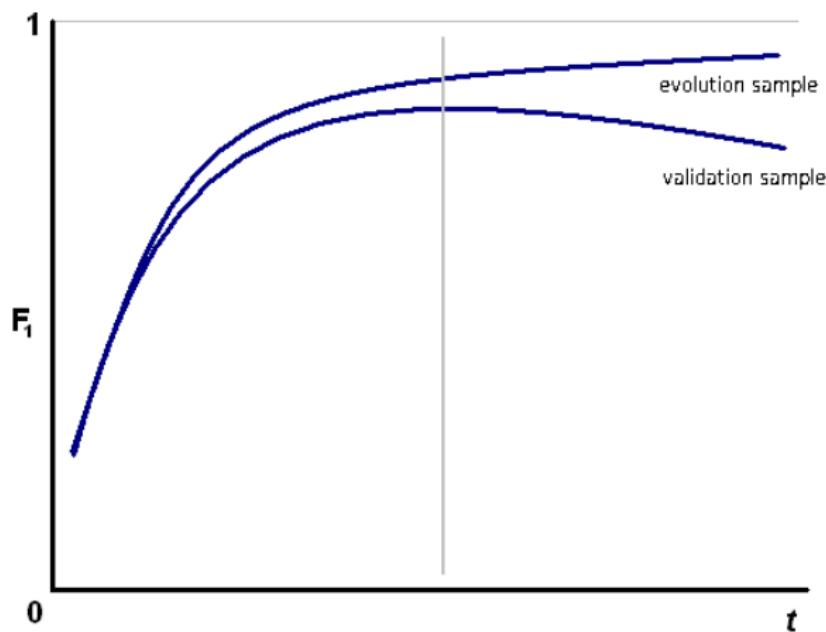
(insertion)

# AM fitness

- ① Fitness measured in terms of AM's  $F_1$  score
  - ▶ measured on a sample of **100 positive/100 negative** examples
- ② **Parsimony pressure**
  - ▶ prevents unlimited growth of the tree
  - ▶ Ocam's razor principle (less is more)

$$\text{fitness}(j) = F_1(j) + \eta \frac{L_{\max} - L(j)}{L_{\max}}$$

## Stopping criterion



- $k$  iterations without fitness improvement on **validation sample**

# Experiment

- **Goal:** evolve 3-gram AMs for legal corpora collocation extraction
- **Corpus:** 7008 legislative documents
- Hand-annotated samples
  - ▶ **evolution:** 100 positives/100 negatives
  - ▶ **validation:** 100 positives/100 negatives
  - ▶ **positives:** compound nouns, terminological expressions, proper names

# Experiment (2)

## Initial population:

- 50 – 50,000 randomly generated solutions (AM trees)
- known measures have a 50% chance of being included in the initial population

## Evolution parameters:

- three-tournament selection
- parsimony  $\eta$ : [0, 0.05]
- mutation probability: [0.0001, 0.3],
- 800 runs with randomly chosen parameters

## Results

- **20%** measures with  $F_1 > 80\%$ 
  - ▶ 23% if known AMs are included in the initial population, 13% if not
- Overall best:  $F_1 = 88.4\%$  with 205 nodes

## Measure $M_{205}$

$$\begin{aligned} & f(abc) f(a) f(c) * / f(abc) f(ab) f(c) - f(c) f(bc) f(b) \\ & -f(abc) + / + / N * f(b) + * \ln f(c) f(b) * * N f(a) * \\ & f(abc) f(a) f(abc) f(a) f(c) * / f(bc) * f(bc) f(b) + / \\ & f(a) N \text{IF}(vr(b)=\{X\}) * (-14.426000) f(b) + / N * f(bc) f(b) \\ & -(2.000000) * \ln \ln / f(a) f(c) * (2.000000) * \ln \ln / N * \\ & \ln * / f(bc) * f(bc) f(b) + / N * (-14.426000) f(b) + / N * \\ & f(abc) N f(a) * f(a) f(abc) f(a) f(c) * / f(bc) * f(abc) \\ & f(b) + / N * (-14.426000) f(b) + / N * f(b) f(c) * \ln \ln / \\ & f(abc) f(a) f(c) * / f(c) * \ln \ln (2.000000) * \ln \ln / N * \\ & / N * / N * \ln f(c) * / f(a) f(b) + * \ln \ln f(abc) f(abc) \\ & f(a) f(a) N \text{IF}(vr(b)=\{X\}) (-14.426000) f(b) + * / N * / N * \\ & \ln f(c) * / f(a) f(b) + * \ln \ln * \ln \ln / f(abc) f(a) f(c) \\ & * / f(a) f(b) + * \ln \ln (2.000000) * \ln \ln / N * \ln \ln \\ & \text{IF}(vr(c)=\{X\}) N * \text{IF}(vr(b)=\{X\}) \end{aligned}$$

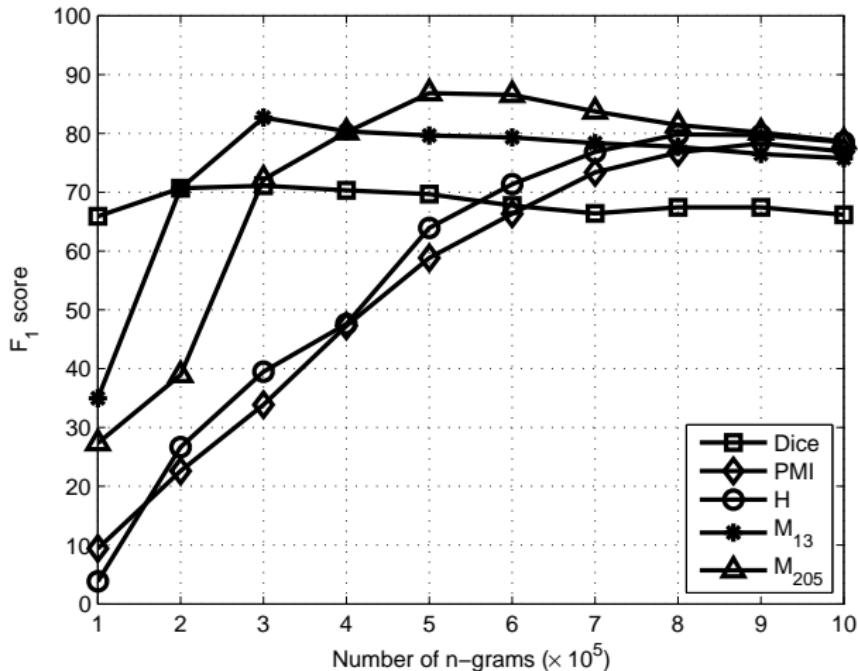
## Results

- 20% measures with  $F_1 > 80\%$ 
  - ▶ 23% if known AMs are included in the initial population, 13% if not
- Overall best:  $F_1 = 88.4\%$  with 205 nodes
- Best among solutions of less than 30 nodes:

$$M_{13}(a, b, c) = \begin{cases} -0.423 \frac{f(a)f(c)}{f^2(abc)} & \text{if } \text{stop}(b) \\ 1 - \frac{f(b)}{f(abc)} & \text{otherwise} \end{cases}$$

- ▶ evolved without measure **H** being included in the initial population!
- ▶ 96% best-ranked measures featured a similar stopword-sensitive condition
- Confirms that trigrams with stopwords should be treated differently

## Results (2)



- Measures  $M_{13}$  and  $M_{205}$  outperform the other three considered measures

# Outline

- 1 Collocations and multiword expression
- 2 Collocation extraction
- 3 Genetic programming of AMs
- 4 Conclusion

# Conclusion

- **MWEs** are important for NLP/IR/TM
- MWEs can be identified based on co-occurrence analysis – **collocations**
- **Association measures** can be used to rank and extract collocations from corpus
- **Genetic programming** can be used to evolve new AMs
- Evolved AM will perform at least as good as any other AM included in the initial population
- Approach not limited to any specific type of collocation, language, or corpus
- FW: collocations in TM/IR, syntactic features, semantic models of collocations, ...

## Selected references

- Delač, D., Krleža, Z., Dalbelo Bašić, B., Šnajder, J., Šarić, F. TermeX: A Tool for Collocation Extraction. (2009) *Lecture Notes in Computer Science (Computational Linguistics and Intelligent Text Processing)*. 5449, 149–157.
- Kolar, M., Vukmirović, I., Dalbelo Bašić, B., Šnajder, J. (2005). Computer-Aided Document Indexing System. *Journal of Computing and Information Technology*, 13(4), 299–305.
- Petrović, S., Šnajder, J., Dalbelo Bašić, B. (2010). Extending lexical association measures for collocation extraction. *Computer Speech and Language* 24(2), 383–394.
- Petrović, S., Šnajder, J., Dalbelo Bašić, B., Kolar, M. (2006). Comparison of Collocation Extraction Measures for Document Indexing. *Journal of Computing and Information Technology*, 14(4), 321–327.
- Šnajder, J., Dalbelo Bašić, B., Petrović, S., Sikirić, I. (2008). Evolving New Lexical Association Measures Using Genetic Programming. *Proceedings of ACL-08: HLT, Short Papers*, 181–184.