# Type-based Idiom Extraction

## Jan Bušta

**based on**

A Measure of Syntactic Flexibility for Automatically Identifying Multiword Expressions in Corpora
by Colin Bannard

and

Automatically Constructing a Lexicon of Verb Phrase Idiomatic Combinations
by Afsaneh Fazly and Suzanne Stevenson

25. 5. 2010, uds

# Outline

1) competition

2) object of interests

3) main goal

4) fixed ness

5) results I

6) mutual information

7) results II

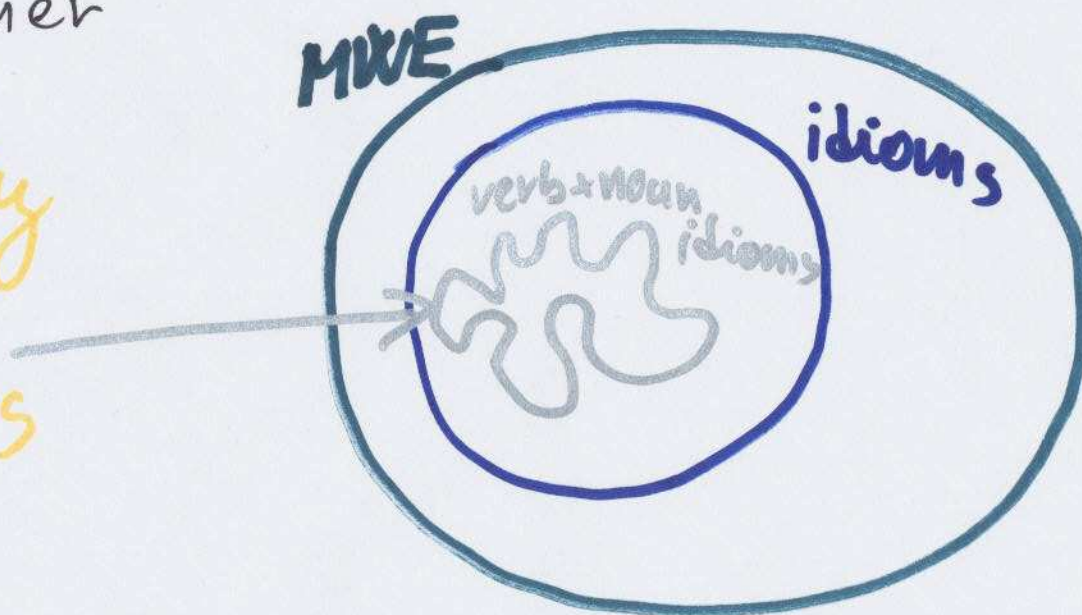8) conclusion

9) evaluation of competition

# Objects of interests

Idiomatic phrases in form of

## verb + noun

- transitive verb (needs an object)
- noun (with determiner

We will deal only
with this type
(verb+noun) in this
presentation

MIXE

idioms

verb+noun
idioms

# Main goal

## Is the phrase an idiom?

- determining fixedness
  → level of idiomacity

- computing mutual information

# Fixedness

(non) morphological variation of the phrase

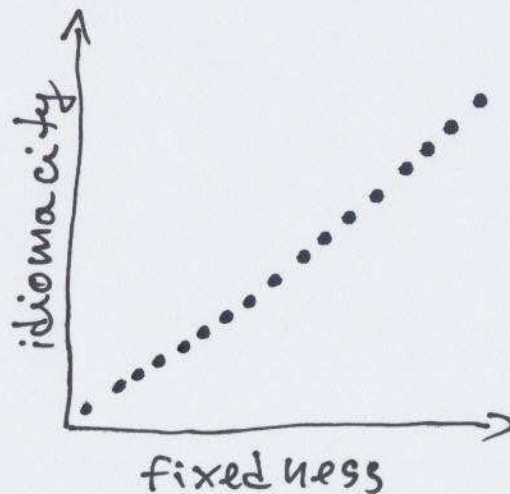- determiner

  run the show → run their show

- internal modification

  break the ice → break the diplomatic ice

- passivisation

  call the shots → the shots were called by ...

- pluralisation



idiomacity

fixedness

# Fixedness $\mathrm{I\!I}$

- lexical

$$\cdots\cdots\cdots\cdots\cdots\cdots\cdots\cdots > PMI(v, n_j) = \log \frac{P(v, n_j)}{P(v) \cdot P(n_j)}$$

$$F_{lex}(v,w) = \frac{PMI(v,w) - \overline{PMI}}{S} \cdots > \text{the mean}$$

$\cdots\cdots >$ the standard deviation

$n_j$ is a synonym for $w$

- syntactic

$\cdots\cdots\cdots\cdots\cdots$ probability of pattern $pt$ given $v$ and $w$

$$F_{syn}(v,w) = D(\overbrace{P(pt \mid v, w)} \| \underbrace{P(pt)})$$

$\uparrow$

$\cdots\cdots$ probability of pattern $pt$

$\cdots\cdots$ Kullback Lieber divergence

- hybrid

$$F_{hyb}(v,w) = \alpha F_{syn}(v,w) + (1-\alpha) F_{lex}(v,w)$$

$\cdots > $ preference of $F_{syn}$ and $F_{lex}$

# Results I

| Measure | Accuracy | Relative error reduction |
|---------|----------|--------------------------|
| Random  | 50%      | —                        |
| PMI     | 64%      | 28%                      |
| $F_{lex}$ | 65%    | 30%                      |
| $F_{sgw}$ | 70%    | 40%                      |
| $F_{nyb}$ | 74%    | 48%                      |

- canonical form

  → by set up threshold border

  → determining from patterns set

$$Cf_k(v,w) = \frac{f(v,w,pt_k) - \overline{f}}{s}$$

$pt_k \in$ set of patterns for the phrase $\underline{v}, \underline{w}$

# Mutual information

- amount of information in bits that $y$ provides about $\underline{x}$ given $\underline{z}$ (and vice versa)

$$I(x;y\,|\,z) = \boxed{\phantom{xxxx}} = \log_2 \frac{p(x\,|\,y,z)}{p(x\,|\,z)}$$

black

box

- syntactic variation

$$\text{Syn Var}\,(W) = \sum_i I(\text{VerbVar}_i\,;\,Obj\,|\,Verb) +$$

$$\sum_j I(ObjVar_j\,;\,Verb\,|\,Obj)$$

# Results II

- using mutual information (with frequency)
  is better than t-score, MI-score, ...

  they are based on frequency

- combining determiner variation, internal modification and
  passivisation goes to best results than a frequency
  based scores, but combining freq. with P,I,D is the best

# Conclusion

## Advantages

- both techniques are robust
- they work independent on dictionary

## Disadvantages

- bigger corpus => better results

- verb + noun phrases ONLY!

# Evaluation of competition

- play the second violin
- sleep on laurel
- have a chicken brain
- have a handle/have a window
- have big eyes
- finish with his own hand
- ask for hand
- draw the same rope

- spit to the water well
- pull at hair
- go to the dog-rose
- have a juice
- be behind the water
- have hare's plans
- admit of color

And the winner is...

Any questions ?

:-)

Thank you for
your attention.

xbusta@fi.muni.cz