

Word Level Analysis

Pavel Šmerk

Natural Language Processing Centre, Faculty of Informatics
Masaryk University, Brno, Czech Republic

16.5.2014

Motivation

chladnička s mrazákem

Přibližný počet výsledků: 1 940 000 (0,16 s)

[Volně stojící chladničky s mrazákem dole | MALL.CZ](#)

Nabízíme zajímavou nabídku volně stojících chladniček s **mrazákem** dole značek AEG, Amica, baumatic, Fagor, LG a další. Vše pohodlně do 24 hodin u vás.

[DATART | Ledničky](#)

Chladničky kombinované s mrazákem nahoře. Americké chladničky. Auto, přenosné ...
NO FROST **chladnička s mrazákem**. Doprava Datart. Přidat do porovnání ...

[Lednice A++ - Heureka.cz](#)

A++, kombinované, volně stojící, 230 l , 180 cm , **Mrazák**. Úsporná kombinovaná
chladnička Gorenje RK 61820 W je řazena do energetické třídy A++. Police ve ...

[Bílé zboží > Chladničky > Mrazák dole - Srovnání cen - Zalevno.cz](#)

Porovnejte si ceny zboží z kategorie Bílé zboží > **Chladničky > Mrazák dole** - najděte nejlevnější obchod.

Motivation

- Many applications need a tool for “clustering” of word forms appearing in texts:
 - *chladniček*
 - *chladničky*
 - *chladničkách* ⇔ *chladnička*
 - *chladničce*
 - ...
- Indexing, searching, keyword extraction, ...
- And almost all NLP tools

Word Level Processing Data for Czech

- For almost 12 M word forms (incl. colloquial forms):
 - lemma (canonical form, dictionary form)
 - grammatical information: part of speech, number, case etc.
- Word form *stroj* has 3 interpretations:
 - lemma *stroj*, nominative
 - lemma *stroj*, accusative
 - noun, masculine animated, singular
 - lemma *strojit*,
 - verb, 2nd person, singular, imperative mood

Possible Applications

- Various types of analyses:
 - word form ⇒ lemma (many types of searching/indexation)
 - *nebral* ⇒ *brát/nebrat* (*úplatky*)
 - *nejstaršího* ⇒ *nejstarší/starý* (*člověk*)
 - *chladnička* ⇒ *chladničky* (as a class)
 - *bavlna* ⇔ *bavlněný* (word derivation)
 - word form/lemma + gram. info. ⇒ word form
 - e.g. salutation generation: *pane Procházko*
 - word form/lemma ⇒ all word forms
 - word form ⇒ lemma + full/partial grammatical information
- The analysis is very fast
 - approx. 1 million word forms per second

Processing Unknown Words

- Some word forms in processed texts are unknown:
 - terms *polydaktylie*, neologisms *klausoviny*, typos *bizardního*, colloquial words *plat'áky*, etc.
- An ending of the word form is able to determine e.g.
 - lemma: *klausoviny* ⇒ *klausovina*
 - grammatical information: *bizardního* ⇒ genitive, etc.
 - derivational relations: *plat'áky* ⇒ *plat'ákový*
- Texts from a particular domain allows grouping of unknown word forms:
 - *polydaktylie*, *polydaktílích*, *polydaktylií*, ... ⇔ *polydaktylie*
 - ⇒ extension of data or more precise “guessing”

Resolving Ambiguities Using Context

- An extreme case *Stroj ženu holí*.
 - *Já stroj ženu holí, ty stroj ženu holí, ten stroj ženu holí.*
- Usual case is e.g. *stát*
 - noun: *Stát jsem já.*
 - verb: *Celá továrna musela hodinu stát.*
 - at the part of speech level, it is a bigger problem for English
- The context of the word determines its interpretation
 - rules and/or statistical data describe typical contexts of nouns, verbs, etc.
 - using such information one can tell that *stát* is noun/verb

Example of Contexts — Word Sketches

stát podstatné jméno

a_modifier	938517	-0.8	gen_2	274456	-0.7	post_verb	143087	-0.8
spojený	223381	12.28	hlava	20922	8.7	dotovat	433	6.3
členský	137993	11.83	zastupování	2716	8.24	mocť	15773	5.93
americký	29942	9.01	složka	5263	7.9	hodlat	528	5.87
demokratický	12202	8.46	majetek	5793	7.85	dlužit	342	5.87

stát sloveso

has_subj	942837	-3.7	post_v	184481	-1.5	is_subj_of	127156	-0.5
zázrak	4433	7.12	čelo	11624	9.36	zavázat	469	6.58
nehoda	4438	6.87	pozadí	2507	7.83	hospodařit	517	6.56
socha	3587	6.72	fronta	2654	7.72	zůstat	3245	6.5
kostel	3714	6.39	přepočet	1098	7.35	přispívat	1021	6.46



Spellchecking and Diacritics Restoration

- Data also allow spellchecking and diacritics restoration:

Result of tool CZ accent

Pred domem zastekal cerny pes.

Před domem zaštěkal černý pes.

Universality

- All the mentioned processes can be
 - tuned for a specific domain
 - using texts from this domain
 - applied to a language other than Czech
 - (Slovak, Polish, German, English, ...)

Latest Applications

- Seznam.cz, Yandex.ru, Aukro.cz, Václav Havel Library
 - indexing and searching
- Information System of Masaryk University
 - other universities and schools (FHS UK, JAMU, VŠFS, ...)
 - affiliate projects (theses.cz, odevzdej.cz, repozitar.cz)
 - indexing, searching and plagiarism detection
- “Internetová jazyková příručka”
 - online source on Czech orthography and grammar
 - NLP Centre data were a starting point for word form tables

Conclusions

- Word level processing of texts allows:
 - various types of base word determining which forms are to be grouped together
 - ambiguity resolution according to the context
 - word form generation
 - spellchecking, diacritics restoration
- The tools/data can be domain specific and for various languages