

Word Level Analysis

Pavel Šmerk

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

Workshop of Natural Language Processing Centre
May 28, 2013

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* \iff *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 \Rightarrow lemma (many types of searching/indexation)
 - *nebral* \Rightarrow *brát/nebrat* (úplatky)
 - *nejstaršího* \Rightarrow *nejstarší/starý* (člověk)
 - *chladnička* \Rightarrow *chladničky* (as a class)
 - *bavlna* \Leftrightarrow *bavlněný* (word derivation)
 - word form/lemma + gram. info. \Rightarrow word form
 - e.g. salutation generation: *pane Procházko*
 - word form/lemma \Rightarrow all word forms
 - word form \Rightarrow 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 *plaťáky*, etc.
- An ending of the word form is able to determine e.g.
 - lemma: *klausoviny* \Rightarrow *klausovina*
 - grammatical information: *bizardního* \Rightarrow genitive, etc.
 - derivational relations: *plaťáky* \Rightarrow *plaťákový*
- Texts from a particular domain allows grouping of unknown word forms:
 - *polydaktylie*, *polydaktiliích*, *polydaktylií*, ... \Leftrightarrow *polydaktylie*
 - \Rightarrow 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

<u>a modifier</u>	<u>938517</u>	<u>-0.8</u>	<u>gen 2</u>	<u>274456</u>	<u>-0.7</u>	<u>post_verb</u>	<u>143087</u>	<u>-0.8</u>
spojený	<u>223381</u>	12.28	hlava	<u>20922</u>	8.7	dotovat	<u>433</u>	6.3
členský	<u>137993</u>	11.83	zastupování	<u>2716</u>	8.24	mocť	<u>15773</u>	5.93
americký	<u>29942</u>	9.01	složka	<u>5263</u>	7.9	hodlat	<u>528</u>	5.87
demokratický	<u>12202</u>	8.46	majetek	<u>5793</u>	7.85	dlužít	<u>342</u>	5.87

stát sloveso

<u>has subj</u>	<u>942837</u>	<u>-3.7</u>	<u>post v</u>	<u>184481</u>	<u>-1.5</u>	<u>is subj_of</u>	<u>127156</u>	<u>-0.5</u>
zázrak	<u>4433</u>	7.12	čelo	<u>11624</u>	9.36	zavázat	<u>469</u>	6.58
nehoda	<u>4438</u>	6.87	pozadí	<u>2507</u>	7.83	hospodařit	<u>517</u>	6.56
socha	<u>3587</u>	6.72	fronta	<u>2654</u>	7.72	zůstat	<u>3245</u>	6.5
kostel	<u>3714</u>	6.39	přepočť	<u>1098</u>	7.35	přispívat	<u>1021</u>	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.

- 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, ...)

- 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

- 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