The Lexico-Semantic Annotation of The Prague **Dependency** Treebank

Eduard Bejček, Petra Möllerová, Pavel Straňák

Department of Formal and Applied Linguistics Charles University, Prague



Outline

- Prague Dependency Treebank
- Annotation
 - Corrections
- Current work and further plans



Prague Dependency Treebank (PDT)

• 3 layers of annotation:

• w-layer: segmentation and tokenization

- m-layer: lemmas and morphological tags
- a-layer: analytical (surface) dependency trees
- t-layer: tectogrammatical (deep) dependency trees



Layers of Annotation







Tectogrammatical tree

Department of Formal and Applied Linguistics Charles University, Prague



Annotation

• Goal:

6 /19

- Manual identification of word senses to obtain a training data for automatic WSD
- Requirements:
 - A semantic lexicon set of all possible meanings (tags) for each word.
 - A method/procedure that assigns a semantic tag to each occurrence of a word.

Department of Formal and Applied Linguistics Charles University, Prague



word = instance of a lemma from the lexicon

Czech WordNet

- Developed at The Masaryk University, Brno
- Originally in EuroWordNet 2, continuing development within the Balkanet project
- Mapped directly to the Princeton WordNet 2.0
- almost XML format
- 17,000 nouns; 2,000 verbs; 4,000 adjectives and adverbs



Annotation

Image: Second state in the second s	 >> představenstvo-wsd-n-1-11799-35906-5319315/n/o ************************************
odmínky ůvodů Představenstvo Jana Gondu najetku Jrodal Společnosti Souladu	

Two annotators in parallel, no interaction (coordination) between them. Instance of each word (lemma) in the CWN (annotation lexicon) annotated. Other words skipped.

Results of Annotation

All words	431447	100.0%	
Autosemantic words	300725	69.7%	100.0%
Annotated words	148744	34.5%	49.5%
Ambiguous words	101703	23.6%	33.8%

- Inter-annotator agreement: 61.7%
- 5,111 sentences fully disambiguated

9 /19

Department of Formal and Applied Linguistics Charles University, Prague



agreement on all words to be annotated

Corrections

- Why?
 - Agreement only 61.7%. We wanted to raise number of correct sentences with small effort.
- How?
 - Third annotator "corrector".
 - Only first 25 lemmas.
 - Why this way?

10/19

Department of Formal and Applied Linguistics Charles University, Prague



the agreement wasn't too good corrector -- inspect cases with non-agreement



3700 of 4700 lemmas have freq. =< 10;

11/19

25 lemmas with freq. > = 200 ... greatest gain for our work:

0.5% of problem lemmas => 7.4\% more annotated words

it's easier to correct let say 100 of occurrences of one lemma then 100 lemmas with only one occurrence

Preparing Data

- Annotations include many error categories to mark the kind of problem. Use them for:
 - splitting of error annotations for each lemma according to the type of disagreement (see the article)
 - faster and more convenient correcting

12/19

Department of Formal and Applied Linguistics Charles University, Prague



many small files we inspected the cases, why annotators doesn't make the same decision we put the same cases to one file for one lemma

Technical Details

- Vim + macros + highlighting
- one key-punch per instance
 - A was right
 - B was right
 - none of them was right, editing of CWN, new sense assigned

13/19

Department of Formal and Applied Linguistics Charles University, Prague



we use only vim editor with macros and highlighting because corrector had only 3 kinds of choice -- and then another instance

- "annotator A was right"
- "annotator B was right"
- "synsets are quite confusing, I'll repair them and assign the sense from upgraded lexicon"

Results of Corrections

- corrected 25 (0.5%) of 4,738 lemmas
- gained 10,971 annotated words (+7.4%)
- Fully annotated sentences:
 - 5,111 => 6,941 (+35.8%) SemEval (Senseval 4)
- improved some difficult lemmas in CWN
- Cost: approx.
 - 320 hours (corrector) + 215 hours (programer)

Department of Formal and Applied Linguistics Charles University, Prague



move from ... to ... these sentences are submitted as one quest in SemEval CWN -- feedback for authors

Near Future

- PDT 2.0 data; format: PML (XML based)
- PML:
 - stand-off annotation (4 layers)
- addition of s-layer ("sense")
 - not a deeper layer, can ref. to a-nodes or t-nodes
 - a list of pairs: lexicon.ref => (t / a)-node.ref.list

15/19

Department of Formal and Applied Linguistics Charles University, Prague



Data (texts) are differently split into files. The files used to be roughly similar in size in PDT 1.0.

A Little Beyond

- the original project resulted in a new one:
- new goal: improvement of *t-lemmas*

(i.e. lemmas of nodes of trees on t-layer)

Department of Formal and Applied Linguistics Charles University, Prague



• Can word sense disambiguation help statistical machine translation?



Methodology

- 2 rounds:
 - multi-word lexemes and named entities
 - there is no point in assigning single-word senses to these
 - remaining single-word lexemes
- New annotation tool for 1st round



File			<u>H</u> elp			
ilename: /home/bejcek/svn/trunk/data/cmpr9410_011.t.gz	Load	Save	Save & Exit			
dva samostatné rámečky Míra nezaměstnanosti by se měla vyvíjet protikladně, než ve s specifických podmínkách české ekonomiky, mj. vzhledem k netržnímu cho podniků, nízkým mzdám jakož i rychlému rozvoji drobné podnikatelské a a šedé ekonomiky), růst nezaměstnanosti v letech 1991 · 1993 značně z Pokračující privatizace a restrukturalizace si však vynutí zvýšení mi koncem roku 1993 na 5 · 6 % ke konci přištího roku. * Saldo běžného účtu platební bilance podle odhadu dosáhlo vlon	tandardn ování nep ktivity aostal z ry nezam	í ekonom rivatizo (včetně a pokles éstnanos 0 mil. U	nice. Ve ovaných tzv. černé sem HDP. sti z 3.5 %			
téměř 2 % HDP. I když letos a příští rok je nutné počítat se zpomalen zrychlením růstu dovozu, prognózujeme, že saldo přesto zůstane kladné ročně (l • 1.6 % HDP). To by umožnilo dále zvýšit devizové rezervy ČN konvertibility koruny. popisky pro ilustrace (na výběr barva a čb): BARVA:	im růstu ve výši B, potře	vývozu <mark>300 – (</mark> bné k za	a 00 mil. USD avedení plné			
NEJMÉNĚ DVOUPROCENTNÍ RŮST ČESKÉ EKONOMIKY JIŽ LETOS:						
Je to radostnější objev, než vidět slunce v temném lese						
			þ			
* "devizové rezervy" [4. odstavec]> wsd#object.						
Značky						
Obecné Pojmenované entity						
Ukázat Odstranit Jméno Instituce Místo Objekt Adresa Č	Čas Bi	blio F	oreign ?			
/19 Departmer	nt of Forma C	l and Appl Charles Uni	ied Linguistics versity, Prague			