

PA153 - NLP

1. Formální gramatiky

- soubor pravidel, jež generují či rozpoznávají věty
gramatiky automaty

- terminální a neterminální soubory

- $S \rightarrow NP VP$, $NP \rightarrow A N$
věta jm. přísl. adj. subst.
frase fr.

- množina: terminály (abeceda), neterminály (synt. kat.), podm. kart. součinu $(N \cup T)^* N$,
počítačové symboly gramatiky

Jazykové inženýrství

- všeobecně NLP - na FI tvorba slovníků (DEBdict,...), korpusy,...

2. Morfologická rovina popisu jazyka

- nejvyšší rovina

- cílem dostat: 1/ lemma, 2/ gramatické významy

- problémy:

- jak popisovat (prašský vs. bruňský popis)

- do jaké míry lemmatizovat - *otcova* \rightarrow *otcův/otec*?

- co s dubletami (*myslet/myslit*)

- manuál pro ruční anotaci musí být jednoduchý a jasný
(ne jako prašský 300str. man.)

- když se na tom neshodnou lidi, stroj to může udělat líp

Anotační systémy

- poziční (prg): význam zn. určen posicí: *UNIS4-----A----* - Hajič, Hlaváčková, Hladký

- atributový (brno): dvojice atribut-hodnota: *k1gInSc4* - Sed líček, Emark
druh rod čís. před

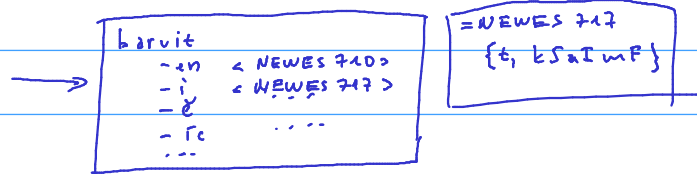
- heterogenní (brno): komb. atr. a pos, posice kódují atr. i hodnotu

- rychle vyčerp. značky

- BNC používá pevnou množinu hotových značek pro ang (Adj - adjective, comparative)

AJKA

- organizuje data do lemmatizovaného vzoru
- pro es asi 1800 vzorů
- odtrhávání a připojování intersegmentů a koncovek → generování tvarů
- db analyzátorů: množ. koncovek, vzory + interseg., přípustné kombinace



AJKA vs. MAJKA <https://nlp.fi.muni.cz/czech-morphology-analyser/majka.html>

- obecně: úkoly morf. analyzátorů:

- vytvořit všechny verze (tvary)
- vybrat nejpravděpodobnější (tagger)
- zanalyzovat neznámá slova (guesser)

- A: - už nevyvíjený, blbě napsaný, velká, nepřehledná

- chyby v práci se složeninami

- "ajka dělá nějakou magii se skládáním slov, blíže jsem nezkoumal"

- špatný handling záponých tvarů "být"

- podivnosti s vidy

2.1. Reprezentace významu

- sémantika nezávislá na ^{užití, situaci} kontextu, ale pragmatika může význam změnit
- repr. predikátovou logikou (omezená vyjadřov. síla), TIL (složitá)
- nutno zpracovat deixi, anafory, katafory, ...
- koupil jsem si auto. [To auto] je bílé.
- pro analýzu nutno mít KB (DB znalosti) - ontologie: klasifikace/vztahy objektů universa (eg. WordNet)
- TIL pracuje i s množinami sufixů, časem, ...
- Levin - klasifikace sloves (pohyb, emoce, změny, ...)
- víceznačnost - problém výběru správného významu - Leskiv alg. (společná slova v kontextu zkoumaného a v det. jednotlivých významů zkoum.)

3. Morf. analýza

- viz #2

3.1 Sémantické role a pády

- role: agens, patiens
- valenční rámce (třeba u sloves v lexikonu); nutno ověřovat na datech v korpusu
- Fl: Verbalex (Pala)
- pády: oblique cases - nesymetrické pády

4. Algoritm. popis deklinace, konj., lemma tisace

- uzony, kmeny - viz #2

4.1. Sloves. val. rámce a třídy

- rámce: E. Rosch: teorie prototypů - obecný zástupce třídy + speciifikátory

- "židle je víc nábytek než sporák" tedy nábytek, ale víc se liší od klasického nábytku

5. Lexikální analýza

viz #2.1.

5.1. Syntaktická analýza, koncepce

- gramatiky (viz #1)

6. Elektronické slovníky

- Brown Corpus, Longman, COBUILD, TEI (guide, jak dělat slovníky - strukturu)

- v XML, DTD popisuje strukturu; XSLT - transformace do jiných formátů

- dict writing systems - app na tvorbu slovníků (propojení s korpusy, morf. analyzátory, ...)

- komerční vs. open - DEB (dict edit & browse)

- DEBdict, DEBvisDict, ISP

- praktická vs. teoretická lexicografie

↳ analýza slov. zřs., teorie tvorby slov.

- lexikální databáze - detailní DB s valencemi, gramatikou, styl. zařazením, ...

- PRALED, DANTE (eng)

↳ Pražská Lexikální Databáze

↳ Database of ANalysed Texts of English

- macrostructure - strukt. většího objektu (asi celého významu vč. příkladů, valencí, etc.)

- microstructure - strukt. jednoho malého významu

6.1. Analýza promluvy

- odkazy v textu: anafora (vzad), katafora (vpřed), endofora (uvnitř textu), exofora (vně textu),
koreference (různé výrazy označující totéž (Klaus, bývalý prezident))
- anaphora resolution - rozpoznání a. - k čemu se odkazuje se impl. pomocí zásobníku (naivní přístup)
- umění mít znalosti o světě - KB (knowledge base)
- SUMO-MILO, WordNet, ConceptNet

7. Víceznanost v NL, desamb., evaluace výsledků desamb.

- WSD - Word Sense Disambiguation - na základě kontextu
- $f(w, c) \rightarrow s$, kde w = množina slov, c = množ. kontextů, s = množ. kontextů
- Leskiv algoritmus - porovnání výskyt slov v kontextu znám. slova se slovy v definicích významů znám. slov, vybere nejvíce shodu (naivní)
- všechny WSD závisí na nějaké DB významů; ty se liší (SSJČ x SSČ x ...)
- komponentová analýza (+HUMAN, +FEMALE) → sémantické třídy (příbuzná slova) - síť. síť → odvozování
- WordNet - významy uloženy podobně jako v mozku (asi?) = ontologie (konceptualizace)
- hyperonymie, hyponymie, tropov, meron, ...
- rozliš. lex. významu (pro stroje): vzdálenost mezi koncepty

7.1. Tagsety

- prg x brno - viz #2.

8. Notace pro synt. analýzu, partiální analyzátor

- viz #1
- partiální synt. anal. - pracuje s úplnými texty, s lexikálním významem
- užito pro rozpoznání klíčových frází - ToPicks

8.1. Logická analýza významu věty

- TIL - množina možných světů, pravdy závisí na světě a na časě
- intense ?
- extense .

9. Desambiguace pro NL

- morfologická (tagger), lexikální (WSD)
- M. současnosti analyzátorů (zřejmě?) - Mojka, TreeTagger, ^{o2} ^{entiten} ^{BNC} CLAWS
- vybírá nejpravděpodobnější kategorii podle kontextu, příp. dle frekvencí výskytů kandidátů v trénovacích datech
- když neví, guesser - dost random výběr
- WSD (viz #7)

9.1. Princip kompozicionality (log./sém. anal., TIL)

- p. kompozicionality = význam skup. slov je složeninou významů jednotlivých slov
- vepnutí u pevných kolokací = MWE (multiword expressions) - nutno detekovat např. pro MT
- MWE by měly mít vlastní slov. hesla (asi?)
- pevné (zakopaný pes) vs. vzory (vzít <několik> na hůl)
- log. anal., TIL, PL1 - ???

10. Sémantické značkování, desambiguace

- WSD viz #7
- sém. značkování: sém. síť jako WordNet užívají XML a rekursivní reference u jednotl. synsetů ^(synonym. řada) ke svým hypo-, hyper-,mero-,...-nýmům

10.1. Pragmatická analýza

- p. jevy jsou mimojazykové
- externí prag.: komunikační situace (mluvčí, posluchač, čas, ...)
- interní prag.: vztah mezi mluvčím a významem
 - oznámení, rozkaz, ...
 - implikatury - intence mluvčího
 - ilokuce (co bylo řečeno), ilokuce (co bylo interpretováno), perlokuce (důsledek, činy)
 - Griceovy maximy (kvality, kvantity, ...)

11. Lexikální DB a nástroje

- viz #6

11.1. Interní a externí pragmatika

- viz # 10.1

12. DC gramatiky a využití Prologu v synt. anal. NL

- DCG = gramatiky vynezaných klauzulí

- umožňují zachovat nekontextovou podobu pravidel a získat kontextovou citlivost

- pravidlo: neterm. \rightarrow predikát s argumenty ($S \rightarrow Np Vp \sim s(s(Np, Vp))$)

- rekurse!

- Prolog: zápis predikátů, inference - logické uvozování faktů

12.1. Form. reprezentace významu, ^{FOL}PL1, TIL

- PL1 - omezení, TIL složitý...

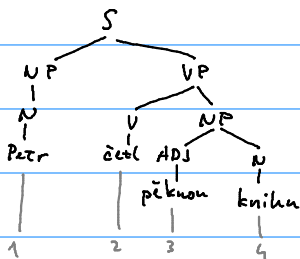
- ne všechno je predikát

...

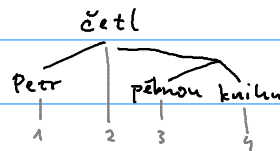
15. Reprezentace pro syntakt. analýzu

- složkově - zobrazení složek věty

(a jejich závislosti)



- závislostní - zobrazení jen synt. závislosti slov



! největší sranda je u stroj. analyzátorů poznat, který z výstupních stromů je správný, bo oni to nepoznají