

Počítačové zpracování přirozeného jazyka – PA153 (Natural Language Processing)

K. Pala et al
Centrum ZPJ FI MU

Podzim 2017

ZPJ (NLP) – přehled, motivace

Proč si PJ zasluhuje naši pozornost?

- jazykové chování představuje jeden z fundamentálních aspektů lidského chování
- PJ je podstatnou složkou našeho života jako hlavní nástroj komunikace,
- pomocí PJ vyjadřujeme a zachycujeme své znalosti, vědecké poznatky, vidění světa
- PJ je východiskem pro umělé (formální) jazyky
- jazykové texty slouží jako paměť lidstva pro předávání znalostí z generace na generaci
- vztahy k technice a počítačům, komunikace h-c

Terminologická poznámka

- Používané termíny
- Kvantitativní a statistická lingvistika
- Algebraická lingvistika (Chomsky)
- Matematická lingvistika (shrnující)
- Počítačová (computational, komputační) lingvistika
- Zpracování přirozeného jazyka (NLP)
- Počítačové zpracování mluvené řeči (ASR)
- Kognitivní věda (lingv., psychol., filos. i logika)

Co je předmětem ZPJ?

PJ – studuje se a zkoumá interdisciplinárně:

- V lingvistice (tradiční, strukturní, matematická)
- V psychologii a psycholingvistice
- Ve filozofii a logice – vztahy k univerzu promluvy, usuzování (inference), pracuje se s výroky
- V algebraické (komputační) lingvistice (60. léta min. stol. je klíčová role N. Chomského)
- Teorie jazyka ve formě algoritmů, datové struktury, empirická data, korpusy
- Vztahy ke kognitivní vědě a umělé inteligenci

ZPJ – vztah k počítačům

- V jazykovém inženýrství
- Potřeba dvoucestné komunikace mezi čl. a poč.
- Zatím je komunikace čl.-poč. jednocestná
- Potřeba komunikačně bohatšího rozhraní
- Rozhraní v PJ musí být chytřejší a pružnější – zejména pro nespecialisty – běžné uživatele
- Výrazné komerční důsledky pro počít. trh (OS)
- Je možný OS s PJ? – pokusy s OS Merlin
- Naše znalosti o struktuře PJ jsou neúplné
- Vztah teorie (výzkumu) a aplikací

ZPJ – aplikace 1

- Zpracování textů – korektory překlepů, gramatické, stylistické korektory
- Dělicí, fulltextové programy (lemmatizátory)
- Morfologické a syntaktické analyzátoři: Majka, synt, SET, NTA (sémantika)
- Prohlížeče, editory – webové, slovníkové nástr.
- Strojově čitelné slovníky (MRD), platforma DEB
- Dialogové a dotazovací systémy
- Turingův test (Eliza, Loebner Prize, listop. 2015)
- Extrakce informací, sumarizace, abstrakty
- Strojový překlad – viz dále

ZPJ – aplikace 2 (SP)

- Strojový překlad – max. snaha o využití v praxi
- EU projekty – EuroMatrix, EUM+, Present aj.
- Google Translator – je dostatečně použitelný?
- Systran – dříve oficiální systém SP v rámci EU
- Systémy s překladovou pamětí – Trados (lokalizační systémy), paralelní korpusy
- Systémy pracující s podjazyky (Taum Meteo)
- Hlasový SP – příklad: systém Verbmobil (1992-2001)
- Budoucnost SP? Firmy: Google, IBM, neur.sítě

ZPJ – aplikace 3 (mluvená řeč)

- Hlasové ovládání počítačů
- Syntéza – systémy TTS, Demosthenes (demo)
- Automatické rozpoznávání řeči (ASR), diktovací stroje (demo)
- Via Voice (IBM), Dragon (Nuance), ang., něm., it.
- Pro češtinu – systém Dictate 4.5..., Newton Technologies (demo)
- Aplikace na soudech, v parlamentu, v medicíně
- Úroveň porozumění u těchto systémů – cca 90 %
- Můžeme si se svým notebookem popovídat?

ZPJ – další aplikace 4 (vztah k AI)

- Expertní systémy – např. Mycin (lék. diagnostika)
- Databázové systémy s PJ rozhraním
- Porozumění příběhům a porozumění PJ
- Abstrakty z novinových článků – konference MUC (Message Understanding Conference)
- Robotické aplikace – soutěže robotů, SHRDLU, 1971 (T. Winograd), první systém obs. Znalosti, inferenci, gram.
- Sémantický web – vyhledávání, uplatnění metadat
- Chytřejší vyhledávání – Google stačí? Seznam? IBM Watson
- Ontologie a konceptuální systémy pro jednotlivé domény, sémantické sítě (WordNet)

Historie ZPJ v ČSR a ČR 1

- Praha – FF UK, seminář SP, 1958
- B. Palek, vztah k N. D. Andrejevovi.
- P. Sgall, P. Novák, D. Konečná, L. Nebeský, E. Hajičová, J. Panevová, P. Piřha, K. Pala
- M. Těšitelová – odd. matemat. lingvistiky, ÚJČ, Frekvenční slovník češtiny, 1961, 1983
- L. Doležel, vedoucí odd. matem. lingvistiky, ÚJČ, konflikty Letenská vs. Malostranské nám.
- J. Štindlová – počátek počítačového zprac. PJ na děrných štítcích

Struktura (roviny) jazyka

- Povaha jazykového systému – jazykové roviny a jejich formální popis – typy teorie
- Fonetika a fonologie
- Morfologie – flexe (ohýbání) a derivace
- Syntax (skladba)
- Sémantika – lexikální, logická
- Pragmatika – vztahy k uživatelům
- Promluva, anaforické vztahy
- Na všech rovinách se budují algoritmické popisy a k nim vhodné počítačové aplikace

Paradigmata v NLP

- Introspektivní – Chomsky, kompetence : performance, generativní a transformační gramatiky
- Gramatiky jsou chápány jako množiny pravidel – jejich neúplnost je klíčová
- Empirická data – počátek korpusů: Brown Corpus, H. Kučera, N. Francis (1960-61), 1M
- Velké počítačové soubory jazykových dat
- Pravidlové vs. statistické přístupy, výhody vs. nevýhody
- Strojové učení, jazykové modely – kdo vede?

Roviny – fonetika, fonologie

- Zvuková stránka jazyka – hlásky (fóny)
- Fyzikální vlastnosti mluvené řeči (signálu)
- Fonologie – fonémy – abstrakce nad hláskami
- Nejmenší jednotky rozlišující význam, *pas – pás*
- Fonologické protiklady: délka – krátkost: *vola/á*
- Souvisí se zpracováním mluvené řeči
- TTS – syntéza řeči, Demosthenes
- ASR (ARŘ, demo)
- Intenzivní výzkum, IBM, Nuance, hodně peněz

Morfologie

- Jednotky – morfémy, nejmenší jednotky nesoucí význam (obvykle menší než slova, uč-)
- Typy morfémů – nesoucí lexikální význam, kořeny či kmeny, morfémy nesoucí gramatické významy
- Slova a jejich segmentace – morfologické analyzátory – algoritmy - *ne/u/věř/i/t/eln/ému*
Flexe (tvarosloví) vs. derivační morfologie
- Čeština je jazyk s bohatou morfologií
- Analyzátory Ajka, Majka, další (Morče) pro češ.

Syntax

- Vztahy mezi slovy ve větě
- Jednotky – větné složky, větné členy, věty
- Algoritmický popis větné struktury – reprezent.
- Formální gramatiky – výsledky N. Chomského
- Hierarchie gramatik, jazyků a automatů
- Koncepce syntaxe – závislostní a složková
- Syntaktická analýza a analyzátory
- Pro češtinu nástroje – Synt, Set, (Va)Dis
- Statistické nástroje (MALT, Collins), n-gramy

Sémantika

- Nemá vlastní jednotky jako takové
- Klíčová otázka – co je to význam?
- Můžeme rozlišovat význam slov a slovních spojení – lexikální význam – lexikální sémantika
- Význam vět – větná či logická sémantika
- Sémantické reprezentace vět
- Používané formalismy – PK1, TIL aj.
- Kombinované techniky – valenční rámce
- Význam – místnost bez oken – nevidíme ven ani dovnitř (podobnost s Platonovými stíny)

Lexikální sémantika

- Významy slov a slovních spojení
- Lexikologie – nauka o slovní zásobě
- Lexikografie – zpracování slovní zásoby (elektr.) - obvykle v podobě slovníků
- Počítačová lexikografie, typy slovníků
- Softwarové nástroje pro práci se slovníky
- Popis významu slov ve slovnících – definice, synonyma,
- DebDict (<https://deb.fi.muni.cz:8005/debdict/>), přístup, platforma DEB

Pragmatika

- Vztahy mezi UJ a jazykovými výrazy
- Interní – postoje k propozici: oznamovací, tázací, rozkazovací, přací (typy vět)
- Externí – komunikační situace a její prvky, vztahy k propozici
- Pragmatická funkce – (*Já mám žízeň.*)
- Deixe a deiktické prvky
- Jejich role v komunikační situaci

Analýza promluvy

- Struktura promluvy
- Anaforické vztahy a jejich rozpoznávání
- Rozpoznávání částí promluvy
- Krabicový model
- Vztah k dialogům

Reprezentace znalostí a inference

- Sémantické sítě (WordNet, ontologie)
- Logické formalismy – PK1, TIL
- Valenční rámce
- Typy inference
- Dedukce, monotonní - nemotonní
- Systémy využívající Common Sense
- Komunikační agenti, model BDI

Strojové učení a NLP

- V současnosti populární techniky
- Přehled – samostatný výklad

Historie ZPJ v ČSR a ČR 2

- Seminář SP na FF UK od r. 1958
- Brno – počátek ZPJ v 1964 (K. Pala)
- Ústav českého jazyka FF UJEP
- V 70. letech experimenty s českými generativ. gramatikami – analýza a syntéza (OVC VUT)
- Implementace syntaktické a sémantické analýzy na počítači Tesla 200 (Čihánek, Palová)
- Havel, Machová, Pala, Sofsem 1978
- V 80. letech spolupráce s ÚVT UJEP, vytvoření českých gramatik v Prologu (počítač PDP 11)

Historie ZPJ v Brně I

- ÚVT – Benešovský, Šmídek, Gerbrich, programovací jazyk Wander (1988-90)
- 1988-9 první PC na FF UJEP MU), vznik morfologického analyzátoru pro češtinu, Xantipa
- Franc, Pala, Osolsobě, gramatický korektor, generátor a analyzátor českých vět v Prologu
- Od r. 1995 dochází k přesunu výzkumu na FI MU
- .V r. 1997 vzniká na FI MU Laboratoř ZPJ
- Umožnily to grantové projekty

ZPJ na FI MU II

- Budování korpusových nástrojů (Rychlý, 1997-8), korpusový manažer Bonito/Manatee
- Vznik české lexikální databáze WordNet, 1999
- Vytvoření nezávislého morfologického analyzátoru Ajka (Sedláček, 1999)
- Pokročilá syntaktická a sémantická analýza češtiny – systém Synt (Horák), Set (Kovář), Dis (Mráková)
- Budování slovesné databáze komplexních valenčních rámců – VerbaLex (Hlaváčková)
- Nový mf. analyzátor Majka, systém Deriv (Šmerk) a Derivancze

ZPJ na FI MU III

- Nové korpusové nástroje – slovní profily (Word Sketches) (Rychlý, Kilgarriff), LCL – ukázat
- Budování velkých webových korpusů
- Soubor nástrojů:
 - Justext – odstraňování smetí z webových stránek (boilerplate)
 - Onion – čištění duplicit z webu
 - Chared – rozpoznávání jazyků na webu
 - WSE, NoSketch, Skell (Suchomel, Jakubíček)