

# PA153 Počítačové zpracování přirozeného jazyka

## 13 – Analýza promluvy, rozpoznávání anafor

Vašek Němčík

NLP Centrum, FI MU, Brno

23. prosince 2013

# Přehled

- Text/diskurs
- Anafora – motivace, definice, úvod
- Typy anafor, potřebné znalosti
- AR algoritmy

- text/**diskurs** – jednotka jazykové komunikace větší než:
- věta/**výpověď** – minimální obsahově úplná jednotka

## **věta**

langue

competence

*produkt*

*struktura*

*nedůležité kdy/kde/jak*

## **výpověď**

parole (de Saussure)

performance (Chomsky)

*proces*

*chování*

*podmínky/okolnosti/způsob*

- referenční výrazy
- reference (odkazování)  
jazykový výraz  $\mapsto$  mimojazyková entita

# Reference

- **exofora** (vnější reference)

výraz odkazuje k entitě ve světě přímo

“Slunce”, “Alpy”, “Václav Havel”, “ten přechod před FI”

- ▶ deixe – odkazování k entitám v rámci komunikační situace (gesta, “tady”, “ted”, “tamto”, ...)

- **endofora** (vnitřní reference)

entita je určena na základě vztahu k jinému výrazu v diskursu (nejen mimojazykový, ale i jazykový kontext ...)

- ▶ anafora – výraz se vztahuje k výrazu dříve v textu
- ▶ katafora – výraz se vztahuje k výrazu dále v textu

méně častá; vysktuje se v beletrii (zvyšuje napětí):

*“Ranní světlo ho probudilo už v pět. Rychle se oblékl a nasnídal. Detektiv Jones věděl, že nemůže ztrácet čas.”*

# Anafora

- **anafora** (anaphor) – anaforický výraz (× Chomsky)
  - ▶ zejména zájmena, ale i “ten muž”, ...
- **antecedent** – předcházející výraz, ke kterému se anafora vztahuje
- **anafora** (anaphora) – anaforická reference (jev)
- **anaphora resolution** – určování anaforických vztahů (hledání vztahů mezi anaforami a antecedenty)

Příklady:

- **[Petr]<sub>i</sub>**; snědl **[koláč]<sub>j</sub>**.  
**[(on)]<sub>i</sub>**; Byl hladový a **[ten koláč]<sub>j</sub>**; vypadal lahodně.
- **[Venus]<sub>i</sub>**; rose at 0930, but I didn't see **[the thing]<sub>i</sub>**.
- **[Jones]<sub>i</sub>**; offered **[[his]<sub>i</sub>; furniture]<sub>j</sub>** for sale,  
but nobody wanted **[the stuff]<sub>j</sub>**.

## Lze udělat úkrok stranou?

- Můžeme se tomu všemu vyhnout, třeba používáním jen přímé reference?
- Nemuseli bychom se zabývat kontextem ...

NE. Z mnoha vážných důvodů:

- Lidé jsou líní.
    - ▶ anafory jsou krátké a snadno se používají
    - ▶ patrně vlastní lidské komunikaci (ve všech jazycích!)
  - diskurs není libovolná sekvence výpovědí
    - ▶ koherence – sémantická návaznost
    - ▶ kohese – gramatické a lexikální vztahy
- ↪ anaforické vztahy drží text pohromadě  
(umožňují nám se držet zamýšleného toku myšlenek)

# Ilustrační příklad

[Jarda]; si koupil Porsche. (On); Rád jezdí rychle.

[Jarda]; si koupil Porsche. [Jarda]\*i,j rád jezdí rychle.

↪ delší/složitější věta zní divně (nutí k zamyšlení)

- **Kooperační princip** (Grice)

Komunikační maximy:

- ▶ kvality
  - ▶ relevance
  - ▶ kvantity
  - ▶ způsobu
- Posлуhač předpokládá, že se jimi mluvčí řídí.
  - Když ne, má to hlubší důvody.
  - více o pragmatice v “IA091 Sémantika a komunikace”

# Proč to učit počítače?

- zásadní úzké hrdlo mnoha NLP aplikací

- **Information Extraction**

- ▶ **[Václav Havel]** was a Czech writer and dramatist.  
**[He]** was the ninth and last President of Czechoslovakia and the first President of the Czech Republic. (*Wikipedia*)

- ▶ “the best doctor in Europe” → Google

Letters from Asia addressed loosely to The Best Doctor in Europe arrived on **[his]** doorstep.

**[His]** own reputation as the best doctor in Europe couldn't save **[him]** from the tragedies of **[his]** life.

- Bez AR nenajdeme to, co hledáme.  
Pouze anaforické výrazy (které jsou samy o sobě prázdné).



# Proč to učit počítače?

- **Strojový překlad**

- CZ  $\mapsto$  EN

**[Sestřička]** mu dala **[pilulku]**. Spolkl **[ji]** a usnul.

**[The nurse]** gave him a pill. He swallowed **[her]** and fell asleep.

- DE  $\mapsto$  EN

Ich suche **[meine Uhr]**. Ich kann **[sie]** nirgendwo finden.

I am looking for **[my watch]**. I can't find **[her]** anywhere.

- nelze překládat přímo (různé gramatické kategorie)
- navíc: různé vlastnosti anafor

# Definice úlohy

- nalézt anaforické výrazy v textu
- určit k nim antecedenty
- určit typ vztahu
  - ▶ koreference  
(dva výrazy se odkazují ke stejnému promluvoému objektu)
  - ▶ bridging (asociativní/nepřímá anafora)  
(jakákoliv sémantická relace)
    - ★ hyperonymie/hyponymie  
“Nábytek je drahý. Židle jsou nejdražší.”
    - ★ část/celek  
“Každý majitel bytu se snaží zabezpečit vchodové dveře.”
    - ★ entita/vlastnost  
“Pepa má nové auto. Barvu určitě vybírala jeho žena.”
    - ★ příčina/následek  
“Včera tu byl požár. Kouř je tu stále cítit.”

# Typy anafor

- **textová vs. gramatická**

[Ben] takes a photo of [himself] every day.

- **pronominální** (pro NLP asi nejrelevantnější)

- **nominální**

Od září bude do [Brna] létat nová letecká linka. Očekává se, že přinese [druhému největšímu městu ČR] nové turisty.

- **slovesná**

John likes cats. So does Bill.

- **one-anaphora**

John has a black Porsche. I would like one too.

- **nulová (zero) anafora**

anafora není povrchově realizována

v češtině (a ostatních pro-drop jazycích) nevyjádřené podmínky

# Typy pronominálních anafor

- osobní zájmena
  - ▶ silná: “jemu”, “on”, “ona”
  - ▶ slabá: “mu”, “ho” (klitika)
  - ▶ nulová: ∅
- demonstrativní zájmena: “ten”, “ta”, “tomu”
- reflexivní zájmena: “se”, “sebe”, “svůj”
- posesivní zájmena: “jeho”, “jejího”
- relativní zájmena: “který”, “jenž”

ALE jsou i neanaforická zájmena:

- deixe: “to”
- expletivní/pleonastická zájmena:  
It's raining. / Es regnet.  
It is the first chapter, I enjoy the most.  
Zdá se, že tu někdo byl.

# Znalosti potřebné pro AR

## • morfologie

- ▶ shoda v  $\Phi$ -atributech (závislé na jazyce)
- ▶ čeština: osoba, číslo, rod
- ▶ angličina: pouze sémantický rod  
⇒ nutnost mít informaci jméno  $\mapsto$  rod

## • syntax

- ▶ posice anafory/antecedentu v syntaktické struktuře věty
- ▶ paralelismus  
tendence k zachování stejných syntaktických rolí:  
[Mary] met [Lucy] at the bus station.  
[She] asked [her] about the new neighbour.

## • pragmatika

- ▶ Griceův kooperační princip ...
- ▶ komunikační situace + kontext
- ▶ scénáře

# Sémantika a znalosti o světě

- hraje při interpretaci anafor často rozhodující roli
- sémantická plausibilita zvyšuje/snižuje pravděpodobnost některé interpretace, některé lze zcela vyloučit

After the [bartender] served [the patron], [he] got a big tip.

After the [bartender] served [the patron], [he] left a big tip.

- iniciální interpretace (hned)
- pokud pozdější informace vedou ke sporu:  
↪ reinterpretace (backtracking)
- **garden-path effect**
- význam slov
- znalosti o světě
- inference

# Sémantika a znalosti o světě

- If the baby does not thrive on raw milk, boil it.
- The FBI's role is to ensure our country's freedom and be ever watchful of those who threaten it.
- Stehlíková ustoupila od sbírky. Romové o ni nestojí.
- Klaus dostal dopis podepsaný Aničkou. Má ho policie.
- A: I ve Veselé vačici by mohla být volná místa.  
B: Jé, tam jsem ještě nebyla. Slyšela jsem, že tam chodí studenti. A že prý dobře vaří.
- 'I said disarm only!' Lockhart shouted in alarm over the heads of the battling crowd, as Malfoy sank to his knees; Harry had hit him with a Tickling Charm, and he could barely move for laughing.  
(*J. Rowling: Harry Potter and the Chamber of Secrets*)

## Sémantika a znalosti o světě

- Genau so sei es ihm vorgekommen, sagte Gauss, schließ ein und wachte bis zum abendlichen Pferdewechsel an der Grenzstation nicht mehr auf. Während die alten Pferde ab- und neue angeschirrt wurden, assen sie Kartoffelsuppe in einer Gastwirtschaft.  
(Daniel Kehlmann: "Die Vermessung der Welt: Die Reise")
- všechny tyto znalosti je obtížné shromáždit
- i kdyby byly k dispozici, bylo by obtížné v nich hledat
- AR je považováno za **"AI-úplný problém"**  
*AR je stejně obtížný problém jako naučit počítače myslet.*  
⇒ nutno si úkol zúžit



# Teoretické problémy

- John loves his wife. So does Bill.
- The man who gave his **[paycheque]** to his wife was wiser than the man who gave **[it]** to his mistress.
- If any man owns **[a donkey]**, he beats **[it]**.
- **[No one]** will be admitted to the examination, unless **[he]** has registered four weeks in advance.
- **[The man who shows he deserves [it]]** will get **[the prize [he] desires]**.

# AR algoritmy

- heuristická pravidla (70. léta)
  - ▶ SHRDLU – “block world” Terryho Windograda
  - ▶ [Hobbsovo syntaktické hledání](#)
  - ▶ jednoduchá pravidla, vzory, časté instance
- sématické teorie
  - ▶ centering, focusing – modelování lokální koherence
  - ▶ [BFP algoritmus](#)
  - ▶ výpočetně problematické
- knowledge-poor (90. léta)
  - ▶ kacířství motivované praktickými potřebami
  - ▶ založené na datech, která lze dostatečně úspěšně spočítat (morfologie, povrchová syntax, jednoduché sémantické třídy)
  - ▶ [RAP](#) – váhování
  - ▶ CoGNIAC (pouze 6 pravidel – vysoká přesnost, malé pokrytí)
  - ▶ MARS – váhování

# Naivní algoritmus – lineární procházení

- za antecedent je považován nejbližší předcházející výraz, který neodporuje zmíněným omezením
  - ▶ osoba, rod, číslo
  - ▶ syntaktická omezení, Chomského principy A, B, C
- předcházející věty lze procházet zleva doprava, nebo na základě syntaktických rolí
- filtrování pomocí sémantických tříd ...

# Hobbs syntactic search

- jako syntaktickou strukturu předpokládá frázové stromy
- X-bar theory (Chomsky, Jackendoff)  
X – complement – X' – adjunct – X' – specifier – XP
- algoritmus je definován jako procházení stromu
- začíná se v listu dané anafory
- podle kategorie aktuálního uzlu se volí další cesta
- prominentnější posice jsou procházeny dříve
- lze adaptovat na jiné formalismy
- jednoduché, ale nefunguje špatně

# BFP algoritmus

- založeno na teorii “Centering”  
(*modelování lokální koherence*)
- každá výpověď:
  - ▶ forward-looking centers (setříděné)
  - ▶ preferred center (ten nejvýše postavený)
  - ▶ backward-looking center
- cílem je nalezení koreferenčních vztahů, které představují nejplynulejší přechod center

- identifikace NP, filtrování nereferenčních, reflexiva atd.
- přidělí se iniciální váhy kandidátům (součet)
- při hledání antecedentu ke konkrétní anafoře se pro danou kombinaci váhy dále upravují (katafora, paralelismus, ...)
- antecedentem je kandidát s nejvyšší vahou
- při zpracovávání nové věty se všechny váhy podělí dvěma

<i>Factor type</i>	<i>Initial weight</i>
Sentence recency	100
Subject emphasis	80
Existential emphasis	70
Accusative emphasis	50
Indirect object and oblique complement emphasis	40
Head noun emphasis	80
Non-adverbial emphasis	50

# Pražské algoritmy

- hned několik algoritmů
- formulovány “na papíře”
- vyhodnocovány ručně
- jako RAP také váhovací princip
- modeluje aktivaci objektu v mysli posluchače
- zohledňuje se informace o AČV
- teoreticky logické, ale prakticky nepotvrzené

# Aktuální členění větné

(Topic-Focus Articulation, Information Structure)

Každá věta obsahuje dvě části:

- **Topic** (základ):  
to o čem věta vypovídá (kontextově zapojené)
- **Focus/Comment** (ohnisko):  
co se vypovídá o základu (nové; kontextově nezapojené)

(toto rozdělení může být triviální – pouze ohnisko)

~ aktivace promluvvého objektu v mysli čtenáře



## Aktuální členění větné II

- je v různých jazycích vyjadřováno různě
- v jazycích s tzv. “volným slovosledem” je hlavním nástrojem slovosled

Rohlíčky prý jsou dneska zvláště vypečené. Je tomu tak?

Ne, není tomu tak, Milosti. Ba naopak.

Vypečené rohlíčky zvláště dnes nejsou.

- slovosled má vliv na sémantiku, není tedy “volný”
- v mluvené řeči větný přízvuk
- Cizí jazyky:
  - ▶ germánské jazyky  
*It was Mary, who John called on the phone.*
  - ▶ finština  
*Poydällä on ruokaa.*

# AR a strojové učení

- statistika a strojové učení dnes v NLP převažují
- AR není klasifikační problém

předefinování umožňující použití std. ML metod:

- **1 instance**: dvojice anafora-antecedent
- **atributy**: knowledge-poor informace
- **cílový atribut**: 1 pro koreferentní dvojici, jinak 0
- velký nepoměr negativních a pozitivních instancí
- nutno část negativních instancí odstranit z trénovacích dat

**Gracie:** Last week my brother went out on a murder case, and you know, he found that man in an hour.

**George:** He found the murderer in an hour?

**Gracie:** No, *the man who was killed*.

**George:** Not only is your brother tall, but he's fast.

**Gracie:** And then Mr. & Mrs. Jones were having matrimonial trouble, and my brother was hired to watch Mrs. Jones.

**George:** Well, I imagine she was a very attractive woman.

**Gracie:** She was, and my brother watched her day and night for six months.

**George:** Well, what happened?

**Gracie:** She finally got a divorce.

**George:** Mrs. Jones?

**Gracie:** No, *my brother's wife*.

*(George Burns and Gracie Allen in "The Salesgirl")*