

Úvod do počítačové lingvistiky

- Počítačová lingvistika = obor, jehož cílem je popsat vlastosti přirozeného jazyka pomocí matematiky
PL
→ cíl: pochopit jak přirozené jazyky "fungují"
- Počítačové zpracování přirozeného jazyka = klasické aplikace metod strojového NLP
NLP
nicméně na velká jazyková data (korpusy)
→ nic nám to neřeká o funkci toho jazyka

• Oblasti PL

- Morfologie

- Syntaxe

- Semantika

- Formalismus

- Korpusová lingvistika

- Statistická lingvistika

- Strojový překlad

- Rozpoznávání a generování mluvené řeči

- Vyhledávání v textu

- Dialogové systémy

• Problémy

- přirozený jazyk je víceznačný

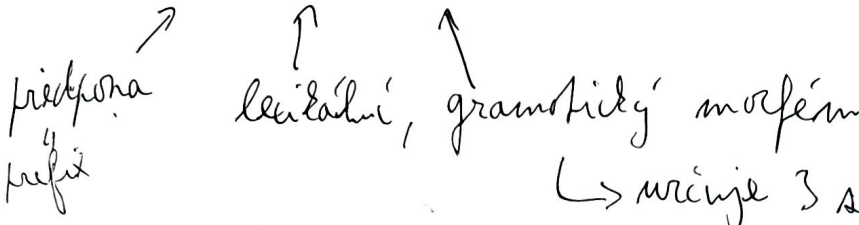
- slovo nese více svůj číslo, rod, pád... jednoznačně

Morfologie

= studium vnitřní struktury slov

morfém = nejmenší znaková jednotka jazyka nesoucí význam

za- hrád- ou = přípona



↳ určuje 3 sémata: pád, číslo a rod

- skloňování = deklinace
- časování = konjugace

Problémy

- homoslovné dublety = jiné tvary, stejný význam: na hradi / hradu
- homonymní tvary = stejný tvar, více rádků: stát (země, sloveso, sloveso)
- alternace = random změna hlásky: žena (podst. jména, přechodník)
- plnovýznamová (byl) × formová slova (jsem doma) - writi změna: viz → roční

Morfologické typy jazyků

- Analytické: slovo = morfém → číňština
- Syntetické: slovo > morfém → slovanské jazyky
- Pozysyntetické: slovo = řada → indické jazyky

máta {
 máte
 máteš
 máteš

Přístupy zpracování morfologie

⇒ rádků jsou dvě

- morfémy: slovo = řada morfémů
- lexémy: slovo = výskledek aplikace pravidel, co není rádků a vyhoví slovní tvar
- slova: hlavní roli hraje vzory

- ⇒ když mám rádků tvar + vzor, tak umím regenerovat ostatní tvary podle toho vzoru
- ⇒ tohle se reálně používá
- ⇒ v číštině je asi 250 jednovýznamových vzorů

angličtina

↓
 breaks
 broke
 broken
 "
 1 lexém

• Two-Level Morphology - 1980

- první obecný model zpracování přezněho jazyka
- V jazyce svůj slovník a pravidla, ale mechanismus morfologie obecný

- 2 úrovně

- lexikální
- fonková

- pravidla se uplatňují paralelně, nikoli seřazeně
- podmínky se mohou vztahovat k 1 úrovní nebo k obou úrovním

lexikální úroveň ↔ morfická úroveň ↔ fonková úroveň

- lexikální vyhledávání v trii a morfologická analýza probíhají současně

Česká morfologie

- příčinné značky, každá pozice ve značce má nějaký význam

- 13 kategorií: slovní druh, rod, číslo, osoba, čas, ...

lemma = jednoznačný identifikátor toho slova = základní tvar / kořen + index

• Morfologická analýza: dle toho slova

p. j. stát → stát-1

- vytvoří seznam lemmat a značek, která popisují jednotlivé možnosti, co to slovo může znamenat

- např. pokud více pádů má stejný tvar → více značek

- pokud to slovo má více významů → více lemmat

+ * lemma má nějaká jiná pravidla na detailní značek

- těch labelů může být spousta

• Morfologické značkování = tagging

- z těch mnoha značek vyberáme tu 1 správnou v daném kontextu

- statistické metody (učíme se o to)

• Částečná morfologická disambiguace - Oliva

- podle gramatických pravidel jsou některé konfigurace nelegální

↳ shoda podmětu s přísudkem, shodný přívlastek

↳ 7. pád po předložce s

- odebrávají pouze ty značky, které jsou 100% špatné

• Lemmatizace

- proces výběru správného rozkladního tvaru

• Stemming

- odříznutí koncovky

- na rozdíl od lemmatizace je různorodým tvarem kmen slova

• Generování

- známe lemma a kombinaci gramatických kategorií
⇒ chceme správný slovní tvar

• Značkování

- hledáme nejpravděpodobnější posloupnost značek pro danou větu

→ vyvíjíme Bayesův vzorec $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$

$P(\text{značka}|\text{slovo}) = P(\text{slovo}|\text{značka}) \cdot \frac{P(\text{značka})}{P(\text{slovo})}$, ale $P(\text{slovo})$ je pro všechny značky stejné

⇒ chceme to maximalizovat

→ čím více po sobě jdoucích slov a jejich značek bereme v úvahu, tím lepší, ale je to drahé a nejsou data

⇒ závislost pouze na jeho značce a pětí výsledku těchto značek na poslední značce

→ nejznačka = $\operatorname{argmax}_{\text{značky}} \prod_{i=1}^n P(\text{slovo}_i|\text{značka}_i) \cdot P(\text{značka}_i|\text{značka}_{i-1})$

• Skryté markovovy modely HMM

- metoda analýzy řád (posloupnosti událostí v čase)

→ my se aplikujeme na posloupnost morf. značek v textu

→ je to posloupnost rozhodnutí, co na sobě nemá závislost (gramot. pravidla)

→ Markovova hypotéza: kontext = # po sobě jdoucích věcí, co spolu souvisí, je možné zbraňovat na nějakou spojitou a pevnou délku

→ kontext délky 2 = bigramy, délky 3 = trigramy

- Skryté, protože některé vlastnosti se posloupnosti nejsou nevíte vidět - jsou tam slova, ne značky

- je to basically stochastický konečný automat

↳ náhodný

Jak se HMM používají

1) rozpoznávání (ohodnocení) statistického modelu

→ jsou dány parametry HMM, cíl je spočítat jaké je pozorovaná posloupnost X .

→ použití: rozpoznávání dráček - registračních značek aut

2) Rekognice

= hledání nejpravděpodobnější posloupnosti skrytých stavů

↳ máme model a posloupnost pozorování

3) Věření statistického modelu

→ máme strukturu modelu a trénovací data

→ chceme zjistit jaké přechody mezi stavy a jaké těch stavů

Kontrola překlepů

Požadavky:

- nalézt a opravit všechny překlepy
- ten výsledek by neměl být nesmysl
- když nějaké slovo neznám, tak to není chyba
- žádné false positives
- co nejvíce automatizace
- co nejrychlejší

} samozřejmě nereálné

Metrika úspěšnosti

Precision $P = \frac{\# \text{ true positives}}{\# \text{ positives}}$... kolik toho hlasím dobře

Recall $R = \frac{\# \text{ true positives}}{\# \text{ chyb v textu}}$... kolik % chyb jsem odholil

→ Precision je reálně důležitější, lidí neradi false positives

$$F\text{-measure } F = \frac{2PR}{P+R}$$

Metody kontroly pílepek

1) Porovnávání řetězců se slovy ve slovníku

- seznam všech možných slovích tvarů = word list
- slovník lemmat + morfologická analýza

⊕ spolehlivé a simple

⊖ pomalé, mávroané na kvalitě slovníku, neznámá slova = chybná slova
K zlepšení musí do slovníku přidat věirostel (ab to ve je bez m-analyz)

2) Porovnávání skupin znaků (dvojice, trojice) + nejake' radované kombinace

⊕ nezávislé na slovníku, rychle

⊖ vhodné chyb to nezachytí - chybná slova mohou být z okolních kombinací

Možná zlepšení

- uvážnat slovoati veniku chyby - blízké klávesy
- zohlednit statistiku chyby - která schvle podm. s přís. je hodně těžká
- zohlednit pravopisné chyby - mne x me, jsem x jme
- heuristika na oddělení chyb a neznámých slov
- zapojení syntaxe a sémantiky
- pracovat s kontextem - porovnávání korpasy

Jak se to reálně dělá

- špatné slovo → jak vybrat možné správné tvary?

→ Levenshteinova vzdálenost řetězců

- důležitá je přesnost (Precision) a rychlost

- kontrola na pozadí, každá úporovní ře se našla chyba

- věirostel vkládá do slovníku jen konkrétní slova tvar - není morfologie

System ASIMUT

= Automatizace Selece Informací Metodou Úplného Textu

- 2 moduly

• Jazykový modul

- není rozsáhlý slovník - používá sítě

- mnoho slov, která mají v základním tvaru stejný koncový segment, se stejně skládají

⇒ retrogradní slovník

→ pravidla kam jít podle rodu + čísla + pádu u konce základního tvaru

↳ náhled na vyplněný suffix, což dává ten tvar

→ jsou různé asi jenom slovy vyjímek

• Vyhledávací modul

- basically to vyhledává regexy

- výrazy z podst. a příd. jmen v základním tvaru + operatory

! významnost slova = jakýkoli tvar

-1- obě slova vedle sebe

-2- měřena měří slovy obsahují nejvýše 2 slova

-3- obě slova ve stejné větě

-4- -4- ve stejné vzdálenosti

- vyhledávání v nejdelších úvodních dokumentech

totaz: vzdálenost!, odstup! -3- rovinný! -1- domeč!

... jak daleko od sebe mají být domky

- problémy jazykového modulu

- proc není vždy víceměrně

- příliš hrubá klasifikace

- malý rozsah retrogradního slovníku ⇒ vylitím vyjímek

- nefunguje to tak spolehlivě pro slova

↳ koncové segmenty základního tvaru sloves jsou vícečetná

- Negativní slovník

- obsahuje negativní slova pro vyhledávání (spojky, citoslovce)
↳ odstraní se při preprocessingu textu

- Kondenzace

* slovník *trou* co není v negativním slovníku dostane
cílel: urychlit hledání

adresu
frekvenci výskytu

System MOZAIKA

= Morphemic Oriented System of Automatic Indexing and Condensation

→ hledá nejdůležitější termíny z nějaké oblasti (technické obrody) v textu

1) Standardní přístup k indexaci

- slovník klíčových slov, elementy indexovací sloz
- v návahu se bere četnost výskytu

2) Přístup MOZAIKY

☺ řada konvok a přípon nes ryčan

- *ic*, *ac*, *ica*, *er*, *or*, *do*, *me*, *graf*, *fon*, *stop* ... nástroj/přístroj
- *ac*, *ke*, *ac* ... procesy
- *ost*, *ita*, *me* ... vlastnosti
- *ac*, *ec* ... účel

angličtina: *er*, *or* = konstel deje

- *tion* = činnost

- *ity*, *ness* = vlastnost

→ pro pokrytí sémotické oblasti elektických obrodů ~ 800 přípon

Algoritmus MOSAIK 4

1. lemmatizace + morfologická analýza vstupního textu → knačky
2. prohledání nalezená lemma a odstranění ty, jejichž Emen nemá vztah k dané sémotické oblasti (například)
→ vznikne malý negovní slovník
3. Chci odhalit několik slovové termíny
operací resitrací TESLA KC 415 → resitrací
↳ podtermíny se také spočítají
⇒ Syntaktická analýza pomocí jednoduché gramatiky v Q-Systemech
↑
kordéji
4. přiřazení váhy na základě místa výskytu v textu
nápis je nejvíc
1. odstavec > další odstavec
1. věta v odstavci > random věta
5. udělám nějakou normalizaci vah vzhledem k délce dokumentu
↳ aby se mohla proměřit relevance různě hlubokých dokumentů
⇒ výstup = 10 nejvýznamějších termínů spolu s četností výskytu

⊕ nepotřebuje slovník odbojích termínů, jen množina přípon a koncovek + negovní slova / pravidla
lokální syntaktická analýza tvoří lepší hledat termíny

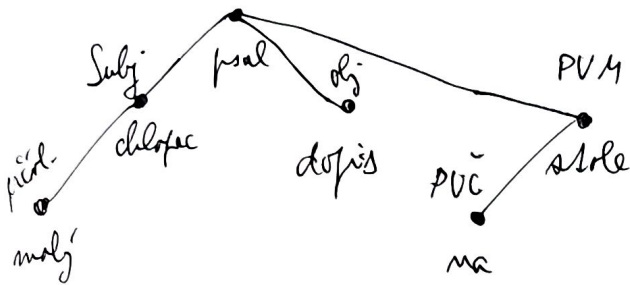
⊖ vymyslet ty přípony a pravidla je pracné
neračítá se rájmena
↳ ten resitrací TESLA tam měřím napsal 1x
a teď se na něj 10x odbozovat rájmeny
⇒ 1 výstup

Syntax

prísudek = koreň

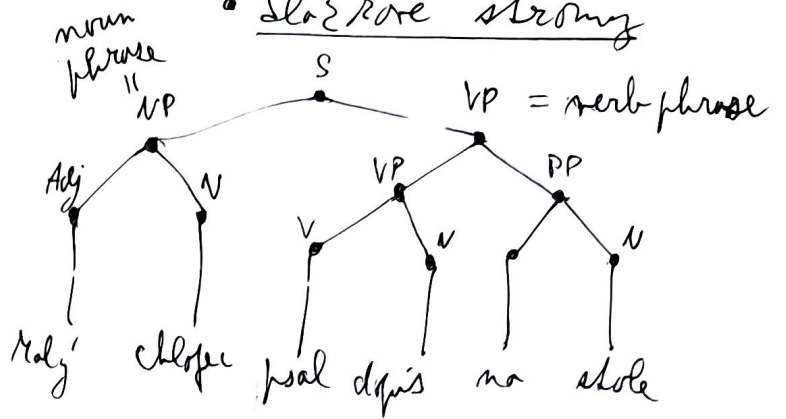
závislostní stromy

Před



- dobře zachycuje vztahy
- strom zachycuje pořadí výpočtu
⇒ není jasné, jak ho reálná
- není jednovácný
- ne všechny vztahy se řeší
se řeší dají popsat → Petra a Pavel

slučkové stromy

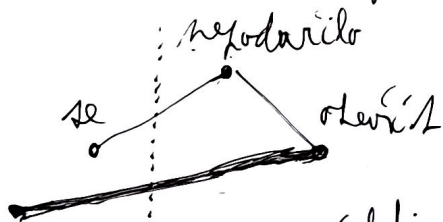


- derivační strom bezkontextové gramatiky
- méně přehledný, zbytečné uzly
- přirozený jazyk není bezkontextový
⇒ spojuje slovní druhy do frází

$V + N \rightarrow VP$

Neprojektivní konstrukce

= dlouhé hrany ... Za hrana v závislostním stromu jde dopředu
řeta: Soubor se nepodařilo otevřít



→ se začíná na konec řety hrana

Soubor → neprojektivní ≡ vertikála má 2 průsečíky
→ složitý strom má problém: ty slova spolu nejsou seči

kontextová gramatika

X

Bezkontextová gramatika

- nějaká pravidla
- se řeší ty měl být
hodně a přísudek

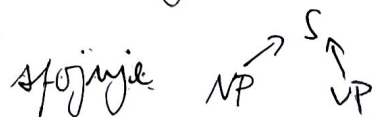
- pravidla na spojování znáčel
- dle terminálních znáčel

Generování

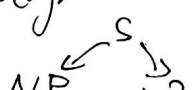
X

Analýza

- z řety dělá strom

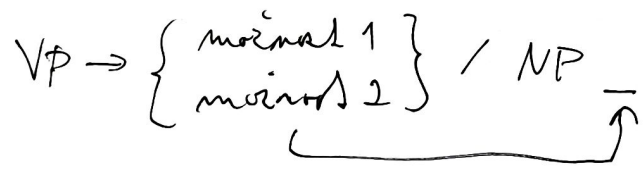


- výrok stromu, do kterého se za
znáčel pak dají dořadit slova
rozvíjím

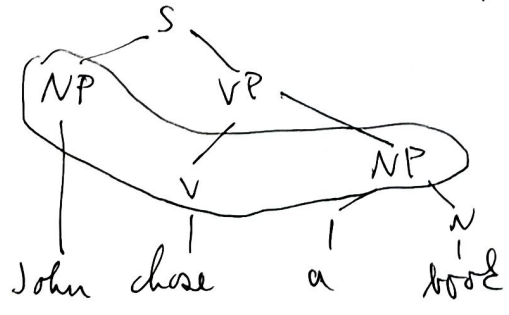


Transformační gramatika - Noam Chomsky ~ 1965

- převratný způsob formálního popisu přirozeného jazyka
- princip: kombinace x hloubková syntaktická struktura
 - parobový strom - klasický složitkový strom
 - hloubkový strom - zachycuje pouze důležitější vztahy
 - ↳ změna pořadí slov ho nezmění
- komponenty
 - bazě - bezkontextová pravidla pro generaci složitkových stromů = frázevé utvářecí
 - Transformační k. - transf. pravidla operující na frázevých utvářecích
 - spojují / mění ty stroměcké
 - pravidla - obligatorní
 - potulbativní
 - fonologická k. - regexy co dávají řečnickému morfému fonetické interpretaci
- Generativní procedura - jak generovat stromy



- můžeme to přepsat jen když předtím je NP
- nedokážeme zachytit různé varianty vět ← řazení / oznamovací
- transf. složka má pravidla jak z předkládajících fráz. utvářecích vytvořit funkční strukturu věty



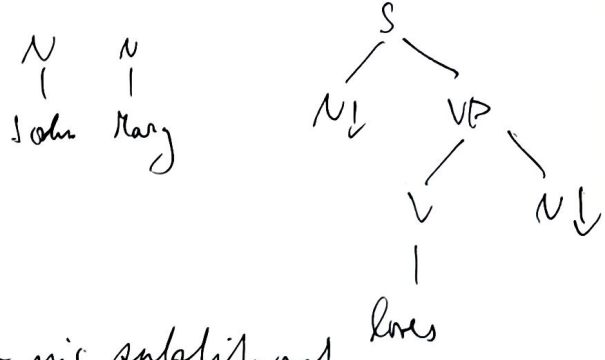
$NP_1 - V - NP_2 \Rightarrow NP_2 - was - V + -en - by + NP_1$

Book was chosen by John

- Ta teorie se hodně změnila, Chomsky vydával další verze
 - 1965 - Standard Theory - sobě
 - 1990 - Teorie minimalismu - pouze rozina logické formy a fonetická realizace

Tree Adjoining Grammars - substituce stromů

- stromečily jako **el.** struktury
- šifera ↓ znací možnou substituci
- strom ~ sym. struktura věty
 - vznik postupem substitucí
 - přes konce, edge ne může nikam nic substituovat



Legal Functional Grammar

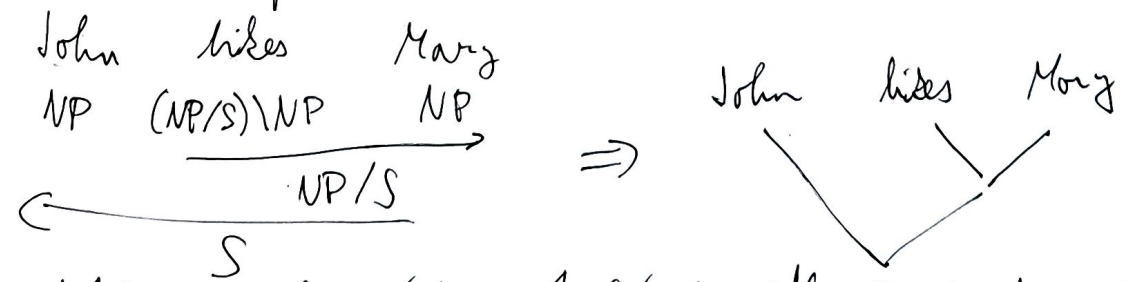
- C-structures ~ spojování slov do frází
- f-structures ~ reprezentace funkční roztoky ve větě (vždy sloves)
- ↳ matice

al	ho
ri	de
ku	no
ky	sy

→ každá matice byt i jiná f- struktura
 ≠ C- struktura má právě 1 f- strukturu
 ale 1 f. s. může být ve více C. S.

Kategoriální grammatiky

- ≠ každý slovo dostane kategorii
- kategorie popisuje syntaktické vztahy slova
- obecný formát α/β nebo β/α ... komitlo říká slovo / spáro
- likes má kategorii $(NP/S)/NP$ → díra na něco být
- 2 hlavní pravidla: $X/Y + Y \rightarrow X$, $Y + Y/X \rightarrow X$



- ⊕ nepřekvapuj obecná gramatická pravidla ≠ spornu výměr
- ⊖ strašně moc kategorií

Unifikáční gramatika

- popis vlastností objektů

→ objekt ~ množina vlastností (máček vl : hadula vl.) = sestava rysů

→ sestavy rysů se dojí unifikovat

↳ lze to, když vlastnosti 1. sestavy neodporují sém 2.

→ sestava rysů většinou popisuje nějaký sign. jer (shoda f. p. v.)

→ hadula vlastnosti může být doba sestava nebo jistě na ni

shoda formátu + přísudek

↳
$$\left[\begin{array}{l} \text{podmět: } \left[\begin{array}{l} \text{osoba 2} \\ \text{rod ž} \end{array} \right] = 1 \\ \text{přísudek: } 1 \end{array} \right]$$

Problém: lze unifikovat nesouvisející rysy

→ sestavy mají sém → ten určuje její vlastnosti

↳ Typované sestavy rysů

↳ např: sloveso, podstatné jméno, ...

HPSG

- obsahuje principy, gramatická pravidla a slovníkové položky

- slovník dobře strukturován, spousta info

- el. jednotka = znak (sign)

- slova a fráze jsou různé podtypy znaků

→ slova má 2 rysy

• fonetický ... znak

• syntaktické a sémantické informace

Nástroje na syntaktickou analýzu

• Augmented Transition Networks

- čím jeká slova fráze co potřebuje

$S \rightarrow NP \text{ a } VP$

$NP \rightarrow Det \text{ a } N \dots Det = \text{člen "a", "the"}$

$VP \rightarrow V [NP]$

The girl saw a boy ... kořen = S a S potřebuje NP a VP

\Rightarrow SEEK NP ... NP potřebuje Det a N

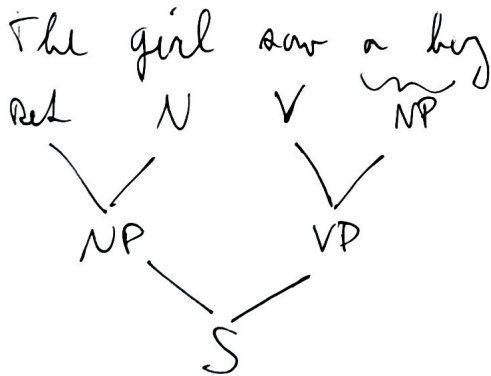
\hookrightarrow CAT Det ... Pokud následuje Det, jdi do dalšího stavu
 CAT N ... Pokud následuje N, jdi do dalšího stavu

\Rightarrow SEEK VP ... potřebuji V a volitelně NP

\hookrightarrow CAT V ... pokud následuje V, pokračuj

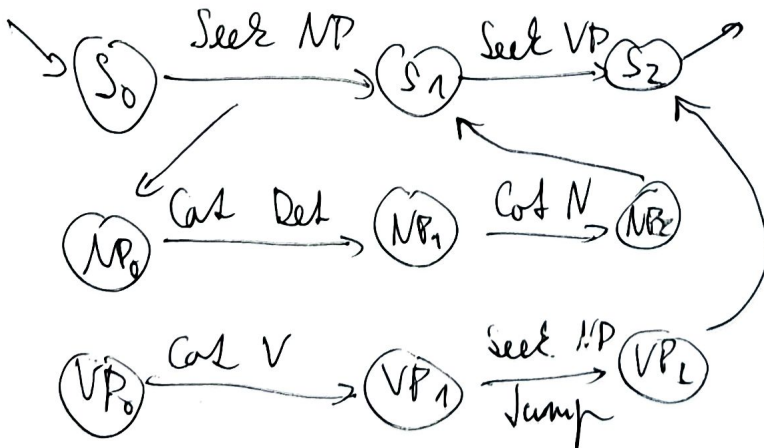
SEEK NP

JUMP ... přejdi do dalšího stavu navíc bys mě nehledal



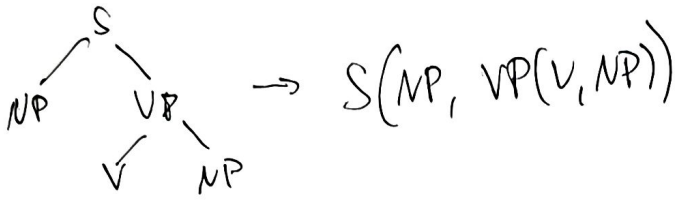
I Syntax beam

- SEEK = přechod do předstí
- CAT = předstí do dalšího stavu pokud najde co hledá
- JUMP = přechod dál bez hledání



Q-Systemy

- formalismus pro konstrukci grafů
- stromy ukládá lineárně



3 typy objektů

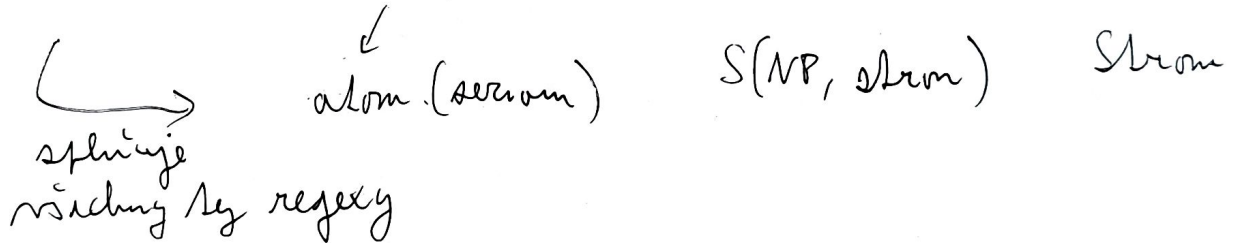
- atomy = konstanty ... písmena A-Z
- stromy = stromy ... písmena L-N
- seznamy stromů ... písmena V-Z

+ má nějaké operátory na těchto objektech

- zapisuje strom jako regex

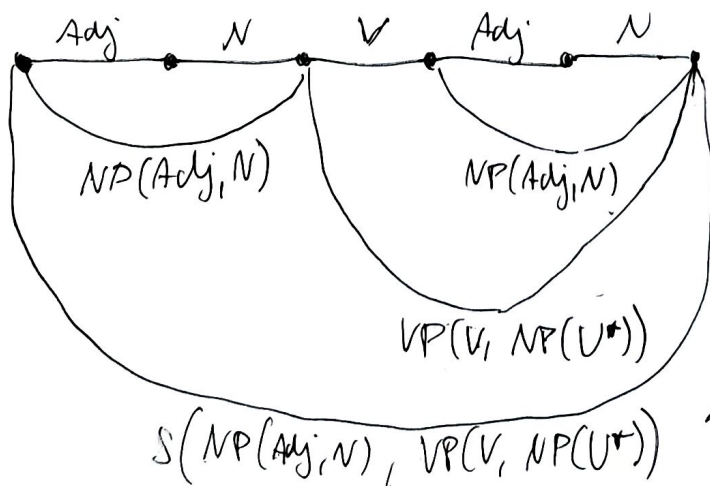
* = proměnná

$$S(NP, VP(V, NP)) \sim A^*(U^*) \text{ nebo } S(NP, L^*) \text{ či } N^*$$



pravidla

- $Adj + N \rightarrow NP(Adj, N)$
- $V + NP(U^*) \rightarrow VP(V, NP(U^*))$
- $NP(U^*) + VP(V^*) \rightarrow S(NP(U^*), VP(V^*))$



ten výstup můžeme poslat do dalšího Q-Systemu

- čistění stromu

- odstraníme nadbytečné hrany - dílčí mezinóžské dty
- vředy spájají hrany ... písmena / výrazy = hrany (na obě strany)
- 1. odstraním vředy hrany, co jsem psal na levé straně nějakého pravidla
- 2. odstraním slepé vředy

Funkční generativní popis - motivy poznět

- stratifikaci teorie - 5 rovin (Telegrafická TG)
- jednota na vyšší rovině reprezentuje funkci jednotky na nižší rovině (Ponchová, Morfémotická, Fonologická, Fonetická)
- teorie valence
 - ↳ vazy jsou vyřadovány nebo problémy říditými slovy

- TG rovina zachycuje strukturu

→ kniha byla vyd. nabl. = nabl. vyd. knihy

→ 2 druhy členů ^{objekt slov}

- Atanty: knožel (agens), Posiel, Adresát, Origo, Efekt

↳ $\&$ může být v 1 větě jen 1x

- volná doplnění klidně vícekrát

→ členy dělíme na

- obligatorní - nesmí na TG chybět

↳ podmět sam musí být, i když na funkční rovině může chybět

- facultativní

→ obligatorní členy jsou pro $\&$ sloveso jiné!

→ Moji přátelé přijeli

- kam? ... to je důležitá ⇒ obligatorní
- odkud? ... facultativní

Valenční rámec = seznam atantů (i facultativních) a obligatorních volných doplnění

Kontrola gramatické správnosti

Také šli do práce. ← Co kdo kam donesl?

- problémy v češtině

- nepřes. konstrukce, shoda podm. s příslovečnou
- interpretace i význam má / nemá

• chybové vzory

- pro jazyky s levným slovesem
↳ chybové konstrukce jsou blízko u sebe

• gramatika

- kontrola podle pravidel
- nelze rozeznat, zda je to chyba, nebo zda je naše gramatika neúplná

} 2 přístupy

• Metoda redukční analýzy

- větu redukují tak, že vždy odstraní něco, co je 100% správné
⇒ neaním chybnou větu opravit ani správnou rozeznat

1. Náš hořejší vícera je velké snaze opit vyhrály.

2. hořejší vícera je snaze vyhrály.

3. hořejší vícera vyhrály.

4. hořejší vyhrály.

→ např. velké snaze, ... shodný přísloveč. ... kontrola shodny
... je to dobře → odstranit

⇒ chyba dostanu blízko k sobě ⇒ chybové vzory

- musí se ošetřit speciální případy

výstup: minimální chybová konfigurace ... vyšetřovaná věta

! formé hledání chybných reálním složitějších konstrukcí nefungují
> které divička chtěla dostat na úroveň? *

↳ min. chybn. konf. i ale dlouhá hrana

! homonymie ... číselní jako orgán x číselní jako měřítka
rod ž rod m

• RFO DG = Robust Tree-Order Dependency Grammar

👁️ 1 gramatické pravidlo může popisovat správnou i chybnou konstrukci zároveň

- výpočet se provádí, interpret gramatiky rozhoduje, jak se bude stejná gramatická pravidla používat

- příklad pravidla: ps předložka A je 7. řád
shoda shodného přívlastku

řádky: pozitivní - tvrdé podmínky → musí to splnit
negativní - měkké podmínky - nemusí to splňovat
projektivní = říci podle sebe

→ řádky: z řádky vždy vypadne nějaký les

1, pozitivní projekční

- pokud je ta věta správná a bez neprog. část
tak to vytvoří celý strom

2, neg. proj. nebo poz. neproj

3, neg. neproj

- takhle se rozkládá všechno

→ rozhodne se více stromů, jak to může být

dílčí a náhodní - neprojektivní strom

↓
strom s chybnou posláním

• Lang GR

- redukuje snailky po morf. analýze → disambiguace
- pozitívni a negativni disambiguační pravidla
- pravidla mohou mít neomezený kontext
- jsou vzájemně nezávislá, nespořádaná a jsou uplatňována v cyklu
- 4 části: kontext, disambiguační část, report, alee
- ⇒ $cont_1 disamb_1 cont_2 disamb_2 \dots cont_n disamb_n cont_{n+1}$ report action
- ⇒ odhali podstatné věty - nějaké slovo na konci nemá řádnou reakci

Korpusová lingvistika

korpus = velké množství textů s nějakou přidanou informací
↳ morf. nebo synt. analýza

• Studie struktury jazyka

- identifikuje jednotky a třídy jazyka (morfémy, slova, fráze...)
- popisuje jak se mohou kombinovat na velké

• Studie užívání jazyka

- když mám hodně náhodně nerovně vybraných textů, tak mi to dává nějakou objektivní informaci o tom, jak se ten jazyk reálně používá
- měl by to být reprezentativní vzorek jazyka

• výběr vzorku a reprezentativnost

- jazyk je nerovinný, ale korpusy konečné
- ⇒ vzorky musí být nějak vyvážené

• konečná velikost

- s různými monitorovacími korpusy, kam se pořád něco přidává
- mají korpusy fixní velikost
- jásem si se jich rozhodl, než se zmenší

• stručná číselná forma

- měly by se dodržovat nějaké standardy - UD

• Brown Corpus of Standard American English

- 1. moderní elektronický korpus
- 1 milion slov textu v americké angličtině vyhledávaných v roce 1961
- náhodný výběr vzorků, reprezentativní
- morfologické značky

• PennTreebank

- 1. a nejznámější syntaktický anotovaný korpus
- ~ 1 milion slov
- články z Wall Street Journal
- ⇒ úplně k ničemu, je to bohemský slang, vůbec nereprezentuje angličtinu

• British National Corpus ~ 100 mil. slov, morfologická anotace

• American National Corpus ~ 22 mil. slov, morfologická anotace

• Corpus of Contemporary American English ~ 410 mil slov —

↳ 20 mil. slov za 4 roky mezi lety 1990 a 2010

• Oxford English Corpus ~ 2 miliardy slov —

↳ gramatické anotace byly nejisté

• Negra Corpus - němčina - syntaktická anotace

• EUROPARL - 11 evropských jazyků, ~ 40 mil slov / jazyk

↳ skrovň paralelní korpus

→ dobré pro trénování překladačů

→ důležitý lidský příklad sítě pro větu, ne jako celek

↳ nezávisle na kontextu

- Česky národní korpus - UK, MUM, Ústav pro jazyk český
 - morfologická anotace
 - dneska ~ 5 miliard slov
 - každých několik let vyjde nový korpus se slovy za ten uplynulý čas
 - ↳ název SYN2020 = výsíl 2020, slova z 2015-2019

• Průběžný závislostní korpus

- chtěli syntaktickou anotaci jako TreeBank
- teoretický základ je Finkel's generativní popis FGP
- 100 000 vět ~ 1.25 mil slov → málo produktivní ČNK
- anotace

- morfologická

- syntaktická pomocí závislostních stromů

↳ chtěli ukázat, že jsou lepší než složitkové

- FGP ⇒ těch 5 rovin

↳ anotace různých jazykových detailů

↳ v TG rovině se kromě stromu kreslí i různé vodorovné reference, např. něco odrozejí rájmen

• Prague Arabic Dependency Treebank

⇒ aby dokázali, že závislostní stromy jsou super, takže pomocí nich udělali korpus arabského

• Universal Dependencies

- díky různým jazykům autorůly stejné konstrukce různě

⇒ snaha anotovat stejné věty v podobných jazycích stejně

⇒ hodně lidí začalo experimentovat závislostní korpusy, ale každý jinak

→ PZK byl moc detailní → UD zjednotěná schéma

→ dneska se běžně používá UD

LINDAT ← všechny korpusy jsou tam

Sémantika

- syntaxe jen říká, jestli je to gramaticky OK
- sémantika říká, jestli to dává smysl

Kulatý obdelník zeleně masírovat kovář vedle minulý.

- význam ≠ pravidla

↳ nepravidla sdílí. mají svůj význam

↳ ne vždy je možné určit pravidla

- není vždy jasné, zda jde o věty se stejným významem

• Pozornosti ho dobruho

• Pyl jimi pozornost dobruho

- vyplývání - pravidla věty mají důsledky

↳ ! Snímací pom. ptáci \Rightarrow mají křídla a létají

formální jazyky \times přirozené jazyky

pravidla = význam

nefunguje to

• Fregelův princip kompozicionality

- význam složeného výrazu je jednorázově určen významy jeho částí a způsobem jejich kombinací.

• Lexikální sémantika

FORMÁLNÍ \rightarrow význam slov popisujeme pomocí nějakého meta-jazyka

\rightarrow hodné kategorie: král: muž+, panovník+

královna: muž-, panovník+

\rightarrow problém: mnoho kategorií nestací

hodné kategorie: některá slova jsou nerozhodnutelná

ani
moc

Přirozený - stejný nebo jiný - jako se slovíkem - lepší

! význam věty závisí i na kontextu

Problém: na významělemí mohou být potíže kvůli nesnáma slova

Rěšení: na popisy se používá listy jen 2000 nejčastějších slov jazyka

👁 význam slova je více než jen jeho kontext, popisuje
tvaru významových tříd

• Ontologie = množina tříd objektů na klasifikaci obecně univerzálně

matematika ∈ přístroj ∈ kuchyňské nádobí ∈ nádobí ∈ fyzické objekty

→ velmi těžké to udělat dobře

→ problém, že význam slova není jednoznačný

↳ schůzka, hlava

kvantifiky
velikost
vlastnosti
⋮

ontologie / veškeré - obecné významy

↓
omezené - více specifické

↳ třeba ontologie vlastností

• Sémantická síť

+ příloha

• Word Net → databáze anglických slovesných a přídavných jmén a sloves

→ dělení na třídy ekvivalence = Synsety

↳ & synset vyjadřuje určitý koncept

→ různé mezi synsety - sémantický ~ lexikální význam

→ slovo může patřit do více synsetů

~ 155 000 hesel

~ 117 000 sloves

• Euro Word Net

- několik jazyků

- veškeré ontologie - 63 nejdůležitějších jazykové sémantických konceptů

↓
mnoho základních konceptů - 1000 základních konceptů

↳ tvoří jádra sítě slova, jazykové sémantické

- dobré pro automatický překládání

- dobré pro práci se sémantickými vztahy - hlavně synonymy

• Reprezentace významu věty

• Predikátová logika 1. řádu

- jakoby to funguje, ale věty nejsou formule

... nezvládná to kvůli kvantifikaci

↳ possible (F), believe (x, F), true-at-some-point-in-future (F)

• Presupozice

= předpoklad, co musí být pravdivé, aby věta vůbec měla predik. hodnotu

"Jupiterův měsíc má oválné pruhy."

↳ Jupiter musí mít právě 1 měsíc

• semantics = Fuzziness - T/F nestací ... Pavel je mladý.

• Existence a intence

- Intence = definice

- Existence = věci, co splňují tu definici

• Cena Big Macu je 20 Kč.

Nepravda
||

pokud to nahradíme za intenci → 125 Kč je 20 Kč.

• Myslím, že cena Big Macu je 20 Kč.

↳ není ekvivalentní k "Myslím, že 125 Kč je 20 Kč."

• Přístupy k sémantice

1) Modelová-teoretická sémantika

pravdivostní podmínky jsou vztahové ke konkrétnímu modelu
→ např. montagueova gramatika

svět



2) Kompoziciální sémantika

- vychází z principu kompozitability = význam slovního výrazu závisí na významu jeho částí a na sémantickém výrazu

• Montaguelovská gramatika

- základem je první logika

- Montague byl přesvědčen, že sémantika přirozených a formálních jazyků se moc neliší

ale: sémantická pravidla jsou svázána se syntaktickými

→ do jisté míry se mu to podařilo

- pravidla připomínají kategoriální gramatiku

↳ slovesa (a se chovají podobně se sestupují do kategorií a zacházejí se s nimi stejně

• TIL = Transparentní intencionální logika

→ logika 1. řádu větací

- vyhodnocuje formou sémantiky možných světů

↳ věta může mít různý význam v různých světech

→ také rozhoduje čas, kdy se to odehrává

- používá lambda kalkul - něco z logiky

ω = svět

τ = čas

$\Delta \omega \wedge \tau (Z_{\omega\tau} F_{\omega\tau} H_{M_{\omega\tau}})$

↳ tohle popisuje velkojazyk věstí

→ kromě možných světů rozhodujeme i měrou individua

• studentka Alena - nálepková individua

* • ministr zahraničí - individua s měrou

☒ • hecni mě - věták mezi dvěma individua

☉ • ministr financí

⇒ měrou vyotřít větu "studentka Alena si myslí, že MF > MZ"

↳ je to nějaký objekt $O_{\omega\tau}$

↳ každému světu a času přiřadí nejvýše 1 podstatný bod

Rozprávania v slohu v textu - Anafory

→ odbovaci rájmena - problém pri príklade

↳ ústie "on" sa môže preložiť na { muž: on → he
stul: on → it

⇒ potrebujeme vedieť na čo sa rájmena odkazujú

Anafora = výraz jebo interpretácia závisí na kontexte

delem:

- Exofora = odkvvaní mimo text. → "Koho ho?"
 - Endofora = odkvvaní v rámci textu
 - Anafora (zpätne) ... ja, fakt se to puzo obje
↳ "Petr vyradil rajenstvi. To nemel dlat"
 - Kofofora (dopredu): "Kdyz se slovi, neni s Petrem zadna rci."
- Anafora ... rici se hlavne ona, vstah predchude & naslednik

1) Zajmena a nevyjadrené vetné členy (prednik)

Petr si koupil vstupenku. Pak ji dr Epsy. Byla divaci.

2) Určité jmenné skupiny

Elektronický rezistor Tesla 1000. Toti rozicev...

negativní

3) Eklipsa

Petr přinesl dva stoly, převrátil a korrý.

Petra puzde dr kina. Sirda stoly.

4) Spejvaní výrazů

napišebol, jedhot - drubel, nejchise - potom

⇒ velmi obtížné: pi křeba aplikovat více metod zároveň

Vyvození:

• rielarání infunaci a textu - MOZAIKA

• automatický příklad - spore příkladů rájmen "on" → it/he

• dialogové systémy

Réšenie úloh

- 1) morfologické znaky - napr. v rámci shody v rodí
- 2) syntaktická štruktúra vety
↳ pro učenie možných kandidátov na predchádzajúce
- 3) statistické metódy
- 4) aktuálne členenie - viz alg. sdílených znalostí
- 5) rozšíriteľné formálne znalosti - ontologie, sémantické siete, ...

• Číslova sdílených znalostí

- modeluje časovú znalosť, čo mluvca predpokladá, že sľubí a poskytnie
↳ mení sa podľa toho, čo je rovná r , centru poznatku
- jednoduchá pravidla
→ veta má najmä obmedzené a jadro
- a^n znamená, že a má dôležitosť n

Príklad pravidel:

Pravidlo referencie v rámci: n sa zmení

Pravidlo a je r obmedzené: $a^n \rightarrow a^0$

Pravidlo a je r jadro: $a^n \rightarrow a^1$

Pravidlo $a^n \rightarrow a^m$: objem asociatív a obdĺžnikového a^{m+2}

Pravidlo a nemá význam ani asociatív: $a^n \rightarrow a^{m+2}$

Skrojinyj priedklad

problém: potrebueme znasť morfologiu, nejde priedklad slova od slova

- pravidla prerozpisu: jeden písl = číslci
- ustálená spojenia: veľký mas = mašín
- ne vždy je možné namapovať slova na druh jazykúch 1-1
- homonyma
- složené výrazy: airport long term car park courtesy vehicle pickup point

Metody prvej generácie

ANALÝZA
TRANSFER
SYNTÉZA

- drojijazyčný slovník, priedklad slova od slova

- ve slovníku nejsou celá slova, ale morfémy

prediction → priedklad → postediting

- jeden jazyk se zrohí jako první pro priedklad do vice jazyku?

- Georgetownský experiment 1955

- velmi stručný ověření doméy a pozicií jazykúch frází
→ ale fungovalo to

- 1955-1965 boiling výraz

- 1966 výzva ALPAC

→ je nutné dlouhodobě investovat do kvalitativního ling. výzkumu

→ důsledkem: konec podpory v USA

↳ ve Francii, SSSR, Kanadě dlel pokračovat

- TAM-METEO - priedklad meteorologických zpráv A → F

↳ dobře definovaná a bodu omezená podmínka jazyka

↳ používá Q-systemy

- SPSTRAN - priedklad dokumentú EU mezi ~20 páry → fungovalo jen A-F-N

- EUROTRA - projekt EU, neúspěšné: 72 páry - negativní efekt jako ALPAC

- VERBMOBIL - priedklad mluvené řeči, německý nástupce EUROTRY

↳ omezení rozboru: plánování píslí schůzky dle dohodnutí

• Skrojem předpřevaný příklad = CAT

- komerčně mnohem úspěšnější než auto-sikar - už od 80. let

- příkladová fanéti

↳ obsahuje páry příbuzných segmentů - typicky vět + metoda, umožňuje efektivní správu fanéti

→ nové vstupní věty se vyhledávají v fanéti

→ před samým koncem, předložitel dostane návrh příkladu se zvýrazněnými odlišnostmi a procenty podobnosti

⇒ vhodné pro opakování příklady méně atraktivních textů - technické manuály
nové verze

- 2 přístupy

1) příklad věta pro větu

- předložitel vidí kontext toho co předkládá

- každý výsledek se musí přeložit znovu

→ IBM Translation Manager

2) příkladová Arbutka

→ protože její opakování

- kontext není vidět, ale v věta se předkládá jen jednou

→ Réjã Vu

• České předložce

• RUSLAN - předložce manuálně sdílených počítačů ~ 1985

- 1 věta ~ 4 minuty

- dvojjazyčný slovník, lokální přístup: analýza, transfer, syntéza

- Q-systemy

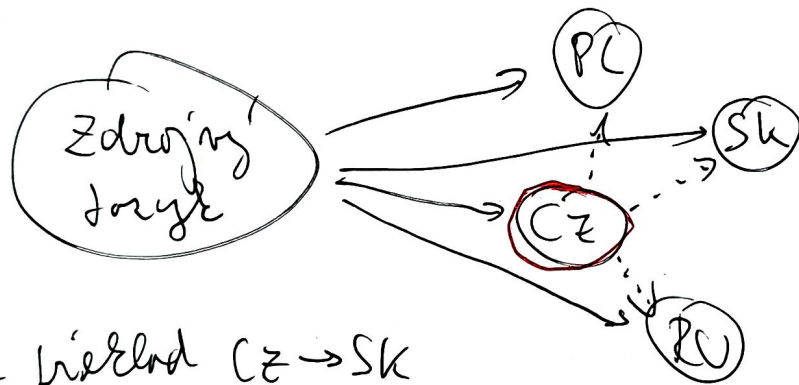
- transkukční slovník

⊗ některá slova s rěcto - lokálním ráčlucem se může přeložit přím

-aa → -acija: industrializace → industrializacia

• System Cosillo

- lokalizace velkých sv. systémů
- nahrazení lokalizace z jazyka Synté odlišného příkladem z jazyka blízkého
- příklady se mezi velmi blízkými jazyky - výsoké podobnosti příklad
- statistické rovnovážné čistoty (hlavně automaticky)
- dvojjazyčný systém



→ : lidský příklad
- - - -> strojový příklad

- příklad CZ → SK

↳ stejná syntaxe

→ většinou shodné pořadí slov ve větě

→ různé slovíčky a morfologie

• PC Translator 2003

- tradiční přístup, porovnání pravidla
- asi nejlepší úroveň komerčního systému

Statistický strojový přelod

- věci se dějí předložit více způsobů : 3 způsobů $\text{im}(EN) \rightarrow F$

→ relativní četnost máš rozjít

→ Bayesův vzorec $P(x|y) = \frac{P(y|x)P(x)}{P(y)}$

$$P(A \cap B) = P(A|B) \cdot P(B)$$

$$P(A \cap B \cap C) = P(A|B \cap C) \cdot P(B|C) \cdot P(C)$$

• Modelování jazyka

úkol: předpovědět další slovo v běžném textu

→ na základě historie h předpovědíme další slovo : $P(w|h)$

cíl: spočítat jaké celé věty $P(\text{věta}) = P(w_1 \wedge w_2 \wedge \dots \wedge w_m)$

→ N-gramy:

$$P(w\text{-tice}) = P(w_1 \wedge \dots \wedge w_m) = P(w_m | w_1 \wedge \dots \wedge w_{m-1}) \cdot P(w_{m-1} | w_1 \wedge \dots \wedge w_{m-2}) \cdot \dots \cdot P(w_2 | w_1) \cdot P(w_1)$$

problém = dlouhá historie \Rightarrow fakticky obří data

→ reálné se používají bigramy, případně trigramy

$$P(W) = P(w_3 | w_2 \wedge w_1) \cdot P(w_2 | w_1) \cdot P(w_1)$$

• Vyhledávání

- problém je velikost data

→ slovník V , $|V| = 40000$ slov \Rightarrow velikost modelu $|V|^2 = 6,4 \cdot 10^{13}$

↳ standardně přečtené křivka na $\sim 10^8$ slov

\Rightarrow obrovská moc nulových slov = řídká data

řešení: nahradíme nulovou část nějakou velmi malou hodnotou

↳ faktor ale nerovnoměrně opakování nevhodných kombinací

• Princip statistického překladačů

= paralelní korpus = trénování musíme dobře přeložit věty
→ ač to nedávána převládající metoda, deska spíš neurvaly

- Metoda rásuměleho kanálu

→ chceme překladač $F \rightarrow A$

⇒ hledáme model $P(a|f)$, a vyjádříme $P(a)$ řetěz a na F vidíme f

$$P(a|f) = \frac{P(f|a) \cdot P(a)}{P(f)} \quad \rightarrow P(f) \text{ je stejná pro } f \text{ a}$$

⇒ obrátili jsme de-facto směr překladačů

⇒ hledáme 2 modely: $P(f|a)$, $P(a)$

→ vlastně předpokládáme, že jsme dostali větu přeloženou z A do F
a hledáme její správný originál

→ językový model $P(a)$ může být trigramový model
rozložený na mnohem rozsáhlejších korpusu určitého jazyka

→ překladačový model $P(f|a)$ je rozložený na paralelním korpusu

- podstata: ↳ mnohem menší

• překladačový model budujeme v opačném směru

• jazykový model odfiltruje neprovozně překlady a
vytvorí duhy překladačového modelu

→ pozná, jestli se anglická věta dáva smysl

→ vyhledává totič "force a heckle" věty, nemá vztah k originálu

• hledání překladačových hypotéz (dekodování) je toby svůj problém

- dostane větu:

- máme na výběr z různých překladačů

- pro f spočítáme $P(a)$

} vybereme \rightarrow max

• Evaluace systémů automatického překladačů

→ metrika BLUE

- přesnost překladačů v n-gramech

⇒ máme 2 kandidáty na překladač ... který je lepší

⇒ pro tu větu máme popsanych několik správných překladačů

• unigramová přesnost

$= \frac{C}{N}$, $N = \#$ slov v daném kandidátu

$C = \#$ slov z kandidátu, co se vyskytly

→ alespoň jednom referenčním překladačů

• n-gramová přesnost

→ dvojice, trojice ...

⇒ $N = \#$ n-tic v daném kandidátu (co na sebe navazují)

$C = \#$ n-tic z kandidátu co jsou v nějakém referenčním překladačů

• BLUE = celkové skóre

• penalizace za strukturu = Breviety penalty (BP)

- pokud bychom mohli v úrovní jen $\#$ n-gramů, tak by se favorizoval kratší věty, protože n-gramy by byly v příkladech, i když by ty příklady reálné světa byly mnohem lepší

$$BLUE = BP \cdot \sqrt[4]{t_1 \cdot t_2 \cdot t_3 \cdot t_4} \in [0, 1]$$

↳ G. průměr n-gramové přesnosti pro $n=1, 2, 3, 4$

⊕ - rychlý výpočet, relativně objektivní míra

- obecně přijímaný standard

⊖ - závislost datového množství referenčních příkladů je draké

- favorizuje statistické systémy, co se více m' ladí

- špatně zachycuje různé varianty - slovosledné, morfologické

- přeceňuje se, není až tak univerzální