

Korpusová lingvistika

Studie jazyka:

— struktury

— užití

Co to je korpus?

Sbírka textů, v ideálním případě kompletní

- je zaučovací nástroj

- Výzdy korpus když původně zdrojení.

Něco se mohou, když v nich mohou.

Vlastnosti:

- reprezentativnost

→ Je důležité v jistý moment zastavit. Protože

- konkrétní velikost

je potřeba pochumat významy v jiných časech.

- sestřípný cílové formy

- standardní reference

Whitelové korpusy

Brown korpus

- první velký korpus

- shazil se nejdříve reprezentant Americkou angličtinu v roce 1961.

Penn Treebank

- první syntakticky anotovaný korpus

- článek z Wall Street Journal → Je to ale burzovní sleny, nemá to spisovnou angličtinu.

Český národní korpus

- vzniklo v 90. letech : UK, MU, AUL

- zahrnuje mnoho morfológické úrovně

- mix novin, knih a technické literatury

↳ 60% ↳ 25% ↳ 15%

Pražský sánskostní korpus (PST) (100 000 ref, 1 200 000 slov)

- syntakticky anotovaný korpus → celý anotovaný ručně (na rozdíl od ostatních ve sítě)

- založeno na teorii funkčního generativního psaní

- úrovňi anotace:

- morfologie

- analytické rovinu (pozorování syntaktické)

- tektotypologické rovinu: závislosti, jidlo, horeference, vše osobní

Prvky Ambie Dependency Treebank

- prvního závislostního korpus použitý na typové úplné odlišení jazyků.

Pravděpodobnostní a statistické metody

- pravděpodobnost je už jen aproximace trénovacích dat

↳ tomu může být velký vliv na výsledky

- empirické trénování příklad na posluchání slovníku / korpusu.

$$\text{Bayesův vzorec: } P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

Cit: občerme přediktovat $p(w|h)$, kde w je hledané slovo a h je historie textu.

- tyto modely by byly ale obtížná většímu přednosti.

Používají se trigramové modely:

jako „robenko“ historie barevné parze kontext délky 3.

- správná konstrukce bude mít 0 pravděpodobnost

Mutaci výhľadoviny:

- nulové pravděpodobnosti se nastavují na velmi malé ϵ .

- některé neexistující kombinace

- překlad bude vnitřně nemít dobrý trénovací dataset, jelikož se preferuje zákonitost významových form 1:1.

- Od 2015 se překládají statisticky, nyní funguje věkověrný překlad.

Mobná základního klasifikace: občerme $f: F \rightarrow A$, hledáme $P(A|F)$

Bayesova miřítko: $P(A|F) = P(F|A) \cdot P(A) / P(F)$ ↳ tady musíme dělit jeho $= 1$ když jsme to dostali na vstupu

problém zhlášť negativní

postřídkové sémantické pochopení

- těch hypotéz nesmí být příliš mnoho

modely: překladový $P(f|a)$

cílového jazyka $P(a)$

↳ obecně směr překladu a hledání věty, ze které to vzniklo.

- neexistuje vztah k originálni, kde se mohou vyskytnout různé

- překlad probíhá obecně

$$\text{BLEU: } \text{BP} * (\rho_1 \rho_2 \rho_3 \rho_n)^{\frac{1}{n}} \in \langle 0, 1 \rangle$$

↳ ře reprezentant jeho %, ale nemá to využití pro parametrizaci oddělovaných systémů

normalizace za středníst

1/1

což je poměrní délka,
aby zvítězila kompletní věta

↳ geometrický průměr n-gramové přesnosti pro $n=1..h$

$$D_i := \# \text{mezinářích i-gramů} / \text{i-gramů celkově}$$

- je to velmi vychlé

- velmi přesné, protože existuje více ref. překladů

- používá synonyma by kouzlo shora, když bych měl jen jeden ref. překlad

- zadání souboru frekvenčního datasetu zavedne BLEU parci o 0,5%

- nemáš jen tak shloubit celý internet, jelikož správné shloubě

je vytvářeno již výplním shloubujícím překladem, když se to nedá použít.