



statistické zpracování přirozeného jazyka

OJ205

- užití: vyhledávání kolokací
- **problémy:**
 - - rozložení přirozeného jazyka
 - - stoplist
 - - nastavení vyhledávání
 - - interpretace výsledků

statistické zpracování přirozeného jazyka

- Jako n-gram označujeme obvykle posloupnost slov délky n , která se vyskytuje v korpusu.
- Známe posloupnost s_1, \dots, s_{n-1} slov v korpusu. Jaké slovo za nimi bude následovat?
- Formálně se n-gramový jazykový model skládá z n podmnožin V_1, \dots, V_n ch pravděpodobností:

$$(p_i | s_1, \dots, s_{n-1})$$

čili pravděpodobností, že se vyskytlo slovo s_n za předpokladu, že před ním se vyskytla slova s_1, \dots, s_{n-1} . Souhrn všech p_i chto pravděpodobností pro všechny možné kombinace slov v korpusu se p_i vá n-gramový jazykový model.

n-gramové jazykové modely

- $1, \dots, n$ nemusí být pouze slova – můžeme vytvořit n -gramový model znaků, fonémů, pádů, rodů, číselných značek apod., případně i komplikovanější modely, kde např. 0 bude morfologická značka a $1, \dots, n-1$ budou slova.
- Využití - n -gramových modelů ve zpracování jazyka: existuje jich velké množství, z nichž nejpoužívanější je pravděpodobně trigramový ($n = 3$) jazykový model.

n-gramové jazykové modely

- $(\neg | \neg) = 3/4$
- $(\neg | \neg) = 1/4$
- $(\neg | \neg) = 2/3$
- $(\neg | \neg) = 1/3$
- $(\neg | \neg) = 1/2$
- $(\neg | \neg) = 1/2$
- $(\neg | \neg) = 1$
- $(\neg | \neg) = 1$
- $(\$ | \neg) = 1$
- $(\neg | \neg) = 0$ *š h* .

rozbor věty: Máma mele maso.
