



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 $n - 1$ slov v korpusu. Jaké slovo za nimi bude následovat?
- Formálně se n-gramový jazykový model skládá z n podmnožin pravděpodobností:

$$P(w_i | w_{i-1}, \dots, w_{i-n})$$

čili pravděpodobností, že se vyskytlo slovo w_i za předpokladu, že před ním se vyskytla slova w_{i-1}, \dots, w_{i-n} . Souhrn všech n podmnožin pravděpodobností pro všechny možné kombinace slov v korpusu se nazývá n-gramový jazykový model.

n-gramové jazykové modely

- $1, \dots, n$ nemusí být pouze slova – můžeme vytvořit n -gramový model znaků, fonémů, pádů, tvarů, slov, značek apod., případně i komplikovanější modely, kde např. n 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

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

rozbor věty: Máma mele maso.
