

SPRÅKMODELLER

En modell som beregner sannsynligheten for en sekvens av ord

N-GRAMMODELLER

Brukes for sannsynlighetsberegning

Minikorpus:

<s> Lise elsker å danse <\s>

<s> Ola hater å sykle <\s>

<s> Peter liker å danse <\s>

Unigrammer (n=1): <s>, <Lise>, <elsker>, <å>, <danse>, <\s>

Bigrammer (n=2): <<s>, Lise>, <Lise, elsker>, <elsker, å>, <å, danse>, <danse, <\s> >

BEREGNE SANNSYNLIGHET FOR EN SETNING

Betinget sannsynlighet- sannsynligheten for at noe skjer, *gitt* noe annet:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \longrightarrow P(\text{første ord} | \langle s \rangle) = \frac{P(\text{første ord} | \langle s \rangle)}{P(\langle s \rangle)}$$

For å beregne sannsynligheten for *resten* av setningen, ganger vi sannsynligheten for hvert bigram sammen:

$$\frac{P(\text{første ord} | \langle s \rangle)}{P(\langle s \rangle)} * \frac{P(\text{ neste ord} | \text{forrige ord})}{P(\text{forrige ord})} * \frac{P(\langle \backslash s \rangle | \text{siste ord})}{P(\text{siste ord})}$$

$$P(< s > \text{ Ola hater å danse } < \backslash s >)$$

Minikorpus:

<s> Lise elsker å danse <\s>

<s> Ola hater å sykle <\s>

<s> Peter liker å danse <\s>

$$P(\text{Ola} | < s >) * P(\text{hater} | \text{Ola}) * P(\text{å} | \text{hater}) * P(\text{danse} | \text{å}) * P(< \backslash s > | \text{danse})$$

$$= \frac{\text{count}(< s >, \text{Ola})}{\text{count}(< s >)} * \frac{\text{count}(\text{Ola}, \text{hater})}{\text{count}(\text{Ola})} * \frac{\text{count}(\text{hater}, \text{å})}{\text{count}(\text{hater})} * \frac{\text{count}(\text{å}, \text{danse})}{\text{count}(\text{å})} * \frac{\text{count}(\text{danse}, < \backslash s >)}{\text{count}(\text{danse})}$$

$$= \frac{1}{3} * \frac{1}{1} * \frac{1}{1} * \frac{2}{3} * \frac{2}{2}$$

$$= 0.33 * 1 * 1 * 0.66 * 1 = 0.217$$

= 21.7% sjanse for setningen

NLTK

- Et genialt verktøy for språktekere!
- Inneholder biblioteker for språkprossesering
 - .. Blant annet en n -gram-metode for oblig 2a
- Installering: <https://www.nltk.org/install.html>
- Litt rustne matteegenskaper? 🤔 → [forkurs](#)