

Syntaks og klassifisering

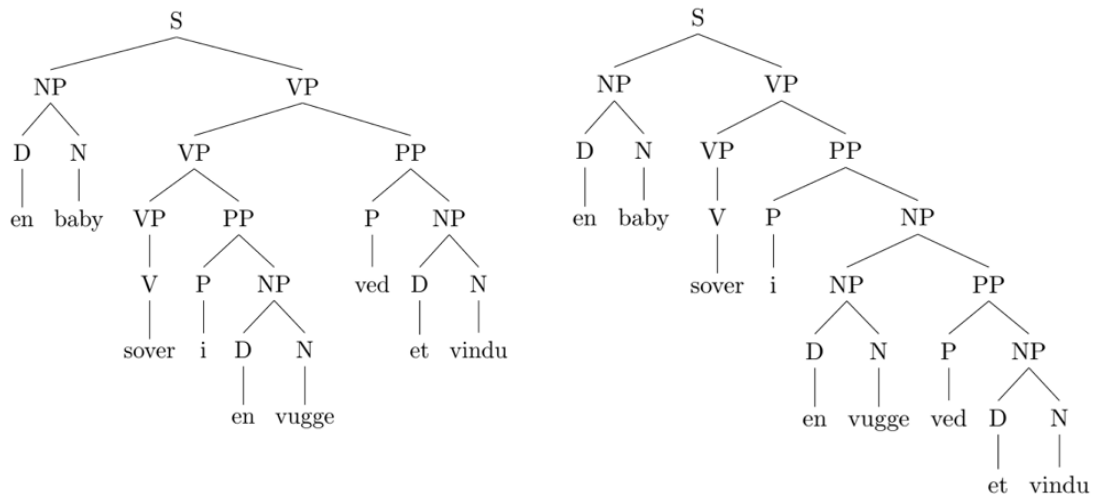
Løsningsforslag

1 Strukturell flertydighet (eksamen V2017)

1.

- a) Grammatikken tildeler 2 analyser til setningen “en baby sover i en vugge ved et vindu” (se trærne nedenfor).
- b) Grammatikken tildeler 1 analyse til setningen “et barn i en vugge sover”.

2.



3. For at grammatikken skal tillate

vilkårlig mange adjektiver må vi utvide den med følgende regler:

- NP → D N'
 N' → N
 N' → Adj N'

(Taggene D og Det er hyppig brukt om hverandre. De er begge forkortelser for “determinativ”.)

4. Den utvidede grammatikken tildeler 1 analyse til setningen “en rund liten søt baby sover”.

2 En grammatikk for norsk

I separat fil.

3 Naive Bayes-klassifisering

3.1 Sannsynlighet for en setning i en kategori

Setningen er "Jeg liker alltid utenlandske filmer"

Sannsynligheten for POS

$$P(POS) = 1/2 = 0.5$$

$$\begin{aligned} P(\text{setning}|POS) &= P(POS) * P(\text{jeg}|POS) * P(\text{liker}|POS) * P(\text{alltid}|POS) * P(\text{utenlandske}|POS) * P(\text{filmer}|POS) \\ &= 0.5 * 0.09 * 0.07 * 0.29 * 0.04 * 0.08 \end{aligned}$$

$$P(\text{setning}|POS) = 0.0000029232$$

Sannsynligheten for NEG

$$P(NEG) = 1/2 = 0.5$$

$$\begin{aligned} P(\text{setning}|NEG) &= P(NEG) * P(\text{jeg}|NEG) * P(\text{liker}|NEG) * P(\text{alltid}|NEG) * P(\text{utenlandske}|NEG) * P(\text{filmer}|NEG) \\ &= 0.5 * 0.16 * 0.06 * 0.06 * 0.15 * 0.11 \end{aligned}$$

$$P(\text{setning}|NEG) = 0.000004752$$

$P(\text{setning}|POS) < P(\text{setning}|NEG)$, så derfor vil Naive Bayes klassifisere setningen "Jeg liker alltid utenlandske filmer" som **negativ**.

3.2 Sannsynlighet for at et dokument er i en kategori

Dokument å klassifisere: *rask, par, skyte, fly*

Ordforråd $|V| = 7$ (gitt av antall unike ord, altså typer, i alle taggedede dokumenter)

Antall tokens i *komedie*: 9

Antall tokens i *action*: 11

Komedie

Sannsynligheter for ord gitt *komedie*:

$$P(\text{rask} \mid \text{komedie}) = (1 + 1) / (9 + 7) = 0,125$$

$$P(\text{par} \mid \text{komedie}) = (2 + 1) / (9 + 7) = 0,1875$$

$$P(\text{skytte} \mid \text{komedie}) = (0 + 1) / (9 + 7) = 0,0625$$

$$P(\text{fly} \mid \text{komedie}) = (1 + 1) / (9 + 7) = 0,125$$

Prior-sannsynligheten for *komedie*:

$$P(\text{komedie}) = 2/5$$

$$(2/5) * 0,125 * 0,1875 * 0,0625 * 0,125 = 0.000073$$

Action

Sannsynligheter for ord gitt *action*:

$$P(\text{rask} \mid \text{action}) = (2 + 1) / (11 + 7) = 0,1667$$

$$P(\text{par} \mid \text{action}) = (0 + 1) / (11 + 7) = 0,0556$$

$$P(\text{skytte} \mid \text{action}) = (4 + 1) / (11 + 7) = 0,2778$$

$$P(\text{fly} \mid \text{action}) = (1 + 1) / (11 + 7) = 0,1111$$

Prior-sannsynligheten for *action*:

$$P(\text{action}) = 3/5$$

$$(3/5) * 0,1667 * 0,0556 * 0,2778 * 0,1111 = 0.00017$$

Dokumentet klassifiseres som **action** fordi $0.000073 < 0.00017$.