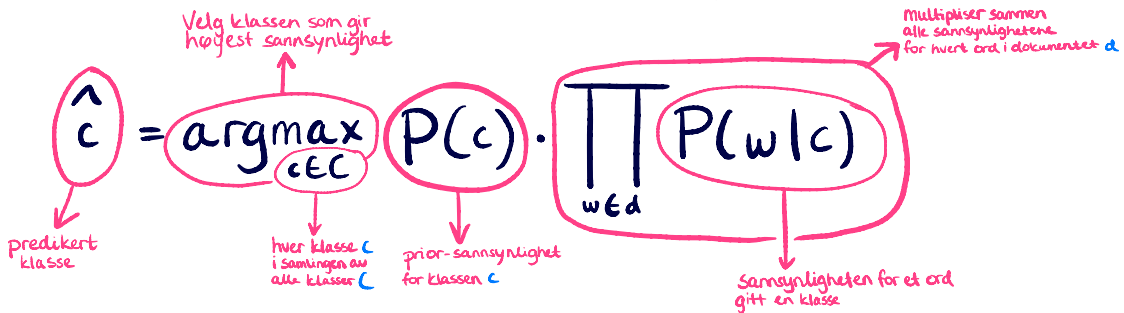


Naive Bayes



$$P(c) = \frac{N_c}{N_{\text{dok}}} \rightarrow \begin{array}{l} \text{antall dokumenter med gitt klasse} \\ \text{totalt antall dokumenter i treningsdatasetten} \end{array}$$

$$P(w_i|c) = \frac{\text{count}(w_i, c)}{(\sum_{w \in V} \text{count}(w, c))} \rightarrow \begin{array}{l} \text{antall ganger ordet forekommer i klassen } c \\ \text{antall tokens i klassen} \end{array}$$

Hovedantagelse

Formelen over er basert på hovedantagelsen:

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c|d)$$

— det vil si —→

Vi kan predikere korrekt klasse ved å velge den som har høyest sannsynlighet for et gitt dokument

Utleddning

• fra hovedantagelsen til Naive Bayes

1) Hovedantagelsen

$$\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c|d)$$

2) Bayes teorem

– betinget sannsynlighet + produktsetningen

$$P(c|d) = \frac{P(d|c)P(c)}{P(d)}$$

3) Fjern $P(d)$

Siden $P(d)$ er uavhengig av k og vil være konstant for hver d , kan den fjernes

$$P(c|d) = \frac{P(d|c)P(c)}{\cancel{P(c)}} = P(d|c)P(c)$$

4) (Naiv) uavhengighetsantagelse

Antar (naivt) at alle trekk er uavhengige av andre trekk, gitt klassen.

→ da kan vi bruke multiplikasjonsregelen for uavhengige hendelser.

$$P(d|c) = P(w_1, w_2, \dots, w_k | c) \approx \prod_{i=1}^k P(w_i | c)$$

erstatt

$$P(c|d) = P(d|c)P(c)$$

$$P(c|d) = P(w_1, w_2, \dots, w_k | c)P(c) \approx \prod_{i=1}^k P(w_i | c)P(c)$$

5) Fullfør

Setter sammen til ferdig formel

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c) \prod_{w \in d} P(w | c)$$