

# Naive Bayes Classifier

## 1. Wet van Bayes

Van alle patiënten in de hematologie afdeling is 1% geïnfecteerd met het Hiv-virus, 97% daarvan testen positief. Dit is 4% van alle niet geïnfecteerde patiënten. Bereken de kans dat een patiënt die positief wordt getest effectief het Hiv-virus heeft.

$C_1$  = "HIV+ patiënten"

$C_2$  = "HIV- patiënten "

$X$  = "test resultaat"

$$P(C_1) = 0.01 \quad P(X = + \mid C_1) = 0.97$$

$$P(C_2) = 0.99 \quad P(X = + \mid C_2) = 0.04$$

$$P(C_1 \mid X = +) =$$

$$\frac{P(C_1) \cdot P(X = + \mid C_1)}{P(C_1) \cdot P(X = + \mid C_1) + P(C_2) \cdot P(X = + \mid C_2)}$$

## Terminologie

- $P(C_1)$ ,  $P(C_2)$ : eerdere waarschijnlijkheid
- $P(C_1 | X = +)$ ,  $P(C_2 | X = +)$ : latere waarschijnlijkheid

In het algemeen

- n klassen:  $C_1, C_2, \dots, C_n$
- k voorspellende variabelen  $X_1, X_2, \dots, X_k$

Vraag

Als  $X_1 = x_1, X_2 = x_2, \dots, X_k = x_k$ , wat is dan de meest waarschijnlijke klasse?

$$\begin{aligned} & P(C_i | X_1 = x_1, X_2 = x_2, \dots, X_k = x_k) ? \\ = & P(C_i | x_1, x_2, \dots, x_k) ? \end{aligned}$$

Wet van Bayes:

Wet van Bayes:

$$\begin{aligned} & P(C_i | x_1, x_2, \dots, x_k) \\ = & \frac{P(C_i) \cdot P(x_1, x_2, \dots, x_k | C_i)}{P(x_1, x_2, \dots, x_k)} \end{aligned}$$

## Voorbeeld

- Titanic dataset
- 4 variabelen, 2201 observaties
- X1: klasse  
( '0' = personeel, '1' = duurste klasse, '2' = middle klasse, '3' = goedkoopste klasse)
- X2: leeftijd  
( '1' = volwassen, '0' = kind)
- X3: geslacht  
( '1' = man, '0' = vrouw)
- C: gered of niet  
( '1' = gered, '0' = niet gered)

Voorbeeld vraag: Zal een meisje in de laagste klasse gered worden of niet?

$P(\text{gered} \mid \text{goedkoopste, kind, vrouw}) = ?$

$P(\text{niet gered} \mid \text{goedkoopste, kind, vrouw}) = ?$

$$\begin{aligned}
 &P(\text{gered} \mid \text{goedkoopste, kind, vrouw}) \\
 &= \frac{P(\text{gered}) \cdot P(\text{goedkoopste, kind, vrouw} \mid \text{gered})}{P(\text{goedkoopste, kind, vrouw})}
 \end{aligned}$$

$$\begin{aligned}
 &P(\text{niet gered} \mid \text{goedkoopste, kind, vrouw}) \\
 &= \frac{P(\text{niet gered}) \cdot P(\text{goedkoopste, kind, vrouw} \mid \text{niet gered})}{P(\text{goedkoopste, kind, vrouw})}
 \end{aligned}$$

Opmerkingen:

- noemer  $P(\text{goedkoopste, kind, vrouw})$  geen reden om te vergelijken.
- Om alle vragen te beantwoorden, moeten de volgende dingen gekend zijn

$P(\text{gered})$  en  $P(\text{niet gered})$   
 $P(x_1, x_2, x_3 \mid \text{gered})$  en  
 $P(x_1, x_2, x_3 \mid \text{niet gered})$   
 Voor alle combinaties van  $x_1, x_2, x_3$   
 $\rightarrow 2 \cdot 4 \cdot 2 \cdot 2 = 32$  mogelijkheden  
 Geschatte waarschijnlijkheden via  
 frequenties in training set

Probleem:

- Grote hoeveelheid van mogelijkheden  $x_1, x_2, x_3$
- Bepaalde combinaties kunnen niet voorkomen in een klasse  $C_i$  in de training set
  - schatting:  $P(x_1, x_2, x_3 | C_i) = 0$ .
  - $P(C_i | x_1, x_2, x_3) = 0$

## 1. Oplossing: Naive Bayes

- Assumptie:  $X_1, X_2, X_3$  niet onderling afhankelijk in elke klasse  $C_i$ 
  - $P(x_1; x_2; x_3 | C_i)$   
 $= P(x_1 | C_i) \cdot P(x_2 | C_i) \cdot P(x_3 | C_i)$
  - $P(x_1 | C_i), P(x_2 | C_i), P(x_3 | C_i)$  Schatting via frequenties in training set
  - Minder snel nul; minder kansen:  $2 \cdot (4+2+2) = 16$