

# INCERTIDUMBRE

## PROBABILIDADES Y PLAUSIBILIDAD

$P(A|I)$  := probabilidad de una aserción  $A$ , dada la información

$$P(\text{Verdad}|I) = 1 \quad P(\text{Falso}|I) = 0$$

$$0 \leq P(A|I) \leq 1$$

Regla de la suma:  $P(A|I) + P(\bar{A}|I) = 1$

Regla del producto:  $P(A, B|I) = P(B|A, I) \cdot P(A|I)$

Teorema de Bayes:  $P(B|A, I) = \frac{P(A|B, I) \cdot P(B|I)}{P(A|I)}$

## PROBABILIDAD CONDICIONAL

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

$$\left. \begin{array}{l} P(A, B) = P(A|B) \cdot P(B) \\ P(B, A) = P(B|A) \cdot P(A) \end{array} \right\} \Rightarrow P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

TEOREMA DE BAYES

$$P(\text{clase}|\text{datos}) = \frac{P(\text{datos}|\text{clase}) \cdot P(\text{clase})}{P(\text{datos})}$$

↑ posterior      ↓ verosimilitud de hipótesis      ← prior      ← evidencia

## Terminología

clase = hipótesis = etiqueta = variables → dependientes  
→ a predecir

datos = atributos = características = variables → independientes  
→ predictoras

- CLASIFICADOR MÁXIMA VEROSIMILITUD (ML): selecciona la hipótesis que maximiza la verosimilitud de la hipótesis dados los datos.

$$H^* = \arg \max_H P(D|H)$$

ML no usa información a priori (equivalente a asumir un priori constante)

- CLASIFICADOR MÁXIMA PROBABILIDAD A POSTERIORI (MAP): selecciona la hipótesis que maximiza la probabilidad a posteriori.

$$H^* = \arg \max_H P(H|D) = \arg \max_H \frac{P(D|H) \cdot P(H)}{\underbrace{P(D)}} = \arg \max_H P(D|H) \cdot P(H)$$

Este es el llamado CLASIFICADOR DE BAYES

↑ igual  $\forall H$  hipótesis

- INFERENCIA BAYESIANA: promediar sobre todas las hipótesis con probabilidades.

$$P(H|D) = \frac{P(D|H) \cdot P(H)}{P(D)}$$

$$E[f(H)] = \int dH f(H) \cdot P(H|D)$$

- CLASIFICADOR A PRIORI DE LA CLASE: selecciona la hipótesis que maximiza la probabilidad de la clase

$$H^* = \arg \max_H P(H)$$

⊛ Prioris uniformes  $\Rightarrow$  Clasificador ML = Clasificador MAP (Bayes)

⊛ Bayes es óptimo (minimiza el error)

## REDES BAYESIANAS ("Modelos gráficos")

Regla de la cadena:

$$\begin{aligned} P(x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}) &= P(x^{(1)} | x^{(2)}, x^{(3)}, x^{(4)}) \cdot P(x^{(2)}, x^{(3)}, x^{(4)}) = \\ &= P(x^{(1)} | x^{(2)}, x^{(3)}, x^{(4)}) \cdot P(x^{(2)} | x^{(3)}, x^{(4)}) \cdot P(x^{(3)}, x^{(4)}) = \\ &= P(x^{(1)} | x^{(2)}, x^{(3)}, x^{(4)}) \cdot P(x^{(2)} | x^{(3)}, x^{(4)}) \cdot P(x^{(3)} | x^{(4)}) \cdot P(x^{(4)}) \end{aligned}$$

Representación gráfica de la regla de la cadena:

Nodos: variables

Aristas dirigidas: dependencias

Una flecha hacia un nodo dado corresponde a variables dependientes de la variable que origina la flecha.

Ej: la representación gráfica no es única

Interpretación grafo:

Nodo  $i$ : variable  $x^{(i)}$

$\Pi(x^{(i)})$ : padres del nodo  $x^{(i)}$

$$P(x^{(1)}, x^{(2)}, \dots, x^{(n)}) = \prod_{i=1}^n P(x^{(i)} | \Pi(x^{(i)}))$$

Nota: Si  $\Pi(x^{(i)}) = \emptyset \Rightarrow P(x^{(i)} | \Pi(x^{(i)})) = P(x^{(i)})$

## INDEPENDENCIA CONDICIONAL

A es condicionalmente independiente de B dado C si:

$$P(A|B,C) = P(A|C)$$

Que es equivalente a:

$$P(A,B|C) = P(A|C) \cdot P(B|C)$$

## CLASIFICADOR DE NAÏVE-BAYES

El clasificador de Naïve-Bayes asume que todos los atributos son condicionalmente independientes dada la clase.

Atributos:  $x = \{x^{(1)}, x^{(2)}, \dots, x^{(D)}\}$

Clase:  $c \in \{1, 2, \dots, G\}$

$$P(c|x) = \frac{P(x|c) \cdot P(c)}{P(x)} = \frac{P(x^{(1)}, x^{(2)}, \dots, x^{(D)}|c) \cdot P(c)}{P(x)} \approx$$

$$\approx \frac{P(x^{(1)}|c) \cdot P(x^{(2)}|c) \cdot \dots \cdot P(x^{(D)}|c) \cdot P(c)}{\sum_{c'=1}^G P(x^{(1)}|c') \cdot P(x^{(2)}|c') \cdot \dots \cdot P(x^{(D)}|c')}$$

## ESTIMADOR DE LAPLACE

- Estimación de probabilidades con frecuencias:

$$P_i = \frac{n_i}{n_{\text{total}}} \quad i=1, 2, \dots, K$$

- Estimador de Laplace: añadir ejemplos ficticios

$$P_i = \frac{n_i + \mu/K}{n_{\text{total}} + \mu} \quad i=1, 2, \dots, K$$

Ventajas:

- evita estimaciones nulas para las probabilidades
- estimaciones más robustas
- asintóticamente pequeño

# INCERTIDUMBRE

## PROBABILIDAD (ENFOQUE BAYESIANO)

$$X_n = \frac{1}{\sqrt{N}} \sum_{i=1}^N \xi_{in}$$

$n=1, \dots, N$

$\xi_{in}$  v.A. independientes e idénticamente distribuidas  
con media cero.

Varianza finita:  $\sigma^2 = E[\xi_{in}^2] < \infty$

TCL:  $X_n \underset{N \rightarrow \infty}{\sim} N(0, \sigma^2)$

## PROBABILIDAD CONDICIONAL

$$P(A|B) = \frac{P(A, B)}{P(B)}$$

$$\left. \begin{aligned} P(A, B) &= P(A|B) \cdot P(B) \\ P(B, A) &= P(B|A) \cdot P(A) \end{aligned} \right\} \Rightarrow$$

TEOREMA DE BAYES

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(\text{clase} | \text{datos}) = \frac{P(\text{datos} | \text{clase}) \cdot P(\text{clase})}{P(\text{datos})}$$

↑ posterior  
↓ verosimilitud de hipótesis  
prior  
↓ evidencia

## Terminología:

clase = hipótesis = etiqueta = variables → dependientes  
→ a predecir

datos = atributos = características = variables → independientes  
→ predictoras

## Ejemplo:

Hipótesis (clase) → llueve  
 → no llueve

Evidencia =  $P(\text{"paraguas"})$

Datos = "alguien llega con paraguas"

PRIOR  $\begin{cases} P(\text{"llueve"}) = 20\% = 0.2 \\ P(\text{"no llueve"}) = 80\% = 0.8 (= 1 - 0.2) \end{cases}$

VEROSIM:  $IP(\text{"paraguas"}|\text{"no llueve"}) = 0.1$

$$IP(\text{"paraguas"}|\text{"llueve"}) = \frac{IP(\text{"llueve"}|\text{"paraguas"}) \cdot P(\text{"paraguas"})}{IP(\text{"llueve"})}$$

sin datos suficientes

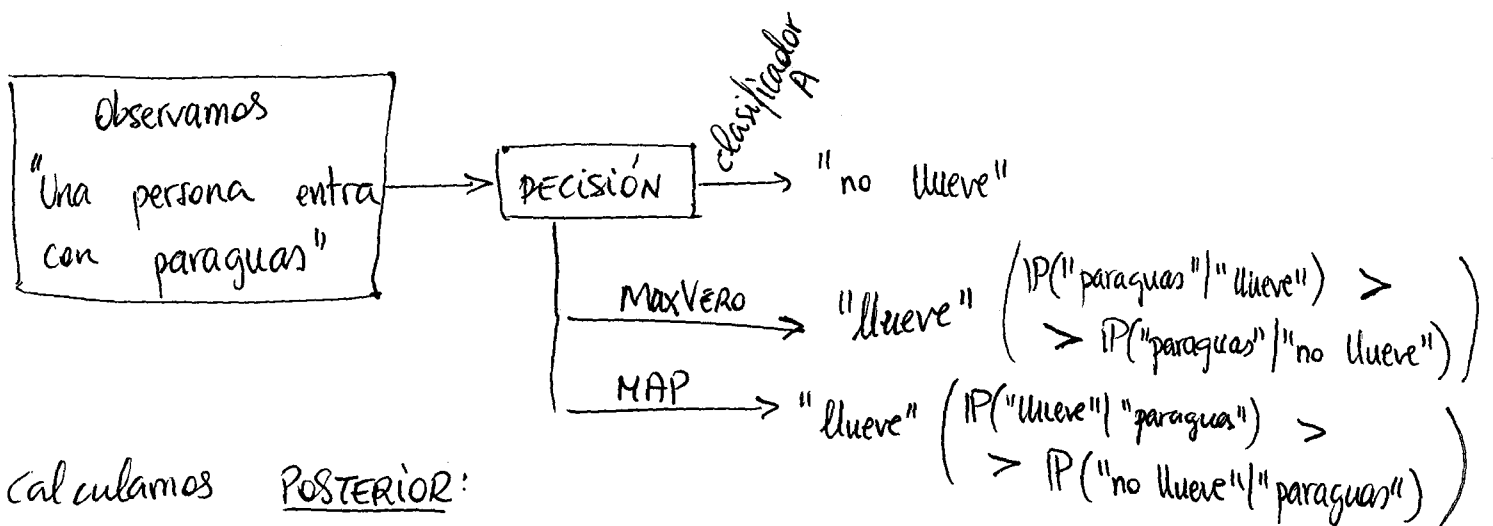
↓ nos lo dan

$$IP(\text{"paraguas"}|\text{"llueve"}) = 0.7$$

CLASIFICADOR A: basado en PRIORES (máximo entre priores)

CLASIFICADOR B: máxima verosimilitud ( $ML \equiv \max \text{likelihood}$ )

CLASIFICADOR C: MAP (máximo entre posteriores)



Calculamos POSTERIOR:

$$P(\text{"llueve"}|\text{"paraguas"}) = \frac{IP(\text{"paraguas"}|\text{"llueve"}) \cdot P(\text{"llueve"})}{IP(\text{"paraguas"})} = \frac{0.7 \cdot 0.2}{IP(\text{"paraguas"})} = 0.64$$

$$P(\text{"no llueve"}|\text{"paraguas"}) = \frac{IP(\text{"paraguas"}|\text{"no llueve"}) \cdot IP(\text{"no llueve"})}{IP(\text{"paraguas"})} = \frac{0.1 \cdot 0.8}{IP(\text{"paraguas"})} = 0.3$$

Observación:

$$IP(\text{"paraguas"}) = IP(\text{"paraguas"}|\text{"llueve"}) \cdot P(\text{"llueve"}) + IP(\text{"paraguas"}|\text{"no llueve"}) \cdot IP(\text{"no llueve"})$$

$$= 0.22$$

OSO! → los costes asimétricos varían el umbral

## ENFOQUE BAYESIANO

$H \in \{H_1, H_2, \dots, H_d\}$  con una probabilidad  $\{P(H_i|D), i=1, \dots, d\}$

$$E_H[f(H)] = \sum_{i=1}^d f(H_i) \cdot P(H_i|D)$$

⊛ Estamos suponiendo costes simétricos en todos estos clasificadores.

NAÏVE - BAYES  $\rightarrow$  probabilidad en espacio  $D$ -dimensional  $\rightarrow D$  prob. en 1 dimensión

$$P(C|\tilde{x}) = \frac{P(\tilde{x}|C) \cdot P(C)}{P(\tilde{x})} \approx \frac{P(x_1|C) \cdot P(x_2|C) \cdot \dots \cdot P(x_D|C) \cdot P(C)}{\text{Normalización}}$$

↑  
VECTOR  
EATORIO

$$\tilde{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_D \end{pmatrix}$$

exacto si  
las estimaciones  
de prob. son correctas

↑  
en la misma aprox.

$$P(\tilde{x}|C) \approx P(x_1|C) \cdot P(x_2|C) \cdot \dots \cdot P(x_D|C)$$

↑ "los atributos son INDEPENDIENTEMENTE CONDICIONADOS a C"

OJO:

$$P(A) = P(A, B) + P(A, \neg B)$$

$$P(G|E) =$$

$$\cancel{P(\neg G | \neg F, \neg T)}$$

$$P(A|B) + P(A|\neg B) \neq 1$$

$$P(\neg G, \neg T | \neg F)$$

$$P(\neg G | \neg T, \neg F)$$

$$P(\neg G, \neg T | \neg F) = P(\neg G | \neg T, \neg F) \cdot \cancel{P(\neg T)}$$

$$P(\bar{G}, \bar{F}, \bar{T}) + P(\bar{G}, F, T) + P(\bar{G}, \bar{F}, T) + P(\bar{G}, F, \bar{T}) = P(\bar{G})$$

$$\downarrow$$
  

$$0'1$$

$$\downarrow$$
  

$$0'05$$

$$\downarrow$$
  

$$0'0$$

$$\downarrow$$
  

$$0'6$$
  

$$N$$

$$P(\bar{G}, E) + P(G, E)$$

$$\Rightarrow P(\bar{G}, \bar{F}, \bar{T}) = 0'6 - 0'1 - 0'05 = 0'45$$

Lo mismo para el otro