



Bayes Learning

↘ Aprendizaje basado en la regla de Bayes

↘ Bases probabilísticas a metodologías comunes en machine learning

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

Sea D un conjunto de datos observados

Queremos encontrar el **mejor** modelo/hipótesis h para los datos



El modelo **más probable** para el fenómeno observado, dados:

- los datos observados
- algún conocimiento previo sobre las probabilidades de las diferentes hipótesis en el espacio H de hipótesis

Ejemplo: Si las etiquetas reflejan cercanía, es más probable un modelo que asigne etiquetas cercanas a datos cercanos!

Denotamos:

Prior probability $P(h)$: Probabilidad inicial de que un modelo ocurra, previa a la observación de los datos

Conocimiento previo
sobre el fenómeno

! Si no tenemos conocimiento previo, asumimos $P(h)$ igual para todas las hipótesis !

$P(D)$: Probabilidad de que los datos D sean observados (sin saber nada sobre la hipótesis)

$p(D|h)$: Probabilidad de observar D en un mundo en donde h se cumple

$P(h|D)$: Probabilidad de que h se cumpla dados los datos observados D .

→ "Posterior Probability"

Queremos encontrar el modelo más probable h dados los datos D

En otras palabras, queremos maximizar $P(h|D)$.

$$\operatorname{argmax}_{h \in H} P(h|D)$$

Contamos con la regla de Bayes

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Cualquier hipótesis que maximice esta expresión se denomina

"Máximo a Posteriori"
(MAP)

hipótesis máxima a posteriori (MAP)

$$h_{\text{MAP}} = \underset{h \in H}{\operatorname{argmax}} P(h|D)$$

$$= \underset{h \in H}{\operatorname{argmax}} \frac{P(D|h)P(h)}{P(D)}$$

$$= \underset{h \in H}{\operatorname{argmax}} P(D|h)P(h)$$

$$= \underset{h \in H}{\operatorname{argmax}} P(D|h) \longrightarrow$$

Maximum Likelihood
hypothesis
(ML)

hipótesis máxima a posteriori (MAP)

$$h_{\text{MAP}} = \arg \max_{h \in H} P(h|D)$$

$$= \arg \max_{h \in H} \frac{P(D|h)P(h)}{P(D)} \quad (\text{Bayes})$$

$$= \arg \max_{h \in H} P(D|h)P(h) \quad (P(D) \text{ no depende de } h)$$

$$= \arg \max_{h \in H} P(D|h) \quad \left(\begin{array}{l} \text{Asumiremos que } P(h) \text{ es} \\ \text{igual para toda } h, \text{ por lo} \\ \text{que no depende de } h \end{array} \right)$$

Utilicemos lo anterior para el caso de variable objetivo continua

bajo algunos supuestos:

X : Espacio de instancias

H : Espacio de hipótesis

Cada $h \in H$ es de la forma $h: X \rightarrow \mathbb{R}$

D : Conjunto con m datos $\{(x_1, d_1) \dots (x_m, d_m)\}$

y aquí un supuesto importante:

Las etiquetas d_i están afectadas por un ruido: Media 0

$$d_i = f(x_i) + e_i$$

en donde el error tiene distribución normal

$$e_i \sim \mathcal{N}(0, \sigma^2)$$

$$d_i \sim \mathcal{N}(f(x), \sigma^2)$$

Nota: La variable objetivo es continua, así que las probabilidades están dadas por una función de densidad.

En particular, recordemos la función de densidad de la distribución normal:

$$f_{N(\mu, \sigma^2)}(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma} \right)^2}$$

Ahora si calculemos la hipótesis más probable

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h) \quad \rightarrow \quad \text{"Maximum Likelihood"}$$

→ Asumiendo un conjunto fijo de instancias de entrenamiento $\{x_1, \dots, x_m\}$ podemos considerar el conjunto D como el conjunto de las etiquetas $D = \{d_1, \dots, d_m\}$ $d_i = f(x_i) + e_i$.

→ Asumiendo también **independencia** entre los datos observados tenemos

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i | h)$$

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i | h)$$

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2}$$

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(d_i - h(x_i))^2$$

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m -\frac{1}{2\sigma^2}(d_i - h(x_i))^2 = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m \frac{1}{2\sigma^2}(d_i - h(x_i))^2$$

$$\underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

||

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i | h)$$

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2}$$

densidad de la dist. normal

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

prop. de logaritmo

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(d_i - h(x_i))^2$$

el primer término no depende de h

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m -\frac{1}{2\sigma^2}(d_i - h(x_i))^2 = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m \underbrace{\frac{1}{2\sigma^2}}_{\text{constante}} (d_i - h(x_i))^2$$

$$\underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

||

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m p(d_i | h)$$

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - \mu)^2}$$

densidad de la dist. normal

$$= \underset{h \in H}{\operatorname{argmax}} \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(d_i - h(x_i))^2}$$

prop. de logaritmo

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2}(d_i - h(x_i))^2$$

el primer término no depende de h

$$= \underset{h \in H}{\operatorname{argmax}} \sum_{i=1}^m -\frac{1}{2\sigma^2}(d_i - h(x_i))^2 = \underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m \underbrace{\frac{1}{2\sigma^2}}_{\text{constante}} (d_i - h(x_i))^2$$

¡ Diferencia de cuadrados ! 

↑

$$\underset{h \in H}{\operatorname{argmin}} \sum_{i=1}^m (d_i - h(x_i))^2$$

||

Lo anterior quiere decir que [bajo los supuestos mencionados]
cualquier algoritmo que minimice el error cuadrático arrojará
una hipótesis de máxima verosimilitud (maximum likelihood)

Ej: regresión
red neuronal

Pero ojo: Los supuestos incluyen ruido (ϵ_i) en la variable objetivo
más no en las variables de entrada !!

Caso de predicción de probabilidades

Ahora consideremos el caso de una función $f: X \rightarrow \{0, 1\}$ para la cual queremos predecir una probabilidad.

Es decir queremos aprender la función

$$\hat{f}: X \rightarrow [0, 1] \quad \text{en donde} \quad \hat{f}(x) = P(f(x)=1)$$

Sea $D = \{ (x_1, d_1), \dots, (x_m, d_m) \}$ en donde $d_i = f(x_i) \in \{0, 1\}$
queremos entonces calcular

$$P(D|h) = \prod_{i=1}^m P(x_i, d_i | h)$$

asumiremos además que x_i es independiente de la hipótesis h . Así;

$$P(D|h) = \prod_{i=1}^m P(x_i, d_i | h) = \prod_{i=1}^m P(d_i | h, x_i) P(x_i)$$

$$P(D|h) = \prod_{i=1}^m P(x_i, d_i | h) = \prod_{i=1}^m \underbrace{P(d_i | h, x_i)}_{???} P(x_i)$$

Recordando que queremos modelar $h(x) = f(x) = P(f(x)=1)$
 es decir $P(d_i=1 | h, x_i) = h(x_i)$ así que

$$P(d_i | h, x_i) = \begin{cases} h(x_i) & \text{si } d_i = 1 \\ 1 - h(x_i) & \text{si } d_i = 0 \end{cases}$$

escrito de otra forma

$$P(d_i | h, x_i) = h(x_i)^{d_i} (1 - h(x_i))^{1-d_i}$$

Luego

$$P(D|h) = \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} P(x_i)$$

$$H_{ML} = \underset{h \in \mathcal{H}}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} p(x_i)$$

independiente de h

$$= \underset{h \in \mathcal{H}}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i}$$

$$= \underset{h \in \mathcal{H}}{\operatorname{argmax}} \sum_{i=1}^m d_i \ln h(x_i) + (1-d_i) \ln (1-h(x_i))$$

$$H_{ML} = \underset{h \in \mathcal{H}}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i} p(x_i)$$

independiente de h

$$= \underset{h \in \mathcal{H}}{\operatorname{argmax}} \prod_{i=1}^m h(x_i)^{d_i} (1 - h(x_i))^{1-d_i}$$

$$= \underset{h \in \mathcal{H}}{\operatorname{argmax}} \sum_{i=1}^m d_i \ln h(x_i) + (1-d_i) \ln (1-h(x_i))$$

ii Cross
entropy !!



Así es como encontramos que minimizar la función de cross entropy en el caso de predicción de probabilidades puede arrojar un modelo de máxima verosimilitud

Clasificador óptimo de Bayes (Bayes optimal classifier)

Hasta ahora pretendíamos responder la pregunta

¿Cuál es la hipótesis más probable dados los datos de entrenamiento? $P(h|D)$

Ahora, preguntémosnos

¿Cuál es la clasificación más probable de una nueva instancia dados los datos de entrenamiento?

¿Cuál es la clasificación más probable de una nueva instancia dados los datos de entrenamiento?

Ejemplo

| | $h_i(x)$ | $Pr(h/D)$ |
|-------|----------|-----------|
| h_1 | + | 0.4 |
| h_2 | - | 0.3 |
| h_3 | - | 0.3 |

h_1 maximum a priori hypothesis

Dada una instancia nueva x ,
¿Cuál es la etiqueta más probable para x ?

Considerando todas las hipótesis, es más probable la etiqueta **negativa** (0.6 vs 0.4)

En general,

$$P(v_j | D) = \sum_{h_i \in H} P(v_j | h_i) P(h_i | D)$$

Clasificación óptima de Bayes

$$\arg \max_{v_i \in V} \sum_{h_i \in H} P(v_j | h_i) P(h_i | D)$$

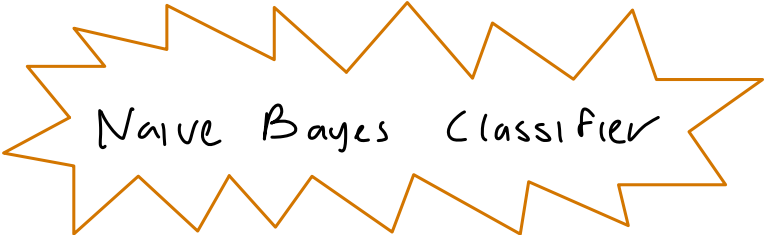
Este clasificador resulta óptimo: Dado el mismo espacio de hipótesis H y el mismo conocimiento previo (a saber, los valores de $P(h_i | D)$), ningún clasificador podría hacerlo mejor !

El inconveniente del "Bayes optimal classifier"....

El costo computacional!

Debe calcular probabilidades para cada $h \in H$ y para cada instancia nueva.

La alternativa?



Naive Bayes Classifier

Clasificador Naive Bayes

$$f: X \longrightarrow V \longrightarrow \text{Conjunto finito}$$

Llamaremos a_1, a_2, \dots, a_n a los atributos de un registro x_i

Dada una nueva instancia con atributos a_1, \dots, a_n , se pretende asignar la etiqueta mas probable "a priori" (MAP)

$$\begin{aligned} v_{\text{MAP}} &= \underset{v_j \in V}{\operatorname{argmax}} P(v_j | a_1, a_2, \dots, a_n) \\ &= \frac{P(a_1, a_2, \dots, a_n | v_j) P(v_j)}{P(a_1, a_2, \dots, a_n)} \\ v_{\text{MAP}} &= P(a_1, a_2, \dots, a_n | v_j) P(v_j) \end{aligned}$$

Bayes

denominador independiente

$$V_{MAP} = P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

$$V_{MAP} = P(a_1, a_2, \dots, a_n | v_j) P(v_j)$$

???

fácil contando frecuencias



asumimos una condición más:
que los a_i son condicionalmente
independientes, dada la etiqueta

$$\text{así, } P(a_1, a_2, \dots, a_n | v_j) = \prod_i P(a_i | v_j)$$

$$V_{NB} = \underset{v_j \in V}{\operatorname{argmax}} P(v_j) \prod_i P(a_i | v_j)$$

etiqueta según Naive Bayes

Ejercicio

Calcular la etiqueta según
Naive Bayes para un día
con características:

(sunny, cool, high, strong)

$$V_{NB} = \underset{V_j \in \{\text{yes}, \text{no}\}}{\operatorname{argmax}} \quad P(V_j) \prod_i P(a_i | V_j)$$

| Outlook | Temperature | Humidity | Wind | Play |
|----------|-------------|----------|--------|------|
| Sunny | Hot | High | Weak | No |
| Sunny | Hot | High | Strong | No |
| Overcast | Hot | High | Weak | Yes |
| Rain | Mild | High | Weak | Yes |
| Rain | Cool | Normal | Weak | Yes |
| Rain | Cool | Normal | Strong | No |
| Overcast | Cool | Normal | Strong | Yes |
| Sunny | Mild | High | Weak | No |
| Sunny | Cool | Normal | Weak | Yes |
| Rain | Mild | Normal | Weak | Yes |
| Sunny | Mild | Normal | Strong | Yes |
| Overcast | Mild | High | Strong | Yes |
| Overcast | Hot | Normal | Weak | Yes |
| Rain | Mild | High | Strong | No |