

EL CLASIFICADOR BAYESIANO ÓPTIMO

ALAN REYES-FIGUEROA APRENDIZAJE ESTADÍSTICO

(AULA 20) 22.ABRIL.2024

Función de costo

Dado un clasificador $\hat{y} = q : S \to \Omega$, definimos una **función de pérdida** o **costo** $L(y, \hat{y}(\mathbf{x}))$.

Ejemplos:

•
$$L(y, \widehat{y}(\mathbf{x})) = \begin{cases} 1, & \text{si } y \neq \widehat{y}(\mathbf{x}); \\ 0, & \text{si } y = \widehat{y}(\mathbf{x}). \end{cases}$$

•
$$L(y,\widehat{y}(\mathbf{x})) = (y - \widehat{y}(\mathbf{x}))^2$$
.

• $L(y, \widehat{y}(\mathbf{x})) = |y - \widehat{y}(\mathbf{x})|$.

•
$$L(y, \widehat{y}(\mathbf{x})) = \begin{cases} c_1, & \text{si } y = 1, \ \widehat{y}(\mathbf{x}) = 0; \\ c_2, & \text{si } y = 0, \ \widehat{y}(\mathbf{x}) = 1; \\ c_3, & \text{si } y = \widehat{y}(\mathbf{x}). \end{cases}$$
 • $L(y, \widehat{y}(\mathbf{x})) = |y - y(\mathbf{x})|.$

Definimos el **error de clasificación** como $\mathbb{E}(L(y, \hat{y}(\mathbf{x})))$. El **error empírico** se define como

$$\frac{1}{n}\sum_{i}L(y_{i},\widehat{y}(\mathbf{x}_{i})).$$

Típicamente, la función de costo satisface $L(y, \hat{y}(\mathbf{x})) \geq 0$.

Dado un conjunto de datos $(X,Y) \sim \mathbb{P}$, y una función de costo $L \geq 0$, queremos encontrar un clasificador $\widehat{y}(\mathbf{x})$ tal que

$$\mathbb{E}\big(L(Y,\widehat{Y}(X))\big) = \mathbb{E}_{X,Y}\big(L(Y,\widehat{Y}(X))\big) \text{ sea minima.} \tag{1}$$

En otros casos, nos puede interesar minimizar la probabilidad $\mathbb{P}(\sum_i L(y_i, \hat{y}(\mathbf{x}_i)) > threshold)$.

De (1)

$$\mathbb{E}_{X,Y}(L(Y,\widehat{Y}(X))) = \mathbb{E}_{X}\mathbb{E}_{Y|X=\mathbf{x}}(L(Y,\widehat{Y}(X)))$$

$$= \int_{\mathbb{R}^{d}} \mathbb{E}_{Y|X=\mathbf{x}} L(Y,\widehat{Y}(\mathbf{x})) f_{X}(\mathbf{x}) d\mathbf{x}$$
(2)

La ecuación en (2) es importante porque de alguna manera indica que el problema de minimización es desacoplado: se puede minimizar de forma separada en X y se puede también minimizar en Y.

Si minimizamos lo anterior sobre \widehat{Y} , es suficiente para cada ${\bf x}$ minimizar la siguiente función

$$\mathsf{argmin}_{\widehat{Y}(\boldsymbol{x})} \, \mathbb{E}_{Y|X=\boldsymbol{x}} \, L(Y,\widehat{Y}(\boldsymbol{x})), \ \, \forall \boldsymbol{x}.$$

Definición

La solución a la ecuación anterior se llama el clasificador bayesiano óptimo.

Obs! Este es un clasificador teórico. En la práctica no es posible calcularlo, a menos que dispongamos de la información de la distribución conjunta $\mathbb{P}(X,Y)$ teórica.

Se utiliza para dar estimados teóricos y cotas del error de clasificación.



Observaciones:

- Se llama *óptimo* porque es lo mejor que podemos hacer en el caso que tenemos la información completa $\mathbb{P}_{X,Y}$.
- En el caso finito, la integral $\mathbb{E}_{X,Y}(L(Y,\widehat{Y}(X)))$ se reduce a una suma.
- Se puede mostrar que en el caso finito, el clasificador bayesiano óptimo es la asignación \hat{y} que minimiza la "probabilidad de cometer un error"

$$n \, \mathbb{E}_{X,Y} \big(L(Y, \widehat{Y}(X)) \big) = \sum_{y} \sum_{\mathbf{x}: \, y(\mathbf{x}) = y} L(y, \widehat{y}(\mathbf{x})) \, \mathbb{P}(y \neq \widehat{y}(\mathbf{x})).$$

Ejemplo: $Y \sim Ber(p)$

Si Y toma solamente dos valores, o y 1, entonces podemos escribir

$$\mathbb{E}_{Y|X=\mathbf{x}} \ L(Y,\widehat{y}(\mathbf{x})) = L(O,\widehat{y}(\mathbf{x})) \ \mathbb{P}(Y=O \mid X=\mathbf{x}) + L(1,\widehat{y}(\mathbf{x})) \ \mathbb{P}(Y=1 \mid X=\mathbf{x}).$$

Denotemos por $\lambda_{ij} = L(i,j) = L(y=i,\widehat{y}=j)$, para $i,j \in \{0,1\}$.

Entonces si tomamos el caso binario y el costo de un falso positivo igual a un falso negativo (costo simétrico):

$$\lambda_{00} = L(0,0) = 0, \quad \lambda_{11} = L(1,1) = 0, \quad \lambda_{01} = L(0,1) = 1 = \lambda_{10} = L(1,0),$$

entonces tenemos los costos

$$\operatorname{si} \widehat{y}(\mathbf{x}) = 0: \qquad L(0, \widehat{y}(\mathbf{x})) = \lambda_{00} = 0, \ L(1, \widehat{y}(\mathbf{x})) = \lambda_{10} = 1,$$

$$\operatorname{si} \widehat{y}(\mathbf{x}) = 1$$
: $L(0, \widehat{y}(\mathbf{x})) = \lambda_{01} = 1, L(1, \widehat{y}(\mathbf{x})) = \lambda_{11} = 0.$

y el error sería

$$\operatorname{si} \widehat{y}(\mathbf{x}) = 0$$
: el error es $\mathbb{P}(Y = 1 \mid X = \mathbf{x})$, si $\widehat{y}(\mathbf{x}) = 1$: el error es $\mathbb{P}(Y = 0 \mid X = \mathbf{x})$.

Así, el clasificador bayesiano óptimo es

$$\widehat{\mathbf{y}}(\mathbf{x}) = \begin{cases} 0, & \text{si } \mathbb{P}(\mathbf{Y} = \mathbf{0} \mid X = \mathbf{x}) > \mathbb{P}(\mathbf{Y} = \mathbf{1} \mid X = \mathbf{x}) \\ 1, & \text{si } \mathbb{P}(\mathbf{Y} = \mathbf{1} \mid X = \mathbf{x}) > \mathbb{P}(\mathbf{Y} = \mathbf{0} \mid X = \mathbf{x}) \end{cases}$$
(3)

Obs! En este caso, \hat{y} asigna **x** a la categoría más probable según $\mathbb{P}(Y \mid X = \mathbf{x})$.

Podemos aún simplificar esto usando la regla de Bayes. Escribimos

$$(3) \iff \mathbb{P}(Y = 0 \mid X = \mathbf{x}) > \mathbb{P}(Y = 1 \mid X = \mathbf{x})$$

$$\iff \frac{\mathbb{P}(X = \mathbf{x} \mid Y = 0) \mathbb{P}(Y = 0)}{\mathbb{P}(X = \mathbf{x})} > \frac{\mathbb{P}(X = \mathbf{x} \mid Y = 1) \mathbb{P}(Y = 1)}{\mathbb{P}(X = \mathbf{x})}$$

$$\iff \mathbb{P}(X = \mathbf{x} \mid Y = 0) \mathbb{P}(Y = 0) > \mathbb{P}(X = \mathbf{x} \mid Y = 1) \mathbb{P}(Y = 1)$$

$$\iff f_{O}(\mathbf{x}) \mathbb{P}(Y = 0) > f_{1}(\mathbf{x}) \mathbb{P}(Y = 1);$$

donde la $f_i(\mathbf{x})$ representa la función de densidad o masa de probabilidad condicional

$$f_i(\mathbf{x}) = \mathbb{P}(X = \mathbf{x} \mid Y = i).$$

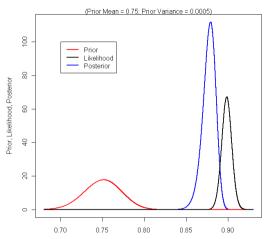
En la ecuación anterior

- $f_i(\mathbf{x}) = \mathbb{P}(X = \mathbf{x} \mid Y = i)$ representa la verosimilitud,
- $\mathbb{P}(Y = i)$ representa la distribución o información a priori de la categoría Y = i,
- mientras que el cociente en la regla de Bayes

$$\frac{f_i(\mathbf{x})\,\mathbb{P}(Y=i)}{\mathbb{P}(X=\mathbf{x})}$$

representa la probabilidad a posteriori.

Entonces, el clasificador Bayesiano óptimo \hat{y} asigna \mathbf{x} a la categoría más probable según la probabilidad a posterior.



La distribución posterior es una mezcla entre la previa y la verosimilitud.

Ejemplo: (Caso general con costo o al clasificar correcto)

El error es

si
$$\widehat{y}(\mathbf{x}) = 0$$
: el error es $\lambda_{10} \mathbb{P}(Y = 1 \mid X = \mathbf{x})$,
si $\widehat{y}(\mathbf{x}) = 1$: el error es $\lambda_{01} \mathbb{P}(Y = 0 \mid X = \mathbf{x})$.

Así, el clasificador bayesiano óptimo es

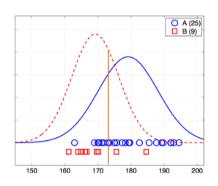
$$\widehat{y}(\mathbf{x}) = \begin{cases} 0, & \text{si } \lambda_{01} \mathbb{P}(Y = 0 \mid X = \mathbf{x}) > \lambda_{10} \mathbb{P}(Y = 1 \mid X = \mathbf{x}) \\ 1, & \text{si } \lambda_{10} \mathbb{P}(Y = 1 \mid X = \mathbf{x}) > \lambda_{01} \mathbb{P}(Y = 0 \mid X = \mathbf{x}) \end{cases}$$

$$(4)$$

Usando la regla de Bayes, obtenemos

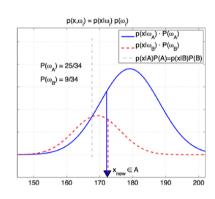
$$\widehat{y}(\mathbf{x}) = \begin{cases} 0, & \text{si } \lambda_{01} f_{0}(\mathbf{x}) \, \mathbb{P}(Y=0) > \lambda_{10} f_{1}(\mathbf{x}) \, \mathbb{P}(Y=1) \\ 1, & \text{si } \lambda_{10} f_{1}(\mathbf{x}) \, \mathbb{P}(Y=1) > \lambda_{01} f_{0}(\mathbf{x}) \, \mathbb{P}(Y=0) \end{cases}$$

$$(5)$$



Maximum Likelihood classifier

Predicted class: Female



Bayes Classifier

Predicted class: Male



Ejemplo: (El caso general)

El error es

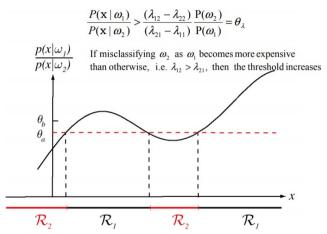
si
$$\widehat{y}(\mathbf{x}) = 0$$
: el error es $\lambda_{00} \mathbb{P}(Y = 0 \mid X = \mathbf{x}) + \lambda_{10} \mathbb{P}(Y = 1 \mid X = \mathbf{x})$, si $\widehat{y}(\mathbf{x}) = 1$: el error es $\lambda_{01} \mathbb{P}(Y = 0 \mid X = \mathbf{x}) + \lambda_{11} \mathbb{P}(Y = 1 \mid X = \mathbf{x})$.

Así, el clasificador bayesiano óptimo es

$$\widehat{\mathbf{y}}(\mathbf{x}) = \begin{cases} \mathbf{0}, & \operatorname{si}(\lambda_{00} - \lambda_{01}) \mathbb{P}(\mathbf{Y} = \mathbf{0} \mid \mathbf{X} = \mathbf{x}) < (\lambda_{11} - \lambda_{10}) \mathbb{P}(\mathbf{Y} = \mathbf{1} \mid \mathbf{X} = \mathbf{x}); \\ \mathbf{1}, & \operatorname{si}(\lambda_{11} - \lambda_{10}) \mathbb{P}(\mathbf{Y} = \mathbf{1} \mid \mathbf{X} = \mathbf{x}) < (\lambda_{00} - \lambda_{01}) \mathbb{P}(\mathbf{Y} = \mathbf{0} \mid \mathbf{X} = \mathbf{x}) \end{cases}$$

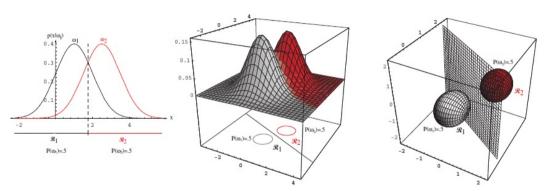
Usando la regla de Bayes, obtenemos

$$\widehat{\mathbf{y}}(\mathbf{x}) = \begin{cases} 0, & \operatorname{si}(\lambda_{00} - \lambda_{01}) f_0(\mathbf{x}) \, \mathbb{P}(\mathbf{Y} = \mathbf{0}) < (\lambda_{11} - \lambda_{10}) f_1(\mathbf{x}) \, \mathbb{P}(\mathbf{Y} = \mathbf{1}); \\ 1, & \operatorname{si}(\lambda_{11} - \lambda_{10}) f_1(\mathbf{x}) \, \mathbb{P}(\mathbf{Y} = \mathbf{1}) < (\lambda_{00} - \lambda_{01}) f_0(\mathbf{x}) \, \mathbb{P}(\mathbf{Y} = \mathbf{0}) \end{cases}$$
(6)



Regla de decisión para el clasificador Bayesiano óptimo.

Ejemplos



Fronteras de decisión del clasificador bayesiano óptimo, para el caso de dos normales $f_i(\mathbf{x})$.

Ejemplos

