

# Ciencia de Datos

## Práctico N°2: Teoría Bayesiana

**Problema 1:** En el caso de dos categorías, según la regla de decisión de Bayes, el error condicional viene dado por la ecuación,

$$P(\text{error}|x) = \min[P(\omega_1|x), P(\omega_2|x)],$$

donde  $\omega_i$  denota los posibles estados del sistema y  $x$  es una variable aleatoria cuyo valor depende del estado del sistema. Incluso si las densidades a posteriori  $P(\omega_i|x)$  son continuas, el error condicional casi siempre conduce a un integrando discontinuo en el calculo del error total

$$P(\text{error}) = \int P(\text{error}|x) p(x) dx.$$

a) Demostrar que para densidades arbitrarias, una cota superior para el error total resulta del hecho de que siempre se cumple que

$$P(\text{error}|x) \leq 2 P(\omega_1|x) P(\omega_2|x).$$

b) Demostrar que si en la expresión para  $P(\text{error})$  se sustituye según  $P(\text{error}|x) = \alpha P(\omega_1|x) P(\omega_2|x)$ , con  $\alpha < 2$ , entonces no puede garantizarse que la integral sea una cota superior para el error.

c) Análogamente, demostrar que puede utilizarse  $P(\text{error}|x) = P(\omega_1|x) P(\omega_2|x)$  para obtener una cota inferior para el error total.

d) Demostrar que si  $P(\text{error}|x) = \beta P(\omega_1|x) P(\omega_2|x)$  con  $\beta > 1$ , entonces la integral puede no ser una cota inferior para el error.

**Problema 2:** Suponer dos variables aleatorias independientes idénticamente distribuídas con la densidad de Laplace,

$$p(x|\omega_i) \propto \exp\left(-\frac{|x - a_i|}{b_i}\right), \text{ con } i = 1, 2 \text{ y } b_i > 0.$$

a) Escribir las expresiones analíticas normalizadas de  $p(x|\omega_i)$ .

b) Calcular el radio de verosimilitud como función de los parámetros.

c) Graficar el radio  $p(x|\omega_1)/p(x|\omega_2)$  para el caso  $a_1 = 0$ ,  $b_1 = 1$ ,  $a_2 = 1$  y  $b_2 = 2$ .

**Problema 3:** Considerar la siguiente regla de decisión para el problema unidimensional con dos categorías: Se decide por  $\omega_1$  si  $x > \theta$  y en otro caso se decide por  $\omega_2$ .

a) Demuestrar que la probabilidad de error para esta regla viene dada por

$$P(\text{error}) = \int P(\text{error}|x)p(x)dx = P(\omega_1) \int_{-\infty}^{\theta} p(x|\omega_1)dx + P(\omega_2) \int_{\theta}^{\infty} p(x|\omega_2)dx.$$

b) Demostrar que una condición necesaria para minimizar el error es  $p(\theta|\omega_1) p(\omega_1) = p(\theta|\omega_2) p(\omega_2)$ .

c) Define esta ecuación un valor de  $\theta$  único?

d) Estudiar como ejemplo el caso en el que la variable  $X$  condicional a  $\omega_i$  tiene distribución normal con media  $\mu_i$  y desvío  $\sigma_i$ ; es decir,  $P(X|\omega_i) \sim N(\mu_i, \sigma_i)$ .

**Problema 4:** Suponer que se sustituye la función de decisión determinística  $\alpha(x)$  por la regla aleatoria dada por la probabilidad  $P(\alpha_i|x)$  de tomar la decisión  $\alpha_i$  dado que se observe  $x$ .

a) Mostrar que el riesgo resultante viene dado por,

$$R = \int \left( \sum_{i=1}^a R(\alpha_i|x) P(\alpha_i|x) \right) p(x) dx.$$

b) Demostrar además que  $R$  se minimiza para  $P(\alpha_i|x) = 1$  para la acción  $\alpha_i$  asociada con el riesgo condicional mínimo  $R(\alpha_i|x)$ , lo que demuestra que no obtenemos ningún beneficio haciendo aleatoria la regla de decisión.

**Problema 5:** En muchos problemas de clasificación multicategoría:  $\omega_i$  con  $i = 1, \dots, c$ , es conveniente trabajar con una función de pérdida pesada. Por ejemplo, puede ocurrir que se rechace un patrón o estado del sistema si este resulta irreconocible,

$$\lambda(\alpha_i|\omega_j) = \begin{cases} 0 & \text{si } i = j, i, j = 1, 2, \dots, c, \\ \lambda_r & \text{si } i = c + 1, \\ \lambda_s & \text{en otro caso.} \end{cases}$$

donde  $\lambda_r$  es la pérdida sufrida por la elección de rechazo,  $\lambda_s$  es la pérdida incurrida por cometer un error.

Mostrar que el riesgo mínimo se obtiene si decidimos  $\alpha_i$  si  $P(\omega_i|x) \geq P(\omega_j|x)$  para todo  $j$ , y si  $P(\omega_i|x) \geq 1 - \frac{\lambda_r}{\lambda_s}$ , caso contrario, rechazar. ¿Que sucede si  $\lambda_r = 0$ ? ¿Que sucede si  $\lambda_r > \lambda_s$ ?

**Problema 6:** Retomar el problema de clasificación con la opción de rechazo del problema anterior.

a) Demostrar que las siguientes funciones discriminantes son óptimas para este tipo de problemas:

$$g_i(x) = \begin{cases} p(x|\omega_i) P(\omega_i) & \text{si } i = 1, 2, \dots, c, \\ \frac{\lambda_s - \lambda_r}{\lambda_s} \sum_{j=1}^c p(x|\omega_j) P(\omega_j) & \text{si } i = c + 1. \end{cases}$$

b) Graficar la función discriminante y las regiones de decisión para el caso del problema unidimensional ( $x \in \mathcal{R}$ ) con dos clases usando los valores

$$p(x|\omega_1) \sim \mathcal{N}(1, 1), p(x|\omega_2) \sim \mathcal{N}(-1, 1), P(\omega_1) = P(\omega_2), \frac{\lambda_r}{\lambda_s} = \frac{1}{4}.$$

c) Describir cualitativamente lo que sucede cuando  $\frac{\lambda_r}{\lambda_s}$  se incrementa desde 0 a 1.

d) Considerar nuevamente este problema, ahora en el caso particular

$$p(x|\omega_1) \sim \mathcal{N}(1, 1), p(x|\omega_2) \sim \mathcal{N}\left(0, \frac{1}{4}\right), P(\omega_1) = \frac{1}{3}, P(\omega_2) = \frac{2}{3}, \frac{\lambda_r}{\lambda_s} = \frac{1}{2}.$$

**Problema 7:** Estudiar la implementación del análisis de discriminante lineal provista por scikit-learn para generar muestras aleatorias de acuerdo a una distribución normal bivariada y calcular la función discriminante para una distribución normal dada y probabilidades a priori  $P(\omega_i)$ .

a) Simular dos variables normales  $(X_1, X_2)$  con  $\Sigma = C^T.C$ , y  $C = \begin{pmatrix} 0 & -0,23 \\ 0,83 & 0,23 \end{pmatrix}$  y vectores de medias  $\mu_1 = (0, 0)$  y  $\mu_2 = (1, 1)$ , respectivamente.

b) Suponer que las probabilidades a priori de las dos categorías son iguales ( $P(\omega_1) = P(\omega_2)$ ), e implementar un clasificador para dos categorías utilizando sólo el valor de característica  $X_1$  especificada en el inciso anterior. El código resultante debe poder clasificar una nueva muestra basado en esta información.

Tener presente que para el diseño del clasificador se estimará la media y varianza a partir de los datos de cada una de las muestras. Si para la muestra  $i$  la media y varianza son  $\mu_i$  y  $\sigma_i^2$  respectivamente, se clasificará un valor  $x$  en la muestra 1 si

$$\frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-(x-\mu_1)^2/2\sigma_1^2} P(\omega_1) > \frac{1}{\sqrt{2\pi\sigma_2^2}} e^{-(x-\mu_2)^2/2\sigma_2^2} P(\omega_2).$$

Tomando logaritmo, y eliminando las probabilidades a priori  $P(\omega_i)$  por ser iguales esto es equivalente a decidir por la clase 1 si

$$-\frac{1}{2} \ln(2\pi) - \ln \sigma_1 - \frac{(x - \mu_1)^2}{2\sigma_1^2} > -\frac{1}{2} \ln(2\pi) - \ln \sigma_2 - \frac{(x - \mu_2)^2}{2\sigma_2^2};$$

es decir, si

$$\ln \sigma_1 + \frac{(x - \mu_1)^2}{2\sigma_1^2} < \ln \sigma_2 + \frac{(x - \mu_2)^2}{2\sigma_2^2}.$$

c) Determinar el error de entrenamiento empírico en la clasificación muestras; esto es, el porcentaje de puntos mal clasificados, dividiendo aleatoriamente el número de muestras  $n = 100$ , en 80 % entrenamiento y en 20 % test. Repetir incrementando los valores de  $n$ , desde 100 a 10000 en pasos de 100 y graficar el error empírico obtenido.

d) Utilizar la cota de Bhattacharyya para acotar el error que obtendrán los nuevos patrones obtenidos muestreando las distribuciones.

e) Repetir todo lo anterior, pero ahora utilice las dos características,  $X_1$  y  $X_2$ .

f) Analizar resultados. ¿Es siempre posible para un conjunto finito de datos que el error empírico resulte mayor al aumentar la dimensión de los datos?

**Problema 8:** La distribución de Poisson para una variable entera no negativa  $x = 0, 1, \dots$  y parámetro real  $\lambda$  viene dada por

$$P(x|\lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

Considerar el problema de clasificación con dos categorías igualmente probables  $P(\omega_1) = P(\omega_2)$  y condicionales con distribuciones de Poisson con diferentes parámetros  $\lambda_1 > \lambda_2$ .

a) Especificar regla de clasificación de Bayes.

b) ¿Cuál es la tasa del error de Bayes?

c) Escribir función discriminante, y determinar qué valores debe tener de entrada para clasificar un nuevo dato.

d) Simular una muestra aleatoria de tamaño 100 con distribuciones de Poisson con  $\lambda_1 = 1,8$ ,  $\lambda_2 = 0,4$ , considerando igual probabilidad a priori. Usar la función de pérdida cero uno y clasificar la muestra acorde a esta función. Estimar el error cometido en la muestra, y compararlo con el error de Bayes calculado.

## Soluciones:

### Problema 1:

a) Suponiendo que para un determinado valor  $x$  se tiene que  $P(\omega_1|x) \leq P(\omega_2|x)$ , entonces  $P(\text{error}|x) = P(\omega_1|x)$ . La condición de normalización para un sistema de dos estados,  $P(\omega_1|x) + P(\omega_2|x) = 1$ , permite escribir  $P(\omega_1|x) = 1 - P(\omega_2|x)$ . Así,  $1 - P(\omega_2|x) \leq P(\omega_2|x)$ , es decir  $2P(\omega_2|x) \geq 1$ . Multiplicando ambos miembros  $P(\omega_1|x)$  se tiene que  $2P(\omega_1|x)P(\omega_2|x) \geq P(\omega_1|x)$ . De esta forma, siempre puede escribirse  $P(\text{error}|x) \leq 2P(\omega_1|x)P(\omega_2|x)$  y para el error total resulta

$$P(\text{error}) \leq \int 2P(\omega_1|x)P(\omega_2|x)p(x)dx$$

b) Como ejemplo tomar  $\alpha = 4/3$  y suponer  $P(\omega_1|x) = 0,4$ , entonces  $P(\omega_2|x) = 0,6$ . De esta forma resulta  $P(\text{error}|x) = 0,4$ , mientras que  $\alpha P(\omega_1|x)P(\omega_2|x) = 0,32 < P(\text{error}|x)$ .

c) Dado que  $\max[P(\omega_1|x), P(\omega_2|x)] < 1$ , resulta  $P(\text{error}|x) = \min[P(\omega_1|x), P(\omega_2|x)] \geq P(\omega_1|x)P(\omega_2|x)$ . Entonces  $P(\text{error}) = \int \min[P(\omega_1|x), P(\omega_2|x)]p(x)dx \geq \int P(\omega_1|x)P(\omega_2|x)p(x)dx$ .

d) Funciona el mismo ejemplo usado en (b).

### Problema 2:

a)

$$\begin{aligned} \int_{-\infty}^{\infty} \exp\left(-\frac{|x-a_i|}{b_i}\right) dx &= \int_{-\infty}^{a_i} \exp\left(-\frac{a_i-x}{b_i}\right) dx + \int_{a_i}^{\infty} \exp\left(-\frac{x-a_i}{b_i}\right) dx \\ &= b_i e^{-a_i/b_i} (e^{a_i/b_i} - 0) - b_i e^{a_i/b_i} (0 - e^{-a_i/b_i}) \\ &= 2b_i. \end{aligned}$$

Entonces,

$$p(x|\omega_i) = \frac{1}{2b_i} \exp\left(-\frac{|x-a_i|}{b_i}\right).$$

b) El radio de verosimilitud resulta

$$\frac{p(x|\omega_1)}{p(x|\omega_2)} = \frac{2b_2}{2b_1} \exp\left(-\frac{|x-a_1|}{b_1} + \frac{|x-a_2|}{b_2}\right).$$

c) En el caso  $a_1 = 0$ ,  $b_1 = 1$ ,  $a_2 = 1$  y  $b_2 = 2$

$$\frac{p(x|\omega_1)}{p(x|\omega_2)} = \begin{cases} 2e^{(x+1)/2} & \text{si } x \leq 0, \\ 2e^{(1-3x)/2} & \text{si } 0 \leq x \leq 1, \\ 2e^{-(x+1)/2} & \text{si } 1 \leq x. \end{cases}$$

### Problema 3:

a) Dada la regla de decisión, se comete error al elegir  $\omega_1$  cuando  $x \leq \theta$  o si se elige  $\omega_2$  cuando  $x > \theta$ . Entonces, usando el teorema de Bayes se obtiene,

$$\begin{aligned} P(\text{error}) &= \int P(\text{error}|x)p(x)dx = \int_{-\infty}^{\theta} P(\omega_1|x)p(x)dx + \int_{\theta}^{\infty} P(\omega_2|x)p(x)dx \\ &= \int_{-\infty}^{\theta} \frac{p(x|\omega_1)P(\omega_1)}{p(x)}p(x)dx + \int_{\theta}^{\infty} \frac{p(x|\omega_2)P(\omega_2)}{p(x)}p(x)dx \\ &= P(\omega_1) \int_{-\infty}^{\theta} p(x|\omega_1)dx + P(\omega_2) \int_{\theta}^{\infty} p(x|\omega_2)dx. \end{aligned}$$

b) La condición necesaria para la existencia de un valor extremo es  $\frac{d}{d\theta}P(\text{error}) = 0$ . Usando el teorema fundamental del cálculo resulta,

$$\frac{d}{d\theta}P(\text{error}) = P(\omega_1)p(\theta|\omega_1) - P(\omega_2)p(\theta|\omega_2) = 0$$

c) En el caso que  $P(\text{error})$  sea una función monótona de  $\theta$ , La ecuación no solución, y los máximos y mínimos se encontrarán en los valores extremos de  $\theta$ . La ecuación puede también tener más de una solución. En el caso de dos soluciones una se corresponderá con el mínimo y otra con el máximo.

d) Teniendo en cuenta que  $p(\theta|\omega_i) = \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{1}{2}\left(\frac{\theta - \mu_i}{\sigma_i}\right)^2\right)$ , sustituyendo en la expresión para  $\frac{d}{d\theta}P(\text{error})$  se obtiene

$$\frac{d}{d\theta}P(\text{error}) = \frac{P(\omega_1)}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{1}{2}\left(\frac{\theta - \mu_1}{\sigma_1}\right)^2\right) - \frac{P(\omega_2)}{\sqrt{2\pi}\sigma_2} \exp\left(-\frac{1}{2}\left(\frac{\theta - \mu_2}{\sigma_2}\right)^2\right)$$

Para simplificar los cálculos supongamos que  $\frac{P(\omega_1)}{\sigma_1} = \frac{P(\omega_2)}{\sigma_2}$ . Para esto tomenos como ejemplo particular  $P(\omega_1) = \frac{1}{3}$ ,  $P(\omega_2) = \frac{2}{3}$  y supongamos que  $\mu_1 = 0$ ,  $\mu_2 = 2$ ,  $\sigma_1^2 = 1$ ,  $\sigma_2^2 = 4$ . La condición de derivada nula implica

$$\exp\left(\frac{1}{2}\left(\frac{\theta - 2}{2}\right)^2 - \frac{\theta^2}{2}\right) = 1 \Rightarrow (\theta - 2)^2 - 4\theta^2 = 0 \Rightarrow 3\theta^2 + 4\theta - 4 = 0.$$

Las soluciones de esta ecuación son **-2** y **2/3**. Por otro lado, la derivada es positiva si se cumple la condición

$$\theta^2 < \left(\frac{\theta - 2}{2}\right)^2 \Rightarrow 3\theta^2 + 4\theta - 4 < 0.$$

Resulta así que  $P(\text{error})$  es una función creciente de  $\theta$  en el intervalo  $(-2, 2/3)$  y entonces tiene un mínimo para  $\theta = -2$  y un máximo en  $\theta = 2/3$ .

#### Problema 4:

a) Si el riesgo fuera determinístico, para cada  $x$  se cumple  $\alpha(x) = \alpha_i$  para algún  $i = 1, \dots, a$ . En este caso se considera un riesgo aleatorio, es decir,  $\alpha(x) = \alpha_i$ , con probabilidad  $P(\alpha_i|x)$ , con  $i = 1, \dots, a$ . Luego el riesgo esperado de  $\alpha$  dado  $x$  es la suma ponderada de sus posibles valores,

$$R(\alpha(x)|x) = \sum_{i=1}^a R(\alpha_i|x) P(\alpha_i|x).$$

Así, el riesgo total resulta,

$$R = \int R(\alpha(x)|x) p(x) dx = \int \sum_{i=1}^a R(\alpha_i|x) P(\alpha_i|x) p(x) dx.$$

b) Para minimizar  $R$  basta con minimizar  $R(\alpha(x)|x)$ . Teniendo en cuenta que

$$R(\alpha_j|x) = \min_{1 \leq i \leq a} R(\alpha_i|x)$$

y dado que  $\sum_i P(\alpha_i|x) = 1$  se obtiene  $\sum_{i=1}^a R(\alpha_i|x) P(\alpha_i|x) \geq \sum_{i=1}^a R(\alpha_j|x) P(\alpha_i|x) = R(\alpha_j|x)$ . En particular, eligiendo  $P(\alpha_j|x) = 1$  se minimiza el riesgo, pero esto le quita la aleatoriedad a la regla.

**Problema 5:** El riesgo total se minimiza, minimizando cada uno de sus términos:  $R(\alpha_i|x)$ ,

$$R(\alpha_i|x) = \sum_{j=1}^c \lambda(\alpha_i|\omega_j) P(\omega_j|x) = \sum_{\substack{j=1 \\ j \neq i}}^c \lambda_s P(\omega_j|x) = \lambda_s (1 - P(\omega_i|x)) \text{ con } i = 1, \dots, c,$$

$$R(\alpha_{c+1}|x) = \sum_{j=1}^c \lambda(\alpha_{c+1}|\omega_j) P(\omega_j|x) = \lambda_r \sum_{j=1}^c P(\omega_j|x) = \lambda_r.$$

Para  $1 \leq i \leq c$ , el riesgo se minimiza si  $P(\omega_i|x)$  es máximo en el conjunto de los  $P(\omega_j|x)$ , ya que  $\lambda(\alpha_i|\omega_i) = 0$ . Ahora,  $R(\alpha_i|x) \leq R(\alpha_{c+1}|x)$  si y sólo si

$$\lambda_s (1 - P(\omega_i|x)) \leq \lambda_r \Leftrightarrow 1 - P(\omega_i|x) \leq \frac{\lambda_r}{\lambda_s} \Leftrightarrow P(\omega_i|x) \geq 1 - \frac{\lambda_r}{\lambda_s}.$$

Si  $\lambda_r = 0$ , entonces el riesgo total se minimiza si se rechaza en todos los casos.

Si  $\lambda_r > \lambda_s$ , entonces en ningún caso conviene rechazar puesto que  $1 - \frac{\lambda_r}{\lambda_s} < 0$ .

Entonces, sólo tiene sentido el caso  $0 < \lambda_r < \lambda_s$ .

**Problema 6:**

a) Según esta función discriminante, dada la característica  $x$  se elige la clase  $\omega_i$  si

$$p(x|\omega_i) P(\omega_i) \geq p(x|\omega_j) P(\omega_j) \text{ para } j \neq i,$$

y además no se la rechaza; es decir,

$$\lambda_s p(x|\omega_i) P(\omega_i) \geq (\lambda_s - \lambda_r) \sum_{j=1}^c p(x|\omega_j) P(\omega_j).$$

Usando el teorema de Bayes y suponiendo  $p(x) > 0$  de la primera desigualdad resulta  $P(\omega_i|x) \geq P(\omega_j|x)$  para  $i \neq j$ , mientras que cancelando el término  $\lambda_s p(x|\omega_i) P(\omega_i)$  de ambos miembros de la segunda desigualdad y teniendo en cuenta que  $\sum_{j=1}^c p(x|\omega_j) P(\omega_j) = p(x)$ , la segunda desigualdad puede reescribirse según

$$\lambda_s \sum_{j \neq i} p(x|\omega_j) P(\omega_j) = \lambda_s (p(x) - p(x|\omega_i) P(\omega_i)) \leq \lambda_r p(x).$$

Dividiendo ambos miembros por  $p(x)$  y observando que  $p(x|\omega_j) P(\omega_j) = P(\omega_j|x) p(x)$  se concluye que  $g_i(x)$  clasifica en  $\omega_i$  si

$$P(\omega_i|x) \geq P(\omega_j|x) \text{ para } j \neq i \text{ y } P(\omega_i|x) \geq 1 - \frac{\lambda_r}{\lambda_s}.$$

Por lo tanto, esta función discriminante produce la elección de la clase  $\omega_i$  que genera el riesgo mínimo, según lo probado en el Problema 1.



FaMAF 2023