

Examen de los temas 1 a 4 de Aprendizaje Automático

ETSINF, Universitat Politècnica de València, 02 de diciembre de 2013

Apellidos:

Nombre:

Cuestiones (2 puntos; tiempo estimado: 30 minutos)

Marca cada recuadro con una única opción de entre las dadas.

- 1 ☐ B Deseamos evaluar un sistema de Aprendizaje Automático utilizando un conjunto de datos de entrenamiento que contiene 1000 muestras y la técnica de *exclusión individual* ("Leaving One Out"), obteniéndose un total de 44 errores. Indicar cuál de las afirmaciones siguientes es correcta:

- A) La talla de entrenamiento efectiva es de 1000 muestras y la talla de test efectiva es 1000 muestras.
- B) La talla de entrenamiento efectiva es de 999 muestras y el error es del 4.4 %
- C) La talla de entrenamiento efectiva es de 900 muestras y el error es del 44 %
- D) La talla de entrenamiento efectiva es de 1000 muestras y la talla de test efectiva es 900 muestras.

- 2 ☐ C Al aplicar la técnica de descenso por gradiente a una modificación de la función de Widrow y Hoff

$$q_S(\boldsymbol{\theta}) = \frac{1}{2} \left(\sum_{n=1}^N (\boldsymbol{\theta}^t \mathbf{x}_n - y_n)^2 + \boldsymbol{\theta}^t \boldsymbol{\theta} \right),$$

el gradiente $\nabla q_S(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}(k)}$ en la iteración $\boldsymbol{\theta}(k+1) = \boldsymbol{\theta}(k) - \rho_k \nabla q_S(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}(k)}$ es

- A) $\sum_{n=1}^N (\boldsymbol{\theta}(k)^t \mathbf{x}_n - y_n) \mathbf{x}_n$
- B) $\sum_{n=1}^N (\boldsymbol{\theta}(k)^t \mathbf{x}_n - y_n) \mathbf{x}_n + \boldsymbol{\theta}(k)^t \boldsymbol{\theta}(k)$
- C) $\sum_{n=1}^N (\boldsymbol{\theta}(k)^t \mathbf{x}_n - y_n) \mathbf{x}_n + \boldsymbol{\theta}(k)$
- D) $\sum_{n=1}^N (\boldsymbol{\theta}(k)^t \mathbf{x}_n - y_n) \mathbf{x}_n + \mathbf{x}_n$

- 3 ☐ C Entre las siguientes propiedades de las funciones discriminantes lineales hay una que es falsa:

- A) La función discriminante lineal aplicada en un punto devuelve un valor proporcional a la distancia del punto al correspondiente hiperplano separador.
- B) La distancia del origen de coordenadas al hiperplano separador asociado a una función discriminante lineal es $\frac{\theta_0}{\|\boldsymbol{\theta}\|}$
- C) Un hiperplano separador tiene asociado una única función discriminante lineal canónica
- D) Un hiperplano separador tiene asociado un número infinito de funciones discriminantes lineales

- 4 ☐ B Se quiere aplicar la técnica esperanza-maximización a un problema de estimación de máxima verosimilitud en el que no hay variables latentes o ocultas. En este caso ¿Cuál de las afirmaciones siguientes es correcta?

- A) En ese caso no se puede aplicar la técnica esperanza-maximización.
- B) En ese caso solo se aplica la etapa de maximización y en una iteración acaba.
- C) En ese caso solo se aplica la etapa de maximización y hay que iterar hasta que converja.
- D) En ese caso solo se aplica la etapa del cálculo de la esperanza.

Problema 1 (4 puntos; tiempo estimado: 40 minutos)

En una tarea de clasificación de correos electrónicos como spam o no-spam se dispone de un conjunto S de 500 correos *no-spam* (clase A) y 300 *spam* (clase B).

- ¿Cuál sería el logaritmo de la verosimilitud de los mensajes de S si las probabilidades a priori $P(A)$ y $P(B)$ fueran iguales?
- Las probabilidades a priori puede estimarse por máxima verosimilitud a partir de S como: $P(A) = 5/8$, $P(B) = 3/8$. Derivar estas probabilidades mediante la técnica de optimización de los multiplicadores de Lagrange
- Calcular el logaritmo de la verosimilitud de S según las probabilidades a priori obtenidas en b) y compararla con la obtenida en a)

- Modelo: $\Theta_0 \equiv (p_A, p_B)^t$: $p_A \equiv P(c = A) = 0.5$, $p_B \equiv P(c = B) = 0.5$
 - El logaritmo de la verosimilitud es S , $P(S | \Theta_0)$

$$\log P(S | \Theta_0) = \log\left(\prod_{i=1}^{500} p_A \prod_{j=1}^{300} p_B\right) = 500 \cdot \log 0.5 + 300 \cdot \log 0.5 = -240.82 \Rightarrow P(S | \Theta_0) = 6.66 \cdot 10^{-240}$$

- Modelo: $\Theta \equiv (p_A, p_B)^t$, con $p_A + p_B = 1$
 - Verosimilitud y logaritmo de la verosimilitud:

$$P(S | \Theta) = \prod_{i=1}^{500} p_A \prod_{j=1}^{300} p_B = p_A^{500} p_B^{300}$$

$$L_S(\Theta) = \log P(S | \Theta) = 500 \log p_A + 300 \log p_B$$

- Estimación de máxima verosimilitud:

$$\Theta^* = \arg \max_{\Theta} L_S(\Theta) = \arg \max_{\substack{p_A, p_B \\ p_A + p_B = 1}} (500 \log p_A + 300 \log p_B)$$

- Lagrangiana: $\Lambda(p_A, p_B, \beta) = 500 \log p_A + 300 \log p_B + \beta(1 - p_A - p_B)$
- Soluciones óptimas en función del multiplicador de Lagrange:

$$\left. \begin{aligned} \frac{\partial \Lambda}{\partial p_A} &= \frac{500}{p_A} - \beta = 0 \\ \frac{\partial \Lambda}{\partial p_B} &= \frac{300}{p_B} - \beta = 0 \end{aligned} \right\} \quad \begin{aligned} p_A^*(\beta) &= \frac{500}{\beta} \\ p_B^*(\beta) &= \frac{300}{\beta} \end{aligned}$$

- Función dual de Lagrange:

$$\Lambda_D(\beta) = 500 \log \frac{500}{\beta} + 300 \log \frac{300}{\beta} + \beta \left(1 - \frac{500}{\beta} - \frac{300}{\beta}\right) = \beta - 800 \log \beta - 800 + 500 \log 500 + 300 \log 300$$

- Valor óptimo del multiplicador de Lagrange: $\frac{d\Lambda_D}{d\beta} = 1 - \frac{800}{\beta} = 0 \Rightarrow \beta^* = 800$

- Solución final: $\theta^* = (p_A^*, p_B^*)^t$: $p_A^* = p_A^*(\beta^*) = \frac{5}{8}$ $p_B^* = p_B^*(\beta^*) = \frac{3}{8}$

- Como en a):

$$L_S(\Theta^*) = 500 \cdot \log(5/8) + 300 \cdot \log(3/8) = -229.85 \Rightarrow P(S | \Theta^*) = 7.09 \cdot 10^{-229} \gg 6.66 \cdot 10^{-240} = P(S | \Theta_0)$$

La verosimilitud es mayor que en a) debido a que se ha maximizado la verosimilitud con respecto a Θ .

Problema 2 (4 puntos; tiempo estimado: 20 minutos)

Para el aprendizaje de una máquina de vectores soporte se dispone de una muestra de entrenamiento linealmente separable

$$S = \{((1, 4), +1), ((1, 6), +1), ((2, 2), +1), ((2, 3), +1), ((4, 2), -1), ((3, 4), -1), ((3, 5), -1), ((5, 4), -1), ((5, 6), -1)\}$$

Los multiplicadores de Lagrange óptimos son: $\boldsymbol{\alpha}^* = (0, 0.25, 0, 1.0, 0, 1.25, 0, 0, 0)^t$.

- Obtener la función discriminante lineal correspondiente
- Calcular el margen óptimo
- Clasificar la muestra $(4, 5)$.

■ La función discriminante lineal

- El vector de pesos:

$$\theta_1^* = +1 \cdot 0.25 \cdot 1 + 1 \cdot 1.0 \cdot 2 - 1 \cdot 1.25 \cdot 3 = -1.5$$

$$\theta_2^* = +1 \cdot 0.25 \cdot 6 + 1 \cdot 1.0 \cdot 3 - 1 \cdot 1.25 \cdot 4 = -0.5$$

- El peso umbral (con la muestra 2) $\theta_0^* = (+1) - (-1.5 \cdot 1 - 0.5 \cdot 6) = 5.5$
- La FDL: $\phi(\mathbf{x}) = -1.5 \cdot x_1 - 0.5 \cdot x_2 + 5.5$

■ Margen óptimo:

$$\frac{2}{\|\boldsymbol{\theta}^*\|} = \frac{2}{\sqrt{0.25 + 1.0 + 0, 1.25}} = 1.26$$

- Clasificación de la muestra $(4, 5)$: $\phi(4, 5) = -1.5 \cdot 4 - 0.5 \cdot 5 + 5.5 = -3 < 0 \Rightarrow \text{clase} = -1$