

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 10 de enero de 2022

Apellidos:  Nombre:  Grupo:

**Cuestiones (2 puntos; tiempo estimado: 30 minutos)**

Marca cada recuadro con una única opción de entre las dadas. Cada acierto suma 1/2 puntos y cada fallo resta 1/6 puntos.

- 1 ☐ A Se ha evaluado un sistema de Aprendizaje Automático mediante la técnica de *validación cruzada en B bloques* (“B-fold Cross Validation”) con  $B = 100$  y utilizando un conjunto de datos etiquetados que contiene 10000 muestras. Se han obtenido un total de 50 errores. Indicar cuál de las afirmaciones siguientes es razonable:

- A) La talla de entrenamiento efectiva es 9900 muestras y el error estimado es  $0.50\% \pm 0.14\%$
- B) La talla de entrenamiento efectiva es de 9900 muestras y el error estimado es  $5.0\%$
- C) La talla de test efectiva es de 10000 muestras y el error estimado es  $5.0\% \pm 0.14\%$
- D) El error estimado es  $0.50\% \pm 0.0014\%$ .

- 2 ☐ C Considerar la siguiente modificación de la función de Widrow y Hoff

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

Al aplicar la técnica de descenso por gradiente, en la iteración  $k$  el vector de pesos,  $\boldsymbol{\theta}$ , se modifica como:  $\boldsymbol{\theta}(k+1) = \boldsymbol{\theta}(k) - \rho_k \nabla q_S(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}(k)}$ . En esta expresión, el gradiente,  $\nabla q_S(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}(k)}$ , es:

- A)  $\sum_{n=1}^N \mathbf{x}_n + \lambda \boldsymbol{\theta}(k) \left( \sum_{n=1}^N \mathbf{x}_n \right)^t \left( \sum_{n=1}^N \mathbf{x}_n \right)$
- B)  $\sum_{n=1}^N \mathbf{x}_n + \lambda \boldsymbol{\theta}(k)$
- C)  $\sum_{n=1}^N \mathbf{x}_n + 2 \lambda \boldsymbol{\theta}(k) \left( \sum_{n=1}^N \mathbf{x}_n \right)^t \left( \sum_{n=1}^N \mathbf{x}_n \right)$
- D)  $\sum_{n=1}^N \boldsymbol{\theta}(k)^t \mathbf{x}_n + \lambda \left( \sum_{n=1}^N \mathbf{x}_n \right)^t \left( \sum_{n=1}^N \mathbf{x}_n \right)$

- 3 ☐ D En el compromiso entre sesgo y varianza indicar qué afirmación es correcta.

- A) Si el modelo presenta una gran varianza quiere decir que el error de entrenamiento es alto.
- B) Cuando el sesgo es muy alto, la varianza suele ser también alta.
- C) Si el modelo presenta una sesgo alto quiere decir que el error de entrenamiento es bajo.
- D) Cuando el sesgo es muy bajo, la varianza suele ser alta.

- 4 ☐ D En el aprendizaje de modelos probabilísticos mediante el algoritmo esperanza-maximización indicar qué afirmación es correcta.

- A) Es el método adecuado cuando las muestra de entrenamiento es completa.
- B) A partir de una cierta estimación de las variables ocultas se busca el óptimo de una función auxiliar como solución final del problema original.
- C) Es el método adecuado cuando no hay variables ocultas.
- D) En cada iteración, a partir de una cierta estimación de las variables ocultas se busca el óptimo de una función auxiliar.



## Problema 1 (3 puntos; tiempo estimado: 30 minutos)

Para entrenar un modelo basado en máquinas de vectores soporte, se dispone de un conjunto de entrenamiento en  $\mathbb{R}^2$ . Estos vectores y los correspondientes multiplicadores de Lagrange óptimos obtenidos con  $C = 10$  son:

$i$	1	2	3	4	5	6	7	8
$x_{i1}$	2	2	2	2	3	4	3	1
$x_{i2}$	2	3	4	1	2	2	1	4
Clase	+1	-1	+1	-1	-1	-1	-1	+1
$\alpha_i^*$	10.0	10.0	3.78	3.11	0.67	0	0	0

- Obtener la función discriminante lineal correspondiente
- Obtener la ecuación de la frontera lineal de separación entre clases y representarla gráficamente junto con los vectores de entrenamiento, indicando cuáles de ellos son vectores soporte.
- Obtener la tolerancia óptima de cada muestra de entrenamiento.
- Clasificar la muestra  $(1, 2)^t$ .

a) Pesos de la función discriminante:

$$\begin{aligned}\theta_1^* &= (+1)(2)(10.0) + (-1)(2)(10.0) + (+1)(2)(3.79) + (-1)(2)(3.11) + (-1)(3)(0.67) = -0.67 \\ \theta_2^* &= (+1)(2)(10.0) + (-1)(3)(10.0) + (+1)(4)(3.79) + (-1)(1)(3.11) + (-1)(2)(0.67) = 0.67\end{aligned}$$

Usando el vector soporte  $\mathbf{x}_4$  (que verifica la condición :  $0 < \alpha_1^* < C$ )

$$\theta_0^* = c_4 - \theta^{*t} \mathbf{x}_4 = 1 - ((-0.67)(2) + (0.67)(1)) = -0.33$$

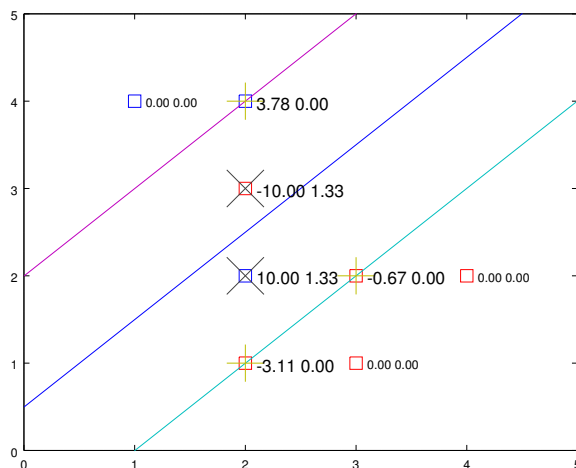
Función discriminante lineal:  $\phi(\mathbf{x}) = -0.33 - 0.67 x_1 + 0.67 x_2$

b) Frontera de separación y representación gráfica:

Ecuación de la frontera lineal de separación:  $-0.33 - 0.67 x_1 + 0.67 x_2 = 0 \rightarrow x_2 = 1.0 x_1 + 0.49$ .

Los vectores de entrenamiento son todos los de la tabla. De ellos, los vectores soporte son:  $(2, 1)^t, (2, 2)^t, (2, 3)^t, (2, 4)^t, (3, 2)^t$ .

Representación gráfica:



Al lado de cada muestra se muestra el valor del multiplicador de lagrange asociado y la tolerancia.

c) Todas las muestras bien clasificadas y fuera del margen ( $i \in \{3, 4, 5, 6, 7, 8\}$ ) tienen una tolerancia  $\zeta_i^* = 0$  y el resto

$$\zeta_1^* = 1 - c_1 (\theta^{*t} \mathbf{x}_1 + \theta_0^*) = 1.33; \quad \zeta_2^* = 1 - c_2 (\theta^{*t} \mathbf{x}_2 + \theta_0^*) = 1.33$$

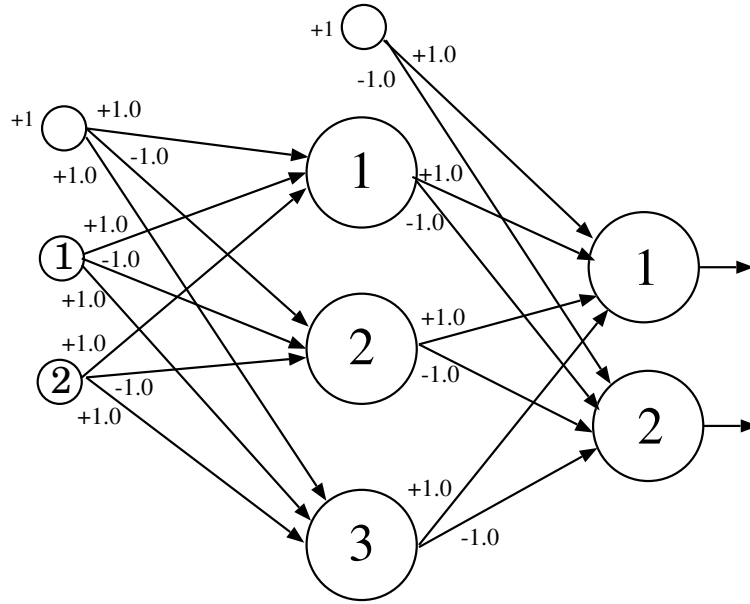
d) Clasificación de la muestra  $(1, 2)^t$ :

El valor de la función discriminante para este vector es:  $\theta_0^* + 1\theta_1^* + 2\theta_2^* = 0.34 > 0 \Rightarrow$  clase +1.



## Problema 2 (3 puntos; tiempo estimado: 30 minutos)

El perceptrón multicapa de la figura se utiliza para resolver un problema de regresión, con función de activación de los nodos de la capa de salida y de la capa oculta de tipo sigmoid, y factor de aprendizaje  $\rho = 1.0$ .



Dado un vector de entrada  $x = (+1, +1)$  y su valor deseado de salida  $t = (+1, 0)$ ,  
Calcular:

- las salidas que faltan
- los correspondientes errores en los nodos de la capa de salida y en los nodos de la capa oculta.
- Los nuevos valores de los pesos de las conexiones  $\theta_{3,2}^1$ , que va del nodo 2 de entrada al nodo 3 de la capa oculta, y  $\theta_{2,3}^2$ , que va del nodo 3 de la capa oculta al nodo 2 de la capa de salida.

- Las salidas de todos los nodos

$$\begin{aligned}
 \phi_1^1 &= \theta_{1,1}^1 x_1 + \theta_{1,2}^1 x_2 + \theta_{1,0}^1 = +3.0; & s_1^1 &= \frac{1}{1+\exp(-\phi_1^1)} = 0.953 \\
 \phi_2^1 &= \theta_{2,1}^1 x_1 + \theta_{2,2}^1 x_2 + \theta_{2,0}^1 = -3.0; & s_2^1 &= \frac{1}{1+\exp(-\phi_2^1)} = 0.047 \\
 \phi_3^1 &= \theta_{3,1}^1 x_1 + \theta_{3,2}^1 x_2 + \theta_{3,0}^1 = +3.0; & s_3^1 &= \frac{1}{1+\exp(-\phi_3^1)} = 0.953 \\
 \phi_1^2 &= \theta_{1,1}^2 s_1^1 + \theta_{1,2}^2 s_2^1 + \theta_{1,3}^2 s_3^1 + \theta_{1,0}^2 = +2.953; & s_1^2 &= \frac{1}{1+\exp(-\phi_1^2)} = 0.950 \\
 \phi_2^2 &= \theta_{2,1}^2 s_1^1 + \theta_{2,2}^2 s_2^1 + \theta_{2,3}^2 s_3^1 + \theta_{2,0}^2 = -2.953; & s_2^2 &= \frac{1}{1+\exp(-\phi_2^2)} = 0.050
 \end{aligned}$$

- los correspondientes errores en los nodos de la capa de salida y en los nodos de la capa oculta:

$$\begin{aligned}
 \delta_1^2 &= (t_1 - s_1^2) s_1^2 (1 - s_1^2) = +0.0023 \\
 \delta_2^2 &= (t_2 - s_2^2) s_2^2 (1 - s_2^2) = -0.0023 \\
 \delta_1^1 &= (\delta_1^2 \theta_{1,1}^2 + \delta_2^2 \theta_{2,1}^2) s_1^1 (1 - s_1^1) = +0.0002 \\
 \delta_2^1 &= (\delta_1^2 \theta_{1,2}^2 + \delta_2^2 \theta_{2,2}^2) s_2^1 (1 - s_2^1) = +0.0002 \\
 \delta_3^1 &= (\delta_1^2 \theta_{1,3}^2 + \delta_2^2 \theta_{2,3}^2) s_3^1 (1 - s_3^1) = +0.0002
 \end{aligned}$$

- El nuevo peso  $\theta_{2,3}^2$  es:  $\theta_{2,3}^2 = \theta_{2,3}^2 + \rho \delta_2^2 s_3^1 = (-1.0) + (1) (-0.0023) (0.953) = -1.0022$   
El nuevo peso  $\theta_{3,2}^1$  es:  $\theta_{3,2}^1 = \theta_{3,2}^1 + \rho \delta_3^1 x_2 = (+1.0) + (1) (0.0002) (+1.0) = +1.0002$



### Problema 3 (2 puntos; tiempo estimado: 20 minutos)

Considerar la red bayesiana  $\mathcal{R}$  definida como  $P(R, T, X, Y, Z) = P(R) P(X | R) P(Y | R) P(Z | R) P(T | Y)$ , cuyas variables  $R$  y  $T$  toman valores en  $\{1, 2, 3\}$  y las variables  $X, Y, Z$ , en el conjunto  $\{\text{"a"}, \text{"b"}\}$ . Las distribuciones de probabilidad asociadas son como sigue:

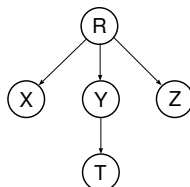
- $P(R)$  es uniforme:  $P(R = 1) = P(R = 2) = P(R = 3)$
- $P(X | R)$ ,  $P(Y | R)$  y  $P(Z | R)$  son idénticas y vienen dadas en la tabla A.
- $P(T | Y)$  viene dada por la tabla B.

A	"a"	"b"
1	2/3	1/3
2	1/4	3/4
3	3/5	2/5

B	1	2	3
"a"	1/3	1/3	1/3
"b"	1/4	2/4	1/4

- Representar gráficamente la red
- Obtener una expresión simplificada de  $P(X, Y, Z | R)$  en función de las distribuciones que definen  $\mathcal{R}$
- ¿Cuál es el mejor valor de  $X$  sabiendo que  $R = 1$ ?

a) Representación gráfica de la red:



b) Expresión simplificada de  $P(X, Y, Z | R)$ :

$$\begin{aligned}
 P(X, Y, Z | R) &= \frac{P(R, X, Y, Z)}{P(R)} \\
 &= \frac{\sum_t P(R, T = t, X, Y, Z)}{P(R)} \\
 &= \frac{\sum_t \cancel{P(R)} P(X | R) P(Y | R) P(Z | R) P(T = t | Y)}{\cancel{P(R)}} \\
 &= P(X | R) P(Y | R) P(Z | R) \sum_t P(T = t | Y) \\
 &= P(X | R) P(Y | R) P(Z | R)
 \end{aligned}$$

c) El mejor valor de  $X$  sabiendo que  $R = 1$  se calcula:

Solución trivial: A partir de la tabla que define  $P(X | R)$ ,  $P(X = \text{"a"} | R = 1) = 2/3$ ,  $P(X = \text{"b"} | R = 1) = 1/3$ , por tanto el mejor valor de  $X$  es "a".

Solución indirecta: A partir del resultado del apartado b)

$$\begin{aligned}
 P(X = \text{"a"} | R = 1) &= \sum_{v, w \in \{\text{"a"}, \text{"b"}\}} P(X = \text{"a"}, Y = v, Z = w | R = 1) \\
 &= \sum_{v, w \in \{\text{"a"}, \text{"b"}\}} P(X = \text{"a"} | R = 1) P(Y = v | R = 1) P(Z = w | R = 1) \\
 &= P(X = \text{"a"} | R = 1) = 2/3 \\
 P(X = \text{"b"} | R = 1) &= 1/3
 \end{aligned}$$

El mejor valor de  $X$  sabiendo que  $R = 1$  es "a"