

Examen de Aprendizaje Automático
ETSINF, Universitat Politècnica de València, 21 de enero de 2020

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 ☐ **D** Para un problema de clasificación en dos clases, sea $\theta(\mathbf{x}; \theta) \stackrel{\text{def}}{=} \theta^t \mathbf{x} + \theta_0$ una función discriminante lineal (FDL) y sea H el hiperplano de decisión definido por $\phi(\mathbf{x}; \theta) = 0$. Entre las siguientes supuestas propiedades hay una que es falsa:
- A) El valor de $\phi(\mathbf{x}; \theta)$ es proporcional a la distancia de \mathbf{x} a H
 - B) La distancia del origen de coordenadas a H es $\frac{\theta_0}{\|\theta\|}$
 - C) H también está definido por un número infinito de FDL $\phi' \neq \phi$
 - D) Solo hay una única FDL que define a H
- 2 ☐ **A** Se ha evaluado un sistema de Aprendizaje Automático mediante un proceso de *exclusion individual* ("Leaving One Out") usando 1000 muestras etiquetadas. En este proceso se han producido 15 errores en total. Indicar cuál de las afirmaciones siguientes es razonable:
- A) La talla de entrenamiento efectiva es de 999 muestras y el error estimado es $1.5 \% \pm 0.75 \%$
 - B) La talla de entrenamiento efectiva es de 1000 muestras y el error estimado es $1.5 \% \pm 0.15 \%$
 - C) La talla de test efectiva es de 999 muestras y el error estimado es $1.5 \% \pm 0.15 \%$
 - D) Las tallas de entrenamiento y de test efectivas son de 1000 muestras y el error estimado es $1.5 \% \pm 0.75 \%$
- 3 ☐ **D** Se desea ajustar por mínimos cuadrados la función $f: \mathbb{R}^2 \rightarrow R$, definida como: $y = f(\mathbf{x}) \stackrel{\text{def}}{=} ax_1x_2 + bx_1 + cx_2$ a una secuencia de N pares entrada-salida: $S = (\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots, (\mathbf{x}_N, y_N)$. La técnica empleada es minimizar por descenso por gradiente la función de error cuadrático:

$$q(a, b, c) = \sum_{n=1}^N (f(\mathbf{x}_n) - y_n)^2$$

Identifica la afirmación acertada de entre las siguientes:

- A) El gradiente es $ax_1 + bx_2 + cx_1x_2$
 - B) El vector gradiente es: $2 \sum_{n=1}^N (f(\mathbf{x}_n) - y_n) \cdot \mathbf{x}_n^t$
 - C) La técnica de descenso por gradiente no es aplicable en este caso ya que la función a ajustar, $f(\cdot)$, no es lineal.
 - D) El vector gradiente es: $2 \sum_{n=1}^N (f(\mathbf{x}_n) - y_n) \cdot (x_{n1}x_{n2}, x_{n1}, x_{n2})^t$
- 4 ☐ **C** Considerar el aprendizaje mediante máquinas de vectores soportes y márgenes blandos con una muestra de aprendizaje $\mathbf{x}_1, \dots, \mathbf{x}_N$ no separable linealmente. Si un multiplicador de Lagrange óptimo α_j^* , asociado a la restricción $c_j (\theta^t \mathbf{x}_j + \theta_0) \geq 1 - \zeta_j$, $1 \leq j \leq N$, es cero, entonces:
- A) La muestra \mathbf{x}_j está mal clasificada
 - B) La muestra \mathbf{x}_j está clasificada correctamente pero θ y θ_0 no es canónico con respecto a la muestra
 - C) La muestra \mathbf{x}_j está clasificada correctamente
 - D) La muestra \mathbf{x}_j es un vector soporte

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

En la siguiente tabla se presenta una muestra de entrenamiento no linealmente separable en \mathbb{R}^2 y los correspondientes multiplicadores de Lagrange óptimos obtenidos al entrenar una máquina de vectores soporte con esta muestra (y $C=10$):

i	1	2	3	4	5	6
x_{i1}	2	4	3	2	1	4
x_{i2}	3	2	2	2	2	1
Clase	+1	-1	+1	-1	+1	-1
α_i^*	0.67	3.78	10.00	10.00	3.11	0.00

- Obtener la función discriminante lineal correspondiente y el valor del margen.
- Representar gráficamente la frontera lineal de separación entre clases, los márgenes y las muestras de entrenamiento, indicando cuáles son vectores soporte.
- Clasificar la muestra $(4, 4)^t$.

a) Pesos de la función discriminante:

$$\theta^* = c_1 \alpha_1^* \mathbf{x}_1 + c_2 \alpha_2^* \mathbf{x}_2 + c_3 \alpha_3^* \mathbf{x}_3 + c_4 \alpha_4^* \mathbf{x}_4 + c_5 \alpha_5^* \mathbf{x}_5$$

$$\theta_1^* = (+1)(0.67)(2) + (-1)(3.78)(4) + (+1)(10)(3) + (-1)(10)(2) + (+1)(3.11)(1) \approx -0.67$$

$$\theta_2^* = (+1)(0.67)(3) + (-1)(3.78)(2) + (+1)(10)(2) + (-1)(10)(2) + (+1)(3.11)(2) \approx +0.67$$

Usando el vector soporte \mathbf{x}_2 (que verifica la condición : $0 < \alpha_2^* < C = 10$)

$$\theta_0^* = c_2 - \theta^{*t} \mathbf{x}_2 = -1 - ((-0.67)(4) + (0.67)(2)) = 0.34$$

$$\text{Margen: } \frac{2}{\|\theta^*\|} \approx 2.11$$

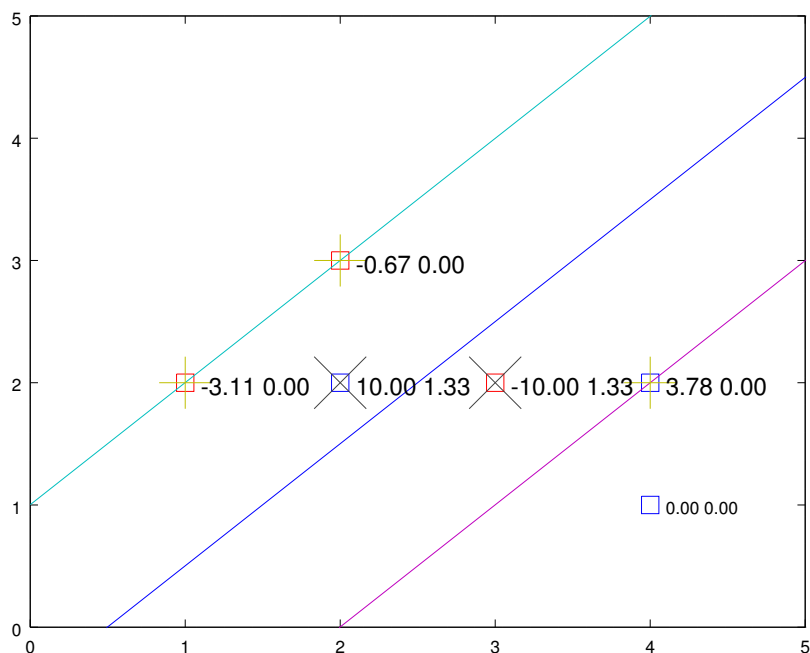
b) Frontera de separación y representación gráfica:

Ecuación de la frontera lineal de separación: $-0.67 x_1 + 0.67 x_2 + 0.34 = 0$

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

El margen lo definen las dos rectas paralelas a la frontera de separación, cada una de ellas situada a una distancia de $2.11/2 \approx 1.05$ y cuyas ecuaciones son: $-0.67 x_1 + 0.67 x_2 + 0.34 = +1$ y $-0.67 x_1 + 0.67 x_2 + 0.34 = -1$

Representación gráfica:



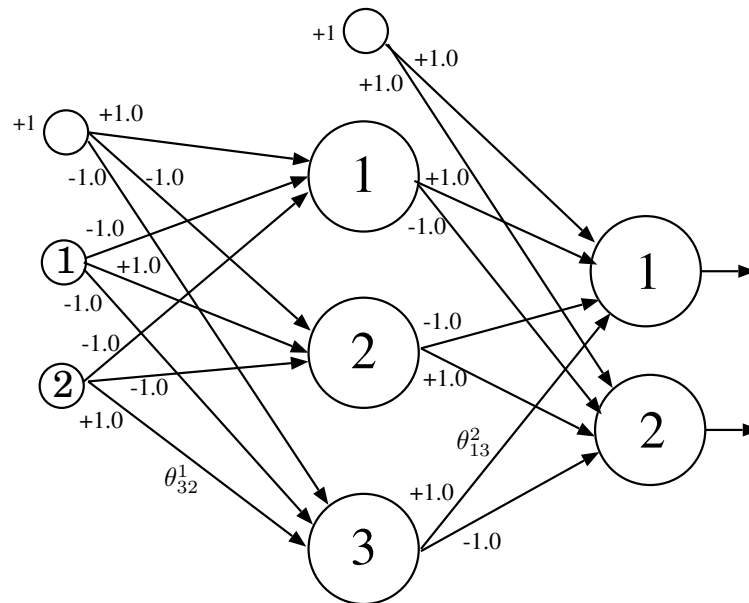
c) Clasificación de la muestra $(4, 4)^t$:

El valor de la función discriminante para este vector es:

$$\theta_0^* + \theta_1^* x_1 + \theta_2^* x_2 = +0.34 + (-0.67) * 4 + (0.67) * 4 = +0.34 > 0 \Rightarrow \text{clase } +1.$$

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 de tipo lineal y de la capa oculta de tipo *sigmoide*, y factor de aprendizaje $\rho = 1.0$.



Dado un vector de entrada $\mathbf{x}^t = (+2, +2)$, las salidas de las unidades de la capa de salida son $s_1^2 = 1.0474$ y $s_2^2 = 0.9526$ y las de las unidades ocultas son $s_1^1 = 0.0474$, $s_2^1 = 0.2689$ y $s_3^1 = 0.2689$. Si el valor deseado de salida es $\mathbf{t}^t = (+1, 0)$, calcular:

- 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 θ_{32}^1 y θ_{13}^2 .

- Los errores en la capa de salida (función de activación lineal) son:

$$\delta_1^2 = (t_1 - s_1^2) = -0.0474$$

$$\delta_2^2 = (t_2 - s_2^2) = -0.9526$$

Los errores en la capa de oculta son:

$$\delta_1^1 = (\delta_1^2 \theta_{11}^2 + \delta_2^2 \theta_{21}^2) s_1^1 (1 - s_1^1) = +0.0409;$$

$$\delta_2^1 = (\delta_1^2 \theta_{12}^2 + \delta_2^2 \theta_{22}^2) s_2^1 (1 - s_2^1) = -0.1780;$$

$$\delta_3^1 = (\delta_1^2 \theta_{13}^2 + \delta_2^2 \theta_{23}^2) s_3^1 (1 - s_3^1) = +0.1780$$

- El nuevo peso θ_{13}^2 es: $\theta_{13}^2 = \theta_{13}^2 + \rho \delta_1^2 s_3^1 = 0.9872$;
El nuevo peso θ_{32}^1 es: $\theta_{32}^1 = \theta_{32}^1 + \rho \delta_3^1 x_2 = 1.3559$

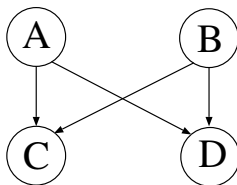
Problema 3 (2 puntos; tiempo estimado: 20 minutos)

Las variables aleatorias A, B, C, D toman valores en el conjunto $\{0, 1\}$. La distribución de probabilidad conjunta de estas variables viene dada por $P(A, B, C, D) = P(A) P(B) P(C | A, B) P(D | A, B)$, donde las distribuciones de probabilidad asociadas son:

$$\begin{aligned} P(A = 1) &= 0.3 & P(A = 0) &= 0.7 \\ P(B = 1) &= 0.4 & P(B = 0) &= 0.6 \\ P(C = 1 | A = 0, B = 0) &= 0.1 & P(C = 0 | A = 0, B = 0) &= 0.9 \\ P(C = 1 | A = 0, B = 1) &= 0.2 & P(C = 0 | A = 0, B = 1) &= 0.8 \\ P(C = 1 | A = 1, B = 0) &= 0.3 & P(C = 0 | A = 1, B = 0) &= 0.7 \\ P(C = 1 | A = 1, B = 1) &= 0.4 & P(C = 0 | A = 1, B = 1) &= 0.6 \\ P(D = 1 | A = 0, B = 0) &= 0.9 & P(D = 0 | A = 0, B = 0) &= 0.1 \\ P(D = 1 | A = 0, B = 1) &= 0.8 & P(D = 0 | A = 0, B = 1) &= 0.2 \\ P(D = 1 | A = 1, B = 0) &= 0.7 & P(D = 0 | A = 1, B = 0) &= 0.3 \\ P(D = 1 | A = 1, B = 1) &= 0.6 & P(D = 0 | A = 1, B = 1) &= 0.4 \end{aligned}$$

- Representar gráficamente la red bayesiana correspondiente
- Obtener una expresión simplificada de $P(A | B, C, D)$ y calcular su valor para $A = 1$ cuando $B = 1, C = 1$ y $D = 1$.
- Dados $B = 1, C = 1$ y $D = 1$, ¿Cuál es el mejor valor de A que se puede predecir?

a) Representar gráficamente la red bayesiana correspondiente



- Obtener una expresión simplificada de $P(A | B, C, D)$ y calcular su valor para $A = 1$ cuando $B = 1, C = 1$ y $D = 1$.

$$\begin{aligned} P(A | B, C, D) &= \frac{P(A, B, C, D)}{P(B, C, D)} = \frac{P(A) \cancel{P(B)} P(C | A, B) P(D | A, B)}{\cancel{P(B)} \sum_a P(A = a) P(C | A = a, B) P(D | A = a, B)} \\ &= \frac{P(A) P(C | A, B) P(D | A, B)}{P(A = 0) P(C | A = 0, B) P(D | A = 0, B) + P(A = 1) P(C | A = 1, B) P(D | A = 1, B)} \\ P(A = 1 | B = 1, C = 1, D = 1) &= \frac{0.3 \cdot 0.4 \cdot 0.6}{0.7 \cdot 0.2 \cdot 0.8 + 0.3 \cdot 0.4 \cdot 0.6} = 0.391 \end{aligned}$$

- Dados $B = 1, C = 1$ y $D = 1$, ¿Cuál es el mejor valor de A que se puede predecir?

$$a^* = \arg \max_{a \in \{0, 1\}} P(A = a | B = 1, C = 1, D = 1)$$

$$P(A = 0 | B = 1, C = 1, D = 1) = 1 - 0.391 = 0.609 \geq 0.391 \Rightarrow \text{valor óptimo de } A \text{ es } a^* = 0$$