

Examen de Aprendizaje Automático
ETSINF, Universitat Politècnica de València, 11 de enero de 2017

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.

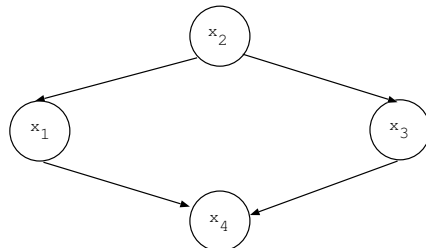
- 1 ☐ C Sea S un conjunto de datos supervisados o etiquetados. Para el diseño de un sistema de reconocimiento de formas, se utilizan datos de S tanto para aprender los parámetros del modelo de reconocimiento, \mathcal{M} , como para estimar la probabilidad de error de reconocimiento esperada para dicho modelo, p_e . Indicar cual de las siguientes afirmaciones es incorrecta.
- A) Si S es suficientemente grande, el método de *validación cruzada en B bloques* puede proporcionar buenas estimaciones de p_e , basadas en todos los datos de S . Una vez estimado p_e , también es recomendable usar todos los datos de S para el aprendizaje final de \mathcal{M} .
 - B) Si la talla de S es 160, y se desea que el intervalo de confianza al 95 % de p_e sea menor que $\pm 1\%$, el método de *partición* sería totalmente inapropiado.
 - C) Si S es suficientemente grande, se puede elegir un valor adecuado de B para que el método de *validación cruzada en B bloques* garantice un entrenamiento de \mathcal{M} que evite tanto el sobreajuste como el sobreentrenamiento.
 - D) Si se usa el método de *exclusión individual* con un conjunto S cuya talla es menor de 100 y se obtiene $p_e = 0.1$, el intervalo de confianza al 95 % de esta estimación será mayor que $\pm 5\%$.
- 2 ☐ D En el problema de optimización con restricciones

$$\begin{array}{ll} \text{minimizar} & q(\Theta), \quad \Theta \in \mathbb{R}^D \\ \text{sujeito a} & v_i(\Theta) \leq 0, \quad 1 \leq i \leq k \\ & u_i(\Theta) = 0, \quad 1 \leq i \leq m \end{array}$$

se cumplen las condiciones complementarias de Karush-Kuhn-Tucker $\alpha_i^* v_i(\Theta^*) = 0$ para $1 \leq i \leq k$. Indicar cuál de las siguientes afirmaciones se deduce de ellas:

- A) Si para un i , $\alpha_i^* < 0$, entonces $v_i(\Theta^*) > 0$
 - B) Si para un i , $u_i(\Theta^*) = 0$, entonces $v_i(\Theta^*) \geq 0$
 - C) Si para un i , $u_i(\Theta^*) = 0$, entonces $\alpha_i^* < 0$,
 - D) Si para un i , $\alpha_i^* > 0$, entonces $v_i(\Theta^*) = 0$
- 3 ☐ C Las siguientes afirmaciones se refieren al método Esperanza Maximización (EM) aplicado a una muestra de entrenamiento S . Identificar cuál de ellas es errónea o inapropiada:
- A) EM es útil para estimar valores maximo-verosímiles de los parámetros de modelos estadísticos a partir de S .
 - B) EM es un método iterativo que garantiza la convergencia a un máximo local de la verosimilitud de S .
 - C) La rapidez de convergencia de EM puede mejorarse eligiendo un factor de aprendizaje adecuado para S .
 - D) La rapidez de convergencia de EM puede mejorarse inicializando los parámetros de forma adecuada para S .

- 4 ☐ D En la red bayesiana



¿cuál de las relaciones siguientes es falsa en general?

- A) $P(x_2, x_4 \mid x_3) = P(x_2 \mid x_3) P(x_4 \mid x_3)$
- B) $P(x_1, x_3 \mid x_2) = P(x_1 \mid x_2) P(x_3 \mid x_2)$
- C) $P(x_1, x_3) = P(x_1) P(x_3)$
- D) $P(x_1, x_3 \mid x_4) = P(x_1 \mid x_4) P(x_3 \mid x_4)$

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

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

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

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

- Obtener la función discriminante lineal correspondiente
- Obtener la ecuación de la frontera de decisión entre clases y representar gráficamente los puntos de entrenamiento y dicha frontera.
- Calcular el margen óptimo
- Clasificar la muestra $(3, 1)$.

a) Función discriminante lineal (FDL):

- Vector de pesos:

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

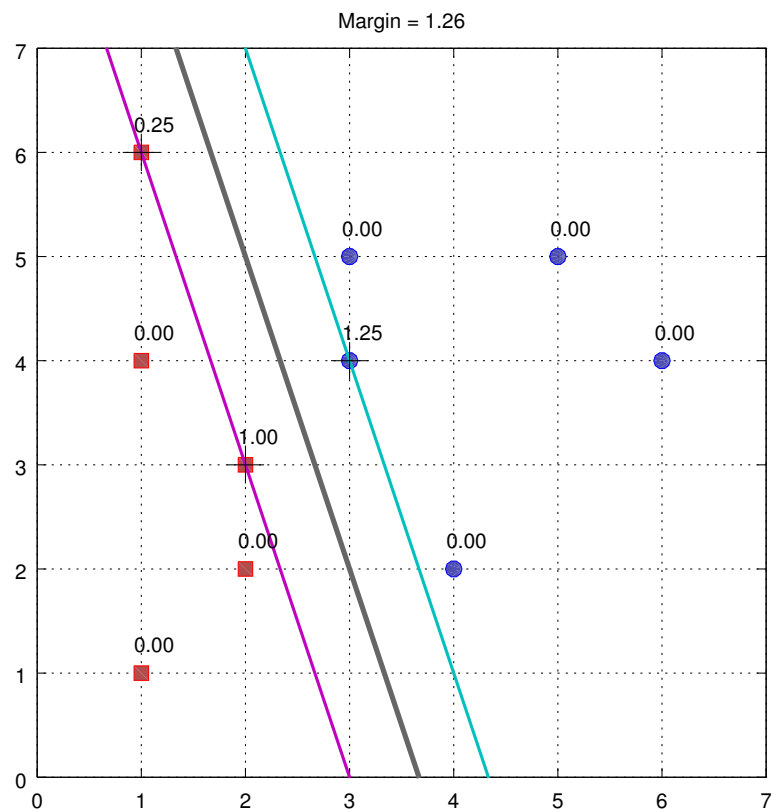
$$\theta_2^* = +1 \cdot 0.25 \cdot 6 + 1 \cdot 1.0 \cdot 3 - 1 \cdot 1.25 \cdot 4 = -0.5$$

- Peso umbral (con el tercer vector de entrenamiento): $\theta_0^* = (+1) - (-1.5 \cdot 1 - 0.5 \cdot 6) = 5.5$
- FDL: $\phi(\mathbf{x}) = -1.5 x_1 - 0.5 x_2 + 5.5$

b) Ecuación de la frontera de decisión:

$$-1.5 x_1 - 0.5 x_2 + 5.5 = 0 \Rightarrow x_2 = -3x_1 + 11$$

Representación gráfica:



c) Margen óptimo:

$$\frac{2}{\|\theta^*\|} = \frac{2}{\sqrt{(-1.5)^2 + (-0.5)^2}} = 1.265$$

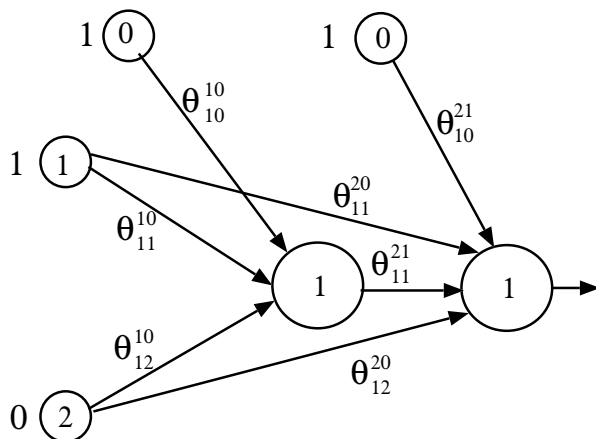
Alternativamente:

$$2 \left(\sum_{n \in \mathcal{V}} \alpha_n^* \right)^{-1/2} = \frac{2}{\sqrt{0.25 + 1.0 + 1.25}} = 1.265$$

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

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

La red hacia adelante (“feedforward”) 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 $\rho = 1.0$.



Dados unos pesos iniciales $\theta_{10}^{10} = \theta_{11}^{10} = \theta_{12}^{10} = \theta_{11}^{20} = \theta_{12}^{20} = \theta_{10}^{21} = \theta_{11}^{21} = 1.0$, un vector de entrada $\mathbf{x}^t = (1, 0)$ y su valor deseado de salida $t = +1$, Calcular:

- las salidas de todas las unidades
- los correspondientes errores en el nodo de la capa de salida y en el de la capa oculta.
- Los nuevos valores de los pesos de las conexiones

Pista: La actualización de pesos en esta red sigue la misma formulación que en el BackProp para el perceptrón multicapa convencional: el incremento de peso es $\Delta\theta = \rho z \delta$, donde ρ es el factor de aprendizaje, z es la entrada del arco asociado al peso θ , y δ es el error que se observa en la salida de la unidad a la que llega ese arco, multiplicado por la derivada de la función de activación.

- Las salidas de todas las unidades

$$\phi_1^1 = \theta_{10}^{10} + \theta_{11}^{10} x_1 + \theta_{12}^{10} x_2 = 2.0$$

$$s_1^1 = \frac{1}{1 + \exp(-\phi_1^1)} = .880797$$

$$\phi_1^2 = \theta_{10}^{21} + \theta_{11}^{20} x_1 + \theta_{12}^{20} x_2 + \theta_{11}^{21} s_1^1 = 2.880797$$

$$s_1^2 = \frac{1}{1 + \exp(-\phi_1^2)} = .946889$$

- El error en la capa de salida es:

$$\delta_1^2 = (t_1 - s_1^2) s_1^2 (1 - s_1^2) = .002671$$

El error en la capa de oculta es:

$$\delta_1^1 = (\delta_1^2 \theta_{11}^{21}) s_1^1 (1 - s_1^1) = .000280$$

- Los nuevos pesos son:

$$\theta_{10}^{21} = \theta_{10}^{21} + \rho \delta_1^2 (+1) = \theta_{10}^{21} + 0.002671 = 1.002671$$

$$\theta_{11}^{21} = \theta_{11}^{21} + \rho \delta_1^2 s_1^2 = \theta_{11}^{21} + 0.002529 = 1.002529$$

$$\theta_{11}^{20} = \theta_{11}^{20} + \rho \delta_1^2 x_1 = \theta_{11}^{20} + 0.002671 = 1.002671$$

$$\theta_{12}^{20} = \theta_{12}^{20} + \rho \delta_1^2 x_2 = \theta_{12}^{20} + 0.0 = 1.0$$

$$\theta_{10}^{10} = \theta_{10}^{10} + \rho \delta_1^1 (+1) = \theta_{10}^{10} + 0.000280 = 1.000280$$

$$\theta_{11}^{10} = \theta_{11}^{10} + \rho \delta_1^1 x_1 = \theta_{11}^{10} + 0.000280 = 1.000280$$

$$\theta_{12}^{10} = \theta_{12}^{10} + \rho \delta_1^1 x_2 = \theta_{12}^{10} + 0.0 = 1.0$$

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

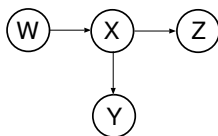
Considerar la red bayesiana \mathcal{R} definida como $P(W, X, Y, Z) = P(W) P(X | W) P(Y | X) P(Z | X)$, cuyas variables aleatorias, W, X, Y, Z , toman valores en el conjunto $\{a, b, c\}$. Las distribuciones de probabilidad asociadas son como sigue:

- $P(W)$ es uniforme: $P(W = a) = P(W = b) = P(W = c)$,
- $P(X | W)$, $P(Y | X)$ y $P(Z | X)$ vienen dadas en las siguientes tablas:

$P(x w)$	$x:$	a	b	c	$P(y x)$	$y:$	a	b	c	$P(z x)$	$z:$	a	b	c
$w: a$		1/2	0	1/2	$x: a$		1/3	0	2/3	$x: a$		1/3	0	2/3
b		1/4	1/2	1/4	b		1/4	1/2	1/4	b		1/4	1/2	1/4
c		1/5	3/5	1/5	c		0	3/5	2/5	c		0	3/5	2/5

- Representar gráficamente la red
- Obtener una expresión simplificada de $P(X, Y, Z | W)$ en función de las distribuciones que definen \mathcal{W} y calcular $P(X = b, Y = b, Z = b | W = b)$
- Obtener una expresión simplificada de $P(W | X, Y, Z)$, y calcular $P(W = c | X = b, Y = b, Z = b)$

- Representación gráfica de la red:



- Expresión simplificada de $P(X, Y, Z | W)$:

$$\begin{aligned}
 P(X, Y, Z | W) &= \frac{P(W, X, Y, Z)}{P(W)} = \frac{\cancel{P(W)} P(X | W) P(Y | X) P(Z | X)}{\cancel{P(W)}} \\
 &= P(X | W) P(Y | X) P(Z | X)
 \end{aligned}$$

$$P(X = b, Y = b, Z = b | W = b) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8}$$

- Expresión simplificada de $P(W | X, Y, Z)$:

$$\begin{aligned}
 P(W | X, Y, Z) &= \frac{P(W, X, Y, Z)}{P(X, Y, Z)} = \frac{P(W) P(X | W) \cancel{P(Y | X)} \cancel{P(Z | X)}}{\sum_{w \in \{a, b, c\}} P(W = w) P(X | W = w) \cancel{P(Y | X)} \cancel{P(Z | X)}} \\
 &= \frac{\cancel{(1/3)} P(X | W)}{\cancel{(1/3)} \sum_{w \in \{a, b, c\}} P(X | W = w)} = \frac{P(X | W)}{\sum_{w \in \{a, b, c\}} P(X | W = w)}
 \end{aligned}$$

$$P(W = c | X = b, Y = b, Z = b) = \frac{P(X = b | W = c)}{\sum_{w \in \{a, b, c\}} P(X = b | W = w)} = \frac{3/5}{0 + 1/2 + 3/5} = \frac{6}{11}$$