

# 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.

# Examen de los temas 5 a 7 de Aprendizaje Automático

ETSINF, Universitat Politècnica de València, 16 de enero de 2014

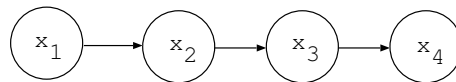
Apellidos:

Nombre:

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

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

- 1 ☒ A Si la función de acción no lineal en la capa de salida de un perceptrón de dos capas fuese lineal, las fórmulas que permiten modificar los pesos de dicha capa de salida en el algoritmo BackProp verifican que (solo una respuesta es correcta):
- A)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) s_j^1$   
B)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) g(\phi_i^2) (1 - g(\phi_i^2)) s_j^1$   
C)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) g(\phi_i^2) s_j^1$   
D)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) (1 - g(\phi_i^2)) s_j^1$
- 2 ☒ B En una presentación de las  $N$  muestras de aprendizaje mediante el algoritmo de retropropagación del error, indicar qué afirmación es la correcta:
- A) Los pesos se modifican una sola vez tanto en la versión “online” o incremental como en la versión batch.  
B) Los pesos se modifican  $N$  veces en la versión “online” o incremental y una en la versión batch.  
C) Los pesos se modifican  $N$  veces en la versión “online” y también en la versión batch.  
D) Los pesos no se modifican en la versión “online” o incremental y una en la versión batch.
- 3 ☒ B En el perceptrón multicapa, la parálisis de la red se produce cuando (solo una respuesta es correcta)
- A) En test, los pesos son muy grandes  
B) En entrenamiento, los valores de las combinaciones lineales de los nodos son muy grandes  
C) En test, los valores de las combinaciones lineales de los nodos son muy grandes  
D) En entrenamiento, los pesos son nulos
- 4 ☒ A En la red bayesiana lineal



¿cuál de las relaciones siguientes es correcta?

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

# Examen de recuperación de Aprendizaje Automático

ETSINF, Universitat Politècnica de València, 28 de enero de 2014

Apellidos:

Nombre:

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

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

- 1 ☐ B Se ha evaluado un sistema de Aprendizaje Automático mediante la técnica de *validación cruzada en B bloques* (“B-fold Cross Validation”) con  $B = 10$  y utilizando un conjunto de datos de entrenamiento que contiene 1000 muestras. Se han obtenido un total de 20 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 900 muestras y el error es del 2 %
- C) La talla de entrenamiento efectiva es de 900 muestras y el error es del 20 %
- D) La talla de entrenamiento efectiva es de 1000 muestras y el error es del 20 %.

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

$$q_S(\theta) = \sum_{n=1}^N (\theta^t x_n - y_n) + \frac{1}{2} \theta^t \theta,$$

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

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

- 3 ☐ A Un clasificador implementado mediante una red neuronal con  $L$  capas ocultas en la que todas las funciones de activación son lineales, es equivalente a:

- A) un clasificador basado en funciones discriminantes lineales
- B) un clasificador basado en funciones discriminantes lineales generalizadas cuyas fronteras de decisión son no-lineales
- C) un clasificador implementado mediante una red neuronal con  $L - 1$  capas ocultas
- D) para que la red pueda usarse para clasificación, al menos las funciones de activación de la capa de salida han de ser no-lineales.

- 4 ☐ A Sea  $\mathcal{C}$  un conjunto de variables aleatorias (VA). Un concepto importante en el que se basan las técnicas de redes bayesianas es:

- A) todas las probabilidades condicionales e incondicionales en las que participan VA's de  $\mathcal{C}$  se pueden obtener mediante las reglas básicas de inferencia estadística, a partir de la probabilidad conjunta de todas las VA de  $\mathcal{C}$ .
- B) la probabilidad incondicional de una VA  $a \in \mathcal{C}$  solo depende de las VA's de los nodos con los que está conectado el nodo de  $a$ .
- C) el grafo que representa las VA's de  $\mathcal{C}$  ha de ser acíclico y conexo.
- D) ha de existir independencia condicional entre al menos dos VA's de  $\mathcal{C}$ .

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 16 de enero de 2015

Apellidos:

Nombre:

**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 Considerar la siguiente modificación de la función de Widrow y Hoff

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

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

- 2 ☐ C En el problema de optimización con restricciones

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

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

- A) Existe un  $i$  tal que  $\alpha_i^* < 0$  y  $v_i(\boldsymbol{\Theta}^*) = 0$   
B) Para todo  $i$ , si  $\alpha_i^* = 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$ ,  
C) Si para un  $i$ ,  $\alpha_i^* > 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$   
D) Existe un  $i$  tal que  $v_i(\boldsymbol{\Theta}^*) > 0$  y  $\alpha_i^* = 0$

- 3 ☐ B Las siguientes afirmaciones se refieren a la estimación por máxima verosimilitud de los parámetros de una mezcla de  $K$  gaussianas (vector-media y peso de cada gaussiana) mediante un conjunto de vectores de entrenamiento cualquiera de dimensión  $D$ . Identifica cuál es *falsa*.

- A) Los parámetros de la mezcla se estiman adecuadamente mediante un algoritmo de *esperanza maximización* (EM)  
B) El algoritmo EM obtiene los valores óptimos de los parámetros a estimar  
C) La verosimilitud del conjunto de entrenamiento, calculada con los parámetros estimados, aumenta en cada iteración del EM  
D) En cada iteración, el algoritmo EM estima los valores de las variables ocultas que, en este caso, son los pesos de las gaussianas.

- 4 ☐ A Sea  $\mathcal{C}$  un conjunto de variables aleatorias. Un concepto importante en el que se basan las técnicas de redes bayesianas es:

- A) el grafo que relaciona a las variables entre si define una distribución de probabilidad conjunta en las variables  $\mathcal{C}$  y permite calcular cualquier probabilidad condicional en la que intervengan variables de  $\mathcal{C}$   
B) los nodos del grafo representan las dependencias entre las variables en  $\mathcal{C}$   
C) el grafo que relaciona a las variables entre si define una distribución de probabilidad condicional entre dos subconjuntos de variables en  $\mathcal{C}$   
D) las probabilidades condicionales se calculan a partir de los cliques (subgrafos completos) que contiene el grafo.

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 28 de enero de 2015

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 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 de entrenamiento que contiene 1000 muestras. Se han obtenido un total de 22 errores. Indicar cuál de las afirmaciones siguientes es razonable:

- A) La talla de entrenamiento efectiva es 990 muestras y el error estimado es  $2.2 \% \pm 0.1 \%$
- B) La talla de entrenamiento efectiva es de 900 muestras y el error estimado es  $2.2 \%$
- C) La talla de test efectiva es de 1000 muestras y el error estimado es  $2.2 \% \pm 0.7 \%$
- D) El error estimado es  $22 \% \pm 7 \%$ .

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

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

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

- 3 ☐ A 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  $\alpha_j^*$ , asociado a la restricción  $c_j (\boldsymbol{\theta}^t \mathbf{x}_{dj} n + \theta_0) \geq 1 - \zeta_j$ ,  $1 \leq j \leq N$ , es cero, entonces:

- A) La muestra  $\mathbf{x}_j$  está clasificada correctamente.
- B) La muestra  $\mathbf{x}_j$  está mal clasificada.
- C) La muestra  $\mathbf{x}_j$  está clasificada correctamente pero  $\boldsymbol{\theta}$  y  $\theta_0$  no es canónico con respecto a la muestra.
- D) La muestra  $\mathbf{x}_j$  es un vector soporte.

- 4 ☐ A La distribución de probabilidad conjunta en una red bayesiana de tres nodos  $A, B$  y  $C$  es  $P(A, B, C) = P(C) P(A | C) P(B | C)$ . Marcar cuál es la afirmación correcta:

- A)  $P(A, B | C) = P(A | C) P(B | C)$
- B) En general,  $P(A, B | C) \neq P(A | C) P(B | C)$
- C)  $P(A, B | C) = P(C | A) P(C | B)$
- D)  $P(A, B | C) = P(A) P(B)$

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 11 de enero de 2016

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** Se ha evaluado un sistema de Aprendizaje Automático mediante la técnica de *Exclusión individual* (“Leaving One Out”) y utilizando un conjunto de datos que contiene 200 muestras. Se han obtenido un total de 10 errores. Indicar cuál de las afirmaciones siguientes es razonable:

- A) La talla de entrenamiento efectiva es 190 muestras, la del test es de 10 muestras y el error estimado es  $5.0\% \pm 0.3\%$
- B) La talla de entrenamiento efectiva es de 199 muestras, la del test es de 1 muestra y el error estimado es  $5.0 \pm 3.0\%$
- C) La talla de entrenamiento efectiva es de 200 muestras, la del test es de 10 muestras y el error estimado es  $5.0 \pm 0.3\%$
- D) La talla de entrenamiento efectiva es de 199 muestras, la del test es de 200 muestras y el error estimado es  $5.0 \pm 3.0\%$

- 2 ☐ **D** En el problema de optimización con restricciones

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

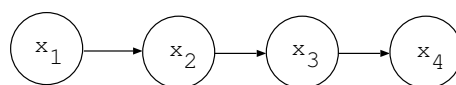
se cumplen las condiciones complementarias de Karush-Kuhn-Tucker  $\alpha_i^* v_i(\Theta^*) = 0$  para  $1 \leq i \leq k$ . Indicar cuál de las siguientes afirmaciones se deduce de ellas:

- A) Si para un  $i$ ,  $\alpha_i^* = 0$ , entonces  $v_i(\Theta^*) > 0$
- B) Si para un  $i$ ,  $\alpha_i^* = 0$ , entonces  $v_i(\Theta^*) = 0$ ,
- C) Si para un  $i$ ,  $v_i(\Theta^*) = 0$ , entonces  $\alpha_i^* = 0$
- D) Si para un  $i$ ,  $\alpha_i^* > 0$ , entonces  $v_i(\Theta^*) = 0$

- 3 ☐ **C** En la estimación por máxima verosimilitud de los parámetros de una mezcla de  $K$  gaussianas de matriz de covarianza común y conocida a partir de  $N$  vectores de entrenamiento, los parámetros a estimar son: el vector-media  $\mu_k$  y el peso  $\pi_k$  de cada gaussiana,  $k, 1 \leq k \leq K$ . Identificar cuál de las siguientes afirmaciones es *correcta*:

- A) Se puede usar *descenso por gradiente*, ya que los valores de  $\mu_k$  no están sujetos a ninguna restricción, lo que hace innecesario recurrir a la técnica de los *multiplicadores de Lagrange*.
- B) La solución se obtiene en un paso, utilizando directamente la *optimización lagrangiana* de la verosimilitud de los  $N$  vectores de entrenamiento. En este caso, hay un único multiplicador de Lagrange,  $\beta$ , asociado a la restricción de igualdad:  $\sum_{k=1}^K \pi_k = 1$ .
- C) El método más adecuado es el de *esperanza-maximización* (EM), el cual garantiza que se cumple la restricción  $\sum_{k=1}^K \pi_k = 1$ . Esto es así gracias a que, en cada iteración de EM, los valores de  $\pi_k, 1 \leq k \leq K$ , se obtienen como medias de valores de variables latentes, usando una expresión que se deriva analíticamente mediante la técnica de los *multiplicadores de Lagrange* con la restricción indicada.
- D) El método más adecuado sería el de *esperanza-maximización* (EM), pero no es posible utilizarlo ya que EM es un método iterativo que no garantiza el cumplimiento de la restricción de igualdad:  $\sum_{k=1}^K \pi_k = 1$ .

- 4 ☐ **B** En la red bayesiana lineal



¿cuál de las relaciones siguientes es falsa en general?

- A)  $P(x_1, x_4 | x_2) = P(x_1 | x_2) P(x_4 | x_2)$
- B)  $P(x_1, x_4 | x_2) = P(x_1) P(x_4)$
- C)  $P(x_1, x_4 | x_2) = P(x_1 | x_2) P(x_4 | x_1, x_2)$
- D)  $P(x_1, x_4 | x_2) = P(x_4 | x_2) P(x_1 | x_4, x_2)$

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 27 de enero de 2016

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 Identifica cuál de las siguientes afirmaciones es *errónea o impropia*:
- A) La *teoría de la decisión estadística* es idónea para problemas de *clasificación*, pero es inadecuada para problemas de *regresión* en aprendizaje de funciones.
  - B) *Clasificación* puede considerarse como un caso particular de *regresión* y los problemas de *sobreajuste* y *sobregeneralización* le afectan de forma similar
  - C) En el *aprendizaje activo* se asume que un agente externo al sistema se encarga de etiquetar datos de entrenamiento seleccionados por el sistema
  - D) El *aprendizaje adaptativo* es esencialmente un modo *supervisado* de aprendizaje
- 2 ☐ A Se ha evaluado un sistema de Aprendizaje Automático mediante un proceso de *validación cruzada en B bloques* ("B-fold Cross Validation") con  $B = 8$  y 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 875 muestras y el error estimado es  $1.5\% \pm 0.75\%$
  - B) La talla de entrenamiento efectiva es de 992 muestras y el error estimado es inferior al 2%
  - C) La talla de test efectiva es de 125 muestras y el error estimado es  $12\% \pm 5.7\%$ .
  - D) Como en cada bloque de test solo hay 125 muestras, el error es muy variable, pudiéndose estimar como  $1.5\% \pm 5.7\%$ .

- 3 ☐ A Se desea ajustar por mínimos cuadrados la función  $f: \mathbb{R}^2 \rightarrow R$ , definida como:  $y = f(\mathbf{x}) \stackrel{\text{def}}{=} ax_1^2 + bx_2^2 + cx_1x_2$  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 vector gradiente es:  $2 \sum_{n=1}^N (f(\mathbf{x}_n) - y_n) \cdot (x_{n1}^2, x_{n2}^2, x_{n1}x_{n2})^t$
  - B) El gradiente es  $2ax_1 + 2bx_2 + cx_1x_2$
  - C) El descenso por gradiente solo es aplicable a funciones convexas, pero  $q(\cdot)$  no lo es.
  - D) La técnica de descenso por gradiente no es aplicable en este caso ya que la función a ajustar,  $f(\cdot)$ , no es lineal.
- 4 ☐ A 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  $\alpha_j^*$ , asociado a la restricción  $c_j (\boldsymbol{\theta}^t \mathbf{x}_{djn} + \theta_0) \geq 1 - \zeta_j$ ,  $1 \leq j \leq N$ , es cero, entonces:
- A) La muestra  $\mathbf{x}_j$  está clasificada correctamente.
  - B) La muestra  $\mathbf{x}_j$  está mal clasificada.
  - C) La muestra  $\mathbf{x}_j$  está clasificada correctamente pero  $\boldsymbol{\theta}$  y  $\theta_0$  no es canónico con respecto a la muestra.
  - D) La muestra  $\mathbf{x}_j$  es un vector soporte.

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.

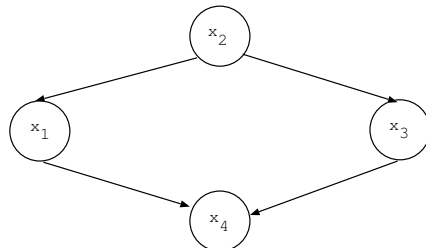
- 1 ☐ C Sea  $S$  un conjunto de datos supervisados o etiquetados. Para el diseño de un sistema de reconocimiento de formas, se utilizan datos de  $S$  tanto para aprender los parámetros del modelo de reconocimiento,  $\mathcal{M}$ , como para estimar la probabilidad de error de reconocimiento esperada 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)  $P(x_2, x_4 | x_3) = P(x_2 | x_3) P(x_4 | x_3)$
- B)  $P(x_1, x_3 | x_2) = P(x_1 | x_2) P(x_3 | x_2)$
- C)  $P(x_1, x_3) = P(x_1) P(x_3)$
- D)  $P(x_1, x_3 | x_4) = P(x_1 | x_4) P(x_3 | x_4)$



Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 25 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 ☐ Sea  $S$  un conjunto de 1000 datos supervisados o etiquetados. En el diseño de un sistema de reconocimiento de formas se utiliza el método de *validación cruzada en 10 bloques* y se obtienen los errores siguientes: (2,4,2,3,0,7,4,4,1). Indicar cual de las siguientes afirmaciones es incorrecta.

- A) El error estimado es  $p_e = 3.1\%$
- B) El intervalo de confianza al 95 % es  $\pm 0.9$
- C) El test efectivo es de 100 muestras
- D) El tamaño de entrenamiento efectivo es de 900 muestras

- 2 ☐ En el problema de aprendizaje de modelos probabilísticos con variables observables  $\mathbf{x}_n$  y latentes  $\mathbf{z}_n$

$$L_S(\Theta) = \sum_{n=1}^N \log \left( \sum_{\mathbf{z}_n} P(\mathbf{x}_n, \mathbf{z}_n | \Theta) \right)$$

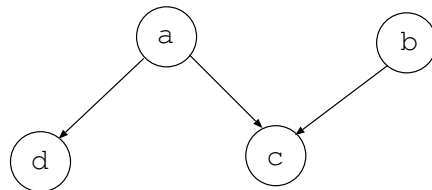
se utiliza la técnica esperanza-maximización (EM). Indicar cuál de las siguientes afirmaciones es cierta:

- A) En cada iteración, el paso E consiste en obtener una estimación de todas las variables  $\mathbf{x}_n$  y  $\mathbf{z}_n$ , y el paso M, obtener los parámetros óptimos de  $\Theta$  utilizando la estimación de las variables  $\mathbf{x}_n$  y  $\mathbf{z}_n$  obtenidas en el paso E.
- B) En cada iteración, el paso E consiste en obtener los valores de las variables latentes  $\mathbf{z}_n$  que maximizan la función objetivo  $L_S(\Theta)$ , y el paso M, obtener los parámetros óptimos de  $\Theta$  utilizando la estimación de las variables  $\mathbf{z}_n$  obtenidas en el paso E.
- C) En cada iteración, el paso E consiste en obtener los valores de todas las variables  $\mathbf{x}_n$  y  $\mathbf{z}_n$  que maximizan la función  $L_S(\Theta)$ , y el paso M, obtener los parámetros óptimos de  $\Theta$  utilizando la estimación de las variables  $\mathbf{x}_n$  y  $\mathbf{z}_n$  obtenidas en el paso E.
- D) En cada iteración, el paso E consiste en obtener una estimación de las variables latentes  $\mathbf{z}_n$ , y el paso M, obtener los parámetros óptimos de  $\Theta$  utilizando la estimación de las variables  $\mathbf{z}_n$  obtenidas en el paso E.

- 3 ☐ 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  $\alpha_j^*$ , 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 la muestra  $\mathbf{x}_j$  está bien clasificada pero ¿cuál de las siguientes afirmaciones es falsa?:

- A)  $\zeta_j = 0$
- B) Se produce un error de margen
- C) No hay error de margen
- D) La muestra  $\mathbf{x}_j$  no es un vector soporte

- 4 ☐ En la red bayesiana



¿cuál de las relaciones siguientes es verdadera?

- A)  $P(a, b) = P(a) P(b)$
- B)  $P(a, d) = P(a) P(d)$
- C)  $P(a, b | d) = P(a | d) P(b | d)$
- D)  $P(a, c | b) = P(a | b) P(c | b)$

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 10 de enero de 2018

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 ☐ Sea  $S$  un conjunto de 1000 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 el error de reconocimiento dicho modelo. Para ello se ha utilizado el método de *validación cruzada en 10 bloques*, obteniéndose 4, 0, 4, 3, 4, 0, 3, 5, 4, 3 errores. Indicar cual de las siguientes afirmaciones es *correcta*.
- A) El test efectivo es de 100 muestras, el conjunto de entrenamiento efectivo es de 900 muestras, el error empírico  $\hat{p}$  es del 3 % y el intervalo de confianza es de  $3.0 \pm 0.1$  %
- B) El test efectivo es de 1000 muestras, el conjunto de entrenamiento efectivo es de 1000 muestras, el error empírico  $\hat{p}$  es del 3 % y el intervalo de confianza es de  $3 \pm 1$  %
- C) El test efectivo es de 1000 muestras, el conjunto de entrenamiento efectivo es de 900 muestras, el error empírico  $\hat{p}$  es del 3 % y el intervalo de confianza es de  $3 \pm 1$  %
- D) El test efectivo es de 1000 muestras, el conjunto de entrenamiento efectivo es de 900 muestras, el error empírico  $\hat{p}$  es del 6 % y el intervalo de confianza es de  $6 \pm 1$  %
- 2 ☐ Sea  $S = \{(\mathbf{x}_1, c_1), \dots, (\mathbf{x}_N, c_N)\}$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $c_n \in \{+1, -1\}$  una muestra de entrenamiento. El algoritmo perceptrón muestra a muestra ("online") trata de encontrar una solución  $\hat{\boldsymbol{\theta}}$  que satisfaga el sistema de  $N$  inecuaciones. Si queremos una solución con *margen*, esto es encontrar una solución  $\hat{\boldsymbol{\theta}}$  que satisfaga el sistema de  $N$  inecuaciones:  $c_n \boldsymbol{\theta}^t \mathbf{x}_n \geq b$ ,  $1 \leq n \leq N$ , para  $b \in \mathbb{R}^{\geq 0}$ . Indicar cuál de los siguientes algoritmos implementa la solución con margen, partiendo de que  $\boldsymbol{\theta}(1) = \text{arbitrario}$ :
- A)  $\boldsymbol{\theta}(k+1) = \begin{cases} \boldsymbol{\theta}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) \geq b \\ \boldsymbol{\theta}(k) + \rho_k c(k) \mathbf{x}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) < b \end{cases}$
- B)  $\boldsymbol{\theta}(k+1) = \begin{cases} \boldsymbol{\theta}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) < b \\ \boldsymbol{\theta}(k) + \rho_k c(k) \mathbf{x}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) \geq b \end{cases}$
- C)  $\boldsymbol{\theta}(k+1) = \begin{cases} b \boldsymbol{\theta}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) \geq 0 \\ b \boldsymbol{\theta}(k) + \rho_k c(k) \mathbf{x}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) < 0 \end{cases}$
- D)  $\boldsymbol{\theta}(k+1) = \begin{cases} \boldsymbol{\theta}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) \geq b \\ \boldsymbol{\theta}(k) + \rho_k b c(k) \mathbf{x}(k) & c(k) \boldsymbol{\theta}^t \mathbf{x}(k) < b \end{cases}$
- 3 ☐ Las siguientes afirmaciones se refieren al método Esperanza Maximización (EM) aplicado a una muestra de entrenamiento  $S$ . Identificar cuál de ellas es *correcta*:
- A) EM es útil para estimar valores máximo-verosímiles de los parámetros de modelos estadísticos a partir de  $S$  cuando hay variables latentes o ocultos.
- B) EM no se puede aplicar en la estimación de valores máximo-verosímiles de los parámetros de modelos estadísticos a partir de  $S$  cuando no hay variables latentes o ocultas.
- C) La rapidez de convergencia de EM puede mejorarse eligiendo un factor de aprendizaje adecuado para  $S$ .
- D) La rapidez de convergencia de EM siempre puede mejorarse inicializando los parámetros a cero.
- 4 ☐ En la red bayesiana cuya distribución conjunta es  $P(x_1, x_2, x_3, x_4) = P(x_1) P(x_2 | x_1) P(x_3 | x_1) P(x_4 | x_2, x_3)$  ¿cuál de las relaciones siguientes es *falsa* en general?
- A)  $P(x_1, x_4 | x_3) = P(x_1 | x_3) P(x_4 | x_3)$
- B)  $P(x_2, x_3 | x_1) = P(x_2 | x_1) P(x_3 | x_1)$
- C)  $P(x_2, x_3) = P(x_2) P(x_3)$
- D)  $P(x_2, x_3 | x_4) = P(x_2 | x_4) P(x_3 | x_4)$

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 24 de enero de 2018

Apellidos:

Nombre:

**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 En el problema de optimización con restricciones

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

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) Existe un  $i$  tal que  $\alpha_i^* < 0$  y  $v_i(\Theta^*) = 0$   
B) Si para algún  $j$ ,  $\alpha_j^* = 0$ , entonces  $v_j(\Theta^*) = 0$   
C) Existe un  $j$  tal que  $v_j(\Theta^*) > 0$  y  $\alpha_j^* = 0$   
D)  $v_j(\Theta^*) = 0 \forall j$  tal que  $\alpha_j^* > 0, 1 \leq j \leq k$
- 2 ☐ D Las siguientes afirmaciones se refieren a la estimación por máxima verosimilitud de los parámetros de una mezcla de  $K$  gaussianas (vector-media y peso de cada gaussiana) mediante un conjunto de vectores de entrenamiento cualquiera de dimensión  $D$ . Identifica cuál es *falsa*.

- A) Los parámetros de la mezcla se estiman adecuadamente mediante un algoritmo de *esperanza maximización* (EM)  
B) En cada iteración, el algoritmo EM estima los valores de las variables ocultas que, en este caso, son los pesos de las gaussianas.  
C) La verosimilitud del conjunto de entrenamiento, calculada con los parámetros estimados no disminuye en cada iteración del EM.  
D) El algoritmo EM obtiene los valores máximos de los parámetros a estimar.

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

$$q_S(\theta) = \sum_{n=1}^N (\theta^t x_n - y_n) + \frac{\lambda}{2} \theta^t \theta,$$

Cual de las siguientes expresiones del gradiente con respecto a  $\theta$  es correcta:

- A)  $\nabla q_S(\theta) = \sum_{n=1}^N x_n$   
B)  $\nabla q_S(\theta) = \theta^t \sum_{n=1}^N x_n$   
C)  $\nabla q_S(\theta) = \sum_{n=1}^N x_n + \lambda \theta$   
D)  $\nabla q_S(\theta) = \sum_{n=1}^N x_n + \lambda \theta^t \theta$

- 4 ☐ C Sea  $\mathcal{A}$  un conjunto de variables aleatorias y  $G$  el grafo que establece las dependencias entre las variables de  $\mathcal{A}$ . Un concepto importante en el que se basan las técnicas de redes bayesianas es:

- A) Los nodos del  $G$  representan las probabilidades incondicionales de las variables de  $\mathcal{A}$   
B)  $G$  define una distribución de probabilidad condicional entre dos subconjuntos de variables en  $\mathcal{A}$   
C)  $G$  define una distribución de probabilidad conjunta de todas las variables de  $\mathcal{A}$ . A partir de esta distribución, por inferencia probabilística puede calcularse cualquier probabilidad condicional o incondicional en la que intervengan dichas variables  
D)  $G$  define una distribución de probabilidad conjunta de todas las variables en  $\mathcal{A}$ . Para calcular las probabilidades de dicha distribución es necesario aplicar reglas de inferencia probabilística tales como la regla de Bayes y la marginalización.

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 7 de enero de 2019

Apellidos:

Nombre:

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

Marca cada recuadro con una única opción de entre las dadas. Cada acierto suma 0.4 puntos y cada fallo resta 1/6 puntos.

- 1 ☒ D Al aplicar el método de partición de datos denominado validación cruzada en 5 bloques a un conjunto de 1000 muestras, el clasificador obtiene, por bloque, 2, 5, 7, 8, 5 errores. Indicar la opción *correcta*:

- A) El error es de  $0.5\% \pm 0.5\%$ .
- B) La talla de entrenamiento efectiva es de 900 muestras.
- C) La talla de entrenamiento efectiva es de 1000 muestras.
- D) El error es de  $2.7\% \pm 1\%$ .

- 2 ☒ C En un clasificador en 3 clases resulta que la probabilidad a-posteriori de cada clase,  $y$ , dada una muestra  $\mathbf{x}$  es:

$y$	$P(Y = y   \mathbf{x})$
A	0.1
B	0.6
C	0.3

Indicar cuál es la opción *errónea*:

- A) La probabilidad de error si se toma la decisión  $Y = B$  es 0.4.
- B) La mínima probabilidad de error es 0.4.
- C) La probabilidad de error si se toma la decisión  $Y = C$  es 0.4.
- D) La peor decisión es  $Y = A$ , cuya probabilidad de error es 0.9.

- 3 ☒ A En el problema de optimización con restricciones

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

sea  $\Theta^*$  la solución óptima y sean  $\alpha_i^*$ ,  $1 \leq i \leq k$ , y  $\beta_i^*$ ,  $1 \leq i \leq m$ , los multiplicadores de Lagrange óptimos para las restricciones de desigualdad e igualdad, respectivamente. Indicar cuál de las siguientes afirmaciones es *falsa*:

- A) Si para algún  $j$ ,  $\alpha_j^* = 0$ , entonces  $v_j(\Theta^*) = 0$ .
- B) Para  $1 \leq i \leq m$   $u_i(\Theta^*) = 0$ .
- C) Para  $1 \leq i \leq k$   $v_i(\Theta^*) \geq 0$ .
- D)  $v_j(\Theta^*) = 0 \forall j$  si  $\alpha_j^* > 0$ ,  $1 \leq j \leq k$ .

- 4 ☒ D Las siguientes afirmaciones se refieren a la estimación por máxima verosimilitud de los parámetros de un modelo mediante el algoritmo de *esperanza maximización* (EM). Identificar cuál es *falsa*.

- A) En el paso E se estiman los valores de las variables ocultas (o se calculan sus probabilidades).
- B) En el paso M se calcula el máximo de una función auxiliar que depende de las estimaciones del paso E.
- C) El algoritmo EM se puede aplicar incluso cuando hay restricciones en los valores de los valores de los parámetros o de las variables ocultas, pero para ello hay que aplicar también la técnica de los multiplicadores de Lagrange.
- D) Si se usa de algoritmo EM es innecesaria la aplicación de la técnica de los multiplicadores de Lagrange.

- 5 ☒ C En una red bayesiana, sea  $\mathcal{A}$  un conjunto de variables aleatorias y  $G$  el grafo que establece las dependencias entre las variables de  $\mathcal{A}$ . Identificar cuál de las siguientes afirmaciones es *cierta*.

- A) Los arcos de  $G$  representan las probabilidades condicionales de las variables de  $\mathcal{A}$ .
- B)  $G$  define una distribución de probabilidad condicional entre las variables en  $\mathcal{A}$ .
- C) Cualquier distribución condicional o conjunta en la que participen todas o cualquier subconjunto de las variables de  $\mathcal{A}$ , se puede deducir a partir de la distribución conjunta definida por  $G$ .
- D) Si el valor de la variable asociada a un nodo  $\nu$  de  $G$  está dada, entonces todas las variables asociadas a los nodos que están directamente conectados con  $\nu$  son independientes entre si.



Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 7 de enero de 2020

Apellidos:

Nombre:

**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 500 muestras. Se han obtenido un total de 55 errores. Indicar cuál de las afirmaciones siguientes es razonable:

- A) Las tallas de entrenamiento y test efectivas son 495 y 500 muestras, respectivamente, y el error estimado es  $11.0 \pm 2.7\%$
- B) Las tallas de entrenamiento y test efectivas son 100 y 400 muestras, respectivamente, y el error estimado es  $13.8 \pm 3.0\%$
- C) Las tallas de entrenamiento y test efectivas son 5 y 495 muestras, respectivamente y el error estimado es  $11.1 \pm 2.8\%$
- D) Ninguna de las anteriores afirmaciones es razonable

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

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

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 este caso, ¿cuál de los siguientes gradientes,  $\nabla q_S(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\boldsymbol{\theta}(k)}$ , es correcto?:

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

- 3 ☐ A En el problema de optimización con restricciones

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

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

- A) Para todo  $i$  tal que  $\alpha_i^* > 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$
- B) Para todo  $i$  tal que  $\alpha_i^* < 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$
- C) Si para un  $i$ ,  $\alpha_i^* = 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$
- D) Para todo  $i$ , si  $\alpha_i^* = 0$ , entonces  $v_i(\boldsymbol{\Theta}^*) = 0$ ,

- 4 ☐ A Las siguientes afirmaciones se refieren a la estimación por máxima verosimilitud de los parámetros de una mezcla de  $K$  gaussianas (vector-media y peso de cada gaussiana) mediante un conjunto de vectores de entrenamiento cualquiera de dimensión  $D$ . Identifica cuál es *falsa*.

- A) El algoritmo *esperanza-maximización* es una alternativa a la técnica de los Multiplicadores de Lagrange en el caso de la estimación de los parámetros de una mezcla de  $K$  gaussianas.
- B) La verosimilitud del conjunto de entrenamiento, calculada con los parámetros estimados, aumenta en cada iteración del *esperanza-maximización*.
- C) En cada iteración, el algoritmo *esperanza-maximización* realiza una estimación de los valores de los pesos de las gaussianas.
- D) Los parámetros de la mezcla se estiman adecuadamente mediante un algoritmo de *esperanza-maximización*

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  $\alpha_j^*$ , 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  $\theta$  y  $\theta_0$  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

Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 18 de enero de 2021

Apellidos:

Nombre:

Grupo:

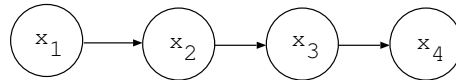
**Cuestiones (2 puntos; tiempo estimado: 30 minutos)**

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 Si la función de activación en la capa de salida de un perceptrón de dos capas es lineal, las fórmulas que permiten modificar los pesos de dicha capa de salida en el algoritmo BackProp verifican que (solo una respuesta es correcta):

- A)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) g(\phi_i^2) (1 - g(\phi_i^2)) s_j^1$   
B)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) g(\phi_i^2) s_j^1$   
C)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) s_j^1$   
D)  $\Delta\theta_{ij}^2 = \rho (t_i - s_i^2) (1 - g(\phi_i^2)) s_j^1$

- 2 ☐ A En la red bayesiana lineal



¿cuál de las relaciones siguientes es correcta?

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

- 3 ☐ D Se ha evaluado un sistema de Aprendizaje Automático mediante la técnica de *validación cruzada en B bloques* (“B-fold Cross Validation”) con  $B = 10$  y utilizando un conjunto de datos etiquetados que contiene 1000 muestras. Se han obtenido un total de 20 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 900 muestras y el error es del  $20.0 \pm 0.2 \%$   
C) La talla de entrenamiento efectiva es de 1000 muestras y el error es del  $20.0 \pm 0.2 \%$   
D) La talla de entrenamiento efectiva es de 900 muestras y el error es del  $2.0 \pm 0.9 \%$

- 4 ☐ A Considerar la siguiente modificación de la función de Widrow y Hoff

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

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



Examen de Aprendizaje Automático  
ETSINF, Universitat Politècnica de València, 3 de febrero de 2021

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 En el problema de optimización con restricciones

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

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  $i$ ,  $\alpha_i^* = 0$ , entonces  $v_i(\Theta^*) = 0$   
C) Si para un  $i$ ,  $\alpha_i^* > 0$ , entonces  $v_i(\Theta^*) > 0$   
D) Si para un  $i$ ,  $\alpha_i^* = 0$ , entonces  $u_i(\Theta^*) = 0$
- 2 ☐ A En la estimación por máxima verosimilitud de los parámetros de una mezcla de  $K$  gaussianas de matriz de covarianza común y conocida a partir de  $N$  vectores de entrenamiento, los parámetros a estimar son: el vector-media  $\mu_k$  y el peso  $\alpha_k$  de cada gaussiana,  $k, 1 \leq k \leq K$ . Identificar cuál de las siguientes afirmaciones es *correcta*:

- A) El método más adecuado es el de *esperanza-maximización* (EM), el cual garantiza que se cumple la restricción  $\sum_{k=1}^K \alpha_k = 1$ . Esto es así gracias a que, en cada iteración de EM, los valores de  $\alpha_k, 1 \leq k \leq K$ , se obtienen como medias de valores de variables latentes, usando una expresión que se deriva analíticamente mediante la técnica de los *multiplicadores de Lagrange* con la restricción indicada.  
B) Se puede usar *descenso por gradiente*, ya que los valores de  $\mu_k$  no están sujetos a ninguna restricción, lo que hace innecesario recurrir a la técnica de los *multiplicadores de Lagrange*.  
C) La solución se obtiene en un paso, utilizando directamente la *optimización lagrangiana* de la verosimilitud de los  $N$  vectores de entrenamiento. En este caso, hay un único multiplicador de Lagrange,  $\beta$ , asociado a la restricción de igualdad:  $\sum_{k=1}^K \alpha_k = 1$ .  
D) El método más adecuado sería el de *esperanza-maximización* (EM), pero no es posible utilizarlo ya que EM es un método iterativo que no garantiza el cumplimiento de la restricción de igualdad:  $\sum_{k=1}^K \pi_k = 1$ .

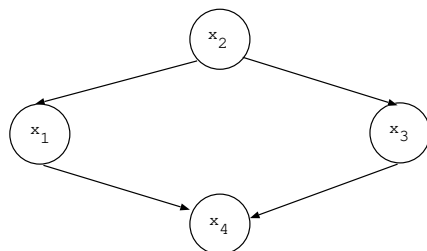
- 3 ☐ D Se desea ajustar por mínimos cuadrados la función  $f: \mathbb{R} \rightarrow \mathbb{R}$ , definida como:  $y = f(x) \stackrel{\text{def}}{=} ax^2 + bx + c$  a una secuencia de  $N$  pares entrada-salida:  $S = ((x_1, y_1), (x_2, y_2), \dots, (x_N, y_N))$ . La técnica empleada es minimizar por descenso por gradiente la función de error cuadrático:

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

Identifica la afirmación acertada de entre las siguientes:

- A) El gradiente es  $2ax + b$   
B) El descenso por gradiente solo es aplicable a funciones convexas, pero  $q(\cdot)$  no lo es.  
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 gradiente es:  $2 \sum_{n=1}^N (f(x_n) - y_n) [x_n^2, x_n, 1]^t$

- 4 ☐ A En la red bayesiana



¿cuál de las relaciones siguientes es falsa en general?

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