

Examen de RFA, MIARFID, DSIC, UPV, 2-11-2023

Apellidos: Nombre:

Cuestiones (1 punto)

Marca cada recuadro con una única opción de entre las dadas o déjalo en blanco si no quieres contestar.

Puntuación: $(A - E/3) \cdot 1/6$ con A aciertos y E errores; $A, E \in \{0, 1, \dots, 6\}$, $A + E \leq 6$.

- 1 ☐ Sea un problema de clasificación en C clases, $y \in \{1, 2, \dots, C\}$, para objetos representados mediante vectores de D características. Los clasificadores naive Bayes son clasificadores generativos basados en la asunción naive Bayes. Indica cuál de las siguientes opciones expresa correctamente dicha asunción.

1. $p(y = c \mid \mathbf{x}, \boldsymbol{\theta}) = \prod_{d=1}^D p(y = c \mid x_d, \boldsymbol{\theta}_{dc})$

2. $p(y = c \mid \mathbf{x}, \boldsymbol{\theta}) = \sum_{d=1}^D p(y = c \mid x_d, \boldsymbol{\theta}_{dc})$

3. $p(\mathbf{x} \mid y = c, \boldsymbol{\theta}) = \prod_{d=1}^D p(x_d \mid y = c, \boldsymbol{\theta}_{dc})$

4. $p(\mathbf{x} \mid y = c, \boldsymbol{\theta}) = \sum_{d=1}^D p(x_d \mid y = c, \boldsymbol{\theta}_{dc})$

- 2 ☐ Indica la afirmación incorrecta (o escoge la última opción si las tres primeras son correctas):

1. Análisis discriminante Gaussiano asume que las densidades condicionales de las clases son Gaussianas independientes, lo que resulta en discriminantes (log-posteriors) cuadráticas.
2. Análisis discriminante lineal asume que las densidades condicionales de las clases son Gaussianas de matriz de covarianzas común y diagonal, lo que resulta en discriminantes (log-posteriors) lineales.
3. LDA coincide con regresión logística (multinomial) en términos de modelo, pero se diferencia mucho en entrenamiento.
4. Todas son correctas. (la 2 no; la matriz común no tiene porqué ser diagonal)

- 3 ☐ Sea un problema de clasificación en dos clases, $y \in \{0, 1\}$, y sea $p(y \mid \mathbf{x}; \boldsymbol{\theta}) = \text{Ber}(y \mid \sigma(\mathbf{w}^t \mathbf{x} + b))$ un modelo de regresión logística binaria con $\mathbf{w} = (1, 1)^t$ y $b = -4$. La probabilidad de que $\mathbf{x} = (1, 1)^t$ pertenezca a la clase 1, $p = p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta})$, es:

1. $p < 0.25$ 0.1192

2. $0.25 \leq p < 0.5$

3. $0.5 \leq p < 0.75$

4. $0.75 \leq p$

- 4 [3] Sea un problema de clasificación en $C = 3$ clases, y sea $p(y | \mathbf{x}; \boldsymbol{\theta}) = \text{Cat}(y | \mathcal{S}(\mathbf{W}^t \mathbf{x} + \mathbf{b}))$, un modelo de regresión logística multinomial con

$$\mathbf{W}^t = \begin{pmatrix} 1 & 1 \\ 1 & 0 \\ 0 & 1 \end{pmatrix} \in \mathbb{R}^{C \times D} \quad \text{y} \quad \mathbf{b} = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}$$

La probabilidad de que $\mathbf{x} = (1, 1)^t$ pertenezca a la clase 2, $p(y = 2 | \mathbf{x}; \boldsymbol{\theta})$, es:

1. $p < 0.25$
2. $0.25 \leq p < 0.5$
3. $0.5 \leq p < 0.75$ 0.5761
4. $0.75 \leq p$

- 5 [4] Indica la afirmación incorrecta sobre perceptrones (redes) multicapa:

1. El problema XOR es un problema famoso del libro Perceptrons con datos 2d no linealmente separables. No puede resolverse con un perceptrón clásico (con activación Heaviside), pero sí con un perceptrón multicapa que combina tres perceptrones clásicos.
2. Los perceptrones multicapa clásicos han caído en desuso ya que emplean una función de activación no diferenciable. Actualmente se emplean redes multicapa con funciones de activación diferenciables como la sigmoide, tanh y ReLU.
3. La revolución del aprendizaje profundo se produce a partir de los 2010, gracias a una mayor disponibilidad de (grandes) conjuntos de datos, potencia de cálculo (en GPU sobre todo) y librerías especializadas de código abierto como tensorflow y pytorch.
4. Las redes con muchas capas (profundas) suelen ofrecer mejores resultados que las redes con pocas capas (superficiales). De hecho, se sabe que un perceptrón de una sola capa oculta se halla muy limitado como aproximador de funciones. (en teoría basta una capa oculta con suficientes unidades)

- 6 [4] Indica la afirmación incorrecta sobre el entrenamiento de redes con descenso por gradiente estocástico (o escoge la última opción si las tres primeras son correctas):

1. El ajuste del factor de aprendizaje es necesario para garantizar la convergencia a una buena solución.
2. Al entrenar modelos (muy) profundos, se suele presentar el problema de los gradientes que desaparecen y explotan. La explosión de gradientes puede evitarse con *gradient clipping*.
3. La desaparición de gradientes puede aliviarse modificando la arquitectura de la red con funciones de activación no saturantes (en lugar de saturantes), estandarización de activaciones (capas batchnorm) e inclusión de conexiones aditivas junto con las multiplicativas (capas con conexiones residuales).
4. Todas son correctas.

Problema (1.5 puntos)

Sea un perceptrón con una capa oculta,

$$h = \sigma(z) \quad \text{con} \quad z = \mathbf{W}\mathbf{x} + \mathbf{b}_1 \quad \mathbf{W} = \begin{pmatrix} -1 & -1 \\ 1 & -1 \end{pmatrix} \quad \mathbf{b}_1 = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

y una capa de salida,

$$\hat{y} = \mathbf{V}\mathbf{h} + b_2 \quad \text{con} \quad \mathbf{V} = \begin{pmatrix} 1 & -1 \end{pmatrix} \quad b_2 = 1$$

Dada la entrada $\mathbf{x} = (2, 2)^t$ y salida $y = -1$, se desea ajustar el modelo con el fin de minimizar la pérdida cuadrática,

$$\mathcal{L} = \frac{1}{2}(y - \hat{y})^2$$

Se pide:

1. (0.5 puntos) *Forward*: halla las pre-activaciones, activaciones y pérdida correspondientes al par entrada-salida dado.

```
1 import numpy as np; np.set_printoptions(precision=4)
2 sigmoid = lambda x: 1.0 / (1.0 + np.exp(-x));
3 x = np.array([2, 2]); y = -1
4 W = np.array([[ -1, -1], [1, -1]]); b1 = np.array([1, -1])
5 V = np.array([1, -1]); b2 = 1
6 z = (W @ x + b1); print("z =", z)
7 h = sigmoid(z); print("h =", h)
8 haty = V @ h + b2; print("haty =", round(haty, 4))
9 loss = .5 * np.square(y-haty).sum(); print("loss =", round(loss, 4))
```

```
1 z = [-3 -1]
2 h = [0.0474 0.2689]
3 haty = 0.7785
4 loss = 1.5815
```

2. (0.75 puntos) *Backward*: usa los cálculos anteriores y halla las Jacobianas de la pérdida con respecto a . . .

$u = \frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y$	la predicción (activación de la capa de salida)
$\mathbf{g}_V = u \frac{\partial \hat{y}}{\partial \mathbf{V}} = \dots$	la transformación lineal de la capa de salida
$g_{b_2} = u \frac{\partial \hat{y}}{\partial b_2} = \dots$	el sesgo de la capa de salida
$\mathbf{u}^t = u \frac{\partial \hat{y}}{\partial \mathbf{h}} = u \mathbf{V}$	la activación de la capa oculta
$\mathbf{u}^t = \mathbf{u}^t \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \mathbf{u}^t \text{diag}(\boldsymbol{\sigma}'(\mathbf{z}))$	la pre-activación de la capa oculta
$\mathbf{g}_W = \mathbf{u}^t \frac{\partial \mathbf{z}}{\partial \mathbf{W}} = \dots$	la transformación lineal de la capa oculta
$\mathbf{g}_{b_1} = \mathbf{u}^t \frac{\partial \mathbf{z}}{\partial b_1} = \dots$	el sesgo de la capa oculta

$u = \frac{\partial \mathcal{L}}{\partial \hat{y}} = \hat{y} - y$	la predicción (activación de la capa de salida)
$\mathbf{g}_V = u \frac{\partial \hat{y}}{\partial \mathbf{V}} = \mathbf{h}u$	la transformación lineal de la capa de salida
$g_{b_2} = u \frac{\partial \hat{y}}{\partial b_2} = u$	el sesgo de la capa de salida
$\mathbf{u}^t = u \frac{\partial \hat{y}}{\partial \mathbf{h}} = u \mathbf{V}$	la activación de la capa oculta
$\mathbf{u}^t = \mathbf{u}^t \frac{\partial \mathbf{h}}{\partial \mathbf{z}} = \mathbf{u}^t \text{diag}(\boldsymbol{\sigma}'(\mathbf{z}))$	la pre-activación de la capa oculta
$\mathbf{g}_W = \mathbf{u}^t \frac{\partial \mathbf{z}}{\partial \mathbf{W}} = \mathbf{x} \mathbf{u}^t$	la transformación lineal de la capa oculta
$\mathbf{g}_{b_1} = \mathbf{u}^t \frac{\partial \mathbf{z}}{\partial b_1} = \mathbf{u}^t$	el sesgo de la capa oculta

```

1 J_haty = haty-y;
2 J_V = np.outer(h, J_haty);
3 J_b2 = J_haty;
4 J_h = J_haty * V;
5 J_z = J_h * sigmoid(z) * sigmoid(-z);
6 J_W = np.outer(x, J_z);
7 J_b1 = J_z;

print("J_haty =", round(J_haty, 4))
print("J_V =", J_V)
print("J_b2 =", round(J_b2, 4))
print("J_h =", J_h)
print("J_z =", J_z)
print("J_W =", J_W)
print("J_b1 =", J_b1);

```

```

1 J_haty = 1.7785
2 J_V = [[0.0843]
3        [0.4783]]
4 J_b2 = 1.7785
5 J_h = [ 1.7785 -1.7785]
6 J_z = [ 0.0803 -0.3497]
7 J_W = [[ 0.1607 -0.6993]
8        [ 0.1607 -0.6993]]
9 J_b1 = [ 0.0803 -0.3497]

```

3. (0.25 puntos) *Actualización de parámetros*: actualiza las transformaciones lineales y sesgos del modelo a partir de las Jacobianas halladas, con factor de aprendizaje $\eta = 1$.

```
1 V = V - 1.0 * J_V.T; print("V =", V)
2 b2 = b2 - 1.0 * J_b2; print("b2 =", b2)
3 W = W - 1.0 * J_W.T; print("W =", W)
4 b1 = b1 - 1.0 * J_b1; print("b1 =", b1)
```

```
1 V = [[ 0.9157 -1.4783]]
2 b2 = -0.7785
3 W = [[-1.1607 -1.1607]
4      [ 1.6993 -0.3007]]
5 b1 = [ 0.9197 -0.6503]
```