



UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA

# Regresión Logística

Jorge Civera  
Alfons Juan  
Albert Sanchis

*DSIC*

Departamento de Sistemas  
Informáticos y Computación

# Objetivos formativos

- Explicar la distribución categórica
- Representar la codificación one-hot
- Describir el modelo probabilístico de clasificación con la función softmax
- Describir el modelo de regresión logística
- Describir el aprendizaje por máxima verosimilitud
- Aplicar descenso por gradiente en regresión logística

# Índice

<b>1</b>	<b>Distribución categórica y codificación one-hot</b>	<b>3</b>
<b>2</b>	<b>Modelo probabilístico de clasificación softmax</b>	<b>5</b>
<b>3</b>	<b>Regresión logística</b>	<b>7</b>
<b>4</b>	<b>Aprendizaje por máxima verosimilitud</b>	<b>9</b>
<b>5</b>	<b>Aprendizaje con descenso por gradiente</b>	<b>13</b>
<b>6</b>	<b>Conclusiones</b>	<b>17</b>

# 1. Distribución categórica y codificación one-hot

- **Variable categórica:** variable aleatoria que toma un valor de un conjunto finito de categorías (no ordenadas)
- **Ejemplos de variables categóricas:** *etiqueta de clase*, palabra de un vocabulario, etc.
- **Codificación one-hot:** de una variable categórica  $y$  que toma un valor entre  $C$  posibles,  $\{1, \dots, C\}$

$$\text{one-hot}(y) = \mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_C \end{pmatrix} \in \{0, 1\}^C \quad \text{con} \quad \sum_c y_c = 1$$

- **Ejemplo:** Codificación one-hot de la muestra  $\mathbf{x}_1 = (0, 0)^t$  de la clase  $c_1 = 1$  y  $\mathbf{x}_2 = (1, 1)^t$  de la clase  $c_2 = 2$ :

$$\mathbf{y}_1 = (1, 0)^t$$

$$\mathbf{y}_2 = (0, 1)^t$$

# Codificación one-hot y distribución categórica

- **Distribución categórica:** distribución de probabilidades entre las  $C$  posibles categorías de una variable categórica, que viene dada por un vector de parámetros  $\theta \in [0, 1]^C$  tal que  $\sum_c \theta_c = 1$

$$p(y \mid \theta) = \text{Cat}(\mathbf{y} \mid \theta) = \prod_{c=1}^C \theta_c^{y_c}$$

- **Convención:**  $0^0 = 1$

- **Ejemplo:**

$$\theta = (0.5, 0.5, 0)^t, \text{Cat}(\mathbf{y} = (1, 0, 0)^t \mid \theta) = 0.5^1 0.5^0 0^0 = 0.5$$

## 2. Modelo probabilístico de clasificación softmax

- **Normalización probabilística de clasificadores:** sea  $G$  cualquier clasificador definido con funciones discriminantes generales  $[g_1, \dots, g_C]$ , se puede definir uno equivalente  $G'$  con funciones discriminantes normalizadas probabilísticamente  $[g'_1, \dots, g'_C]$

$$\begin{aligned} c(\mathbf{x}) &= \operatorname{argmax}_c g_c(\mathbf{x}) \\ &= \operatorname{argmax}_c e^{g_c(\mathbf{x})} \text{ con } h(z) = e^z \in \mathbb{R}^{\geq 0} \text{ estrictamente creciente} \\ &= \operatorname{argmax}_c \frac{e^{g_c(\mathbf{x})}}{\sum_{c'} e^{g_{c'}(\mathbf{x})}} \text{ con } h(z) = kz, k \text{ constante positiva} \end{aligned}$$

Por tanto,  $g'_c(\mathbf{x}) = \frac{e^{g_c(\mathbf{x})}}{\sum_{c'} e^{g_{c'}(\mathbf{x})}}$  define un clasificador equivalente

Esta transformación es conocida como **función softmax**

- En este modelo probabilístico se asume que los valores  $g_c(\mathbf{x})$  son log-probabilidades no normalizadas denominados **logits**

# Modelo probabilístico de clasificación softmax

- **La función softmax:** transforma un vector de **logits** (log-probabilidades no normalizadas)  $G \in \mathbb{R}^C$  en uno de probabilidades  $G' \in [0, 1]^C$

$$G' = \mathcal{S}(G) = \left[ \frac{e^{g_1}}{\sum_{c'} e^{g_{c'}}}, \dots, \frac{e^{g_C}}{\sum_{c'} e^{g_{c'}}} \right]$$

donde se cumple

$$0 \leq \mathcal{S}(G)_c \leq 1 \quad \text{y} \quad \sum_c \mathcal{S}(G)_c = 1$$

### 3. Regresión logística

- **Regresión logística:** modelo con softmax y funciones discriminantes lineales (en notación homogénea)

$$p(\mathbf{y} \mid \mathbf{x}, \mathbf{W}) = \text{Cat}(\mathbf{y} \mid \boldsymbol{\mu})$$

donde

$$\boldsymbol{\mu} = \mathcal{S}(\mathbf{a}), \quad \mathbf{a} = f(\mathbf{x}; \mathbf{W}) = \mathbf{W}^t \mathbf{x}, \quad \mathbf{W} \in \mathbb{R}^{D \times C} \quad \text{y} \quad \mathbf{x} \in \mathbb{R}^D$$

- No hay diferencia con los clasificadores basados en funciones discriminantes lineales, a excepción de que ahora predecimos las probabilidades de todas las clases



## Ejemplo

Sea un modelo de regresión logística en notación homogénea para un problema de clasificación en  $C = 2$  clases y datos representados mediante vectores de dimensión  $D = 2$

$$\mu = \mathcal{S}(\mathbf{a}), \quad \mathbf{a} = f(\mathbf{x}; \mathbf{W}) = \mathbf{W}^t \mathbf{x} \quad \text{con}$$

$$\mathbf{W}^t = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \end{pmatrix} \quad \text{y} \quad \mathbf{x} = \begin{pmatrix} 1 \\ x_1 \\ x_2 \end{pmatrix}$$

Calcula la probabilidad de que  $\mathbf{x} = (1, 0, 0)^t$  y  $\mathbf{x} = (1, 1, 1)^t$  pertenezcan a cada clase:

$\mathbf{x}^t$	$\mathbf{a}^t$	$\mu_1 = \mathcal{S}(\mathbf{a})_1$	$\mu_2 = \mathcal{S}(\mathbf{a})_2$
$(1, 0, 0)$	$(1, -1)$	$\frac{e^1}{e^1 + e^{-1}} = 0.8808$	$\frac{e^{-1}}{e^1 + e^{-1}} = 0.1192$
$(1, 1, 1)$	$(-1, 1)$	$\frac{e^{-1}}{e^{-1} + e^1} = 0.1192$	$\frac{e^1}{e^{-1} + e^1} = 0.8808$

## 4. Aprendizaje por máxima verosimilitud

- **Objetivo:** establecer un criterio para aprender  $\mathbf{W}$  a partir de un conjunto de datos de entrenamiento,  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$
- **Log-verosimilitud (condicional):** log-probabilidad de  $\mathcal{D}$  interpretada como función de  $\mathbf{W}_N$

$$\begin{aligned} \text{LL}(\mathbf{W}) &= \log p(\mathcal{D} \mid \mathbf{W}) = \log \prod_{n=1}^N p(\mathbf{y}_n \mid \mathbf{x}_n, \mathbf{W}) \\ &= \sum_{n=1}^N \log \text{Cat}(\mathbf{y}_n \mid \boldsymbol{\mu}_n) \quad \text{con} \quad \boldsymbol{\mu}_n = \mathcal{S}(\mathbf{a}_n) \quad \text{y} \quad \mathbf{a}_n = \mathbf{W}^t \mathbf{x}_n \\ &= \sum_{n=1}^N \log \prod_{c=1}^C \mu_{nc}^{y_{nc}} = \sum_{n=1}^N \sum_{c=1}^C y_{nc} \log \mu_{nc} \end{aligned}$$

- **Aprendizaje por máxima verosimilitud:** elegimos una  $\mathbf{W}$  que otorgue máxima probabilidad a  $\mathcal{D}$

$$\mathbf{W}^* = \underset{\mathbf{W}}{\operatorname{argmax}} \text{LL}(\mathbf{W})$$

## Ejemplo

Calcula la log-verosimilitud de  $\mathbf{W}^t = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \end{pmatrix}$  con dos datos  $\mathcal{D} = \{((1, 0, 0)^t, (1, 0)^t), ((1, 1, 1)^t, (0, 1)^t)\}$

$$\text{LL}(\mathbf{W}) = \log p(\mathcal{D} \mid \mathbf{W}) = \sum_{n=1}^N \sum_{c=1}^C y_{nc} \log \mu_{nc}$$

$$= y_{11} \log \mu_{11} + y_{12} \log \mu_{12} + y_{21} \log \mu_{21} + y_{22} \log \mu_{22}$$

$$= \log \mu_{11} + \log \mu_{22}$$

$$= \log 0.8808 + \log 0.8808 = -0.1269 - 0.1269 = -0.2538$$

# Planteamiento como problema de minimización

- **Neg-log-verosimilitud:** log-verosimilitud con el signo cambiado y normalizada por el número de datos

$$\text{NLL}(\mathbf{W}) = -\frac{1}{N} \text{LL}(\mathbf{W}) = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_{nc} \log \mu_{nc}$$

- **Ejemplo:** neg-log-verosimilitud del ejemplo anterior

$$\text{NLL}(\mathbf{W}) = -\frac{1}{2} \text{LL}(\mathbf{W}) = 0.1269$$

- **Aprendizaje por mínima NLL:** aprendizaje por máxima verosimilitud planteado como un problema de minimización

$$\mathbf{W}^* = \underset{\mathbf{W}}{\text{argmin}} \text{NLL}(\mathbf{W})$$

## Ejemplo

Supongamos que tenemos que elegir por mínima NLL entre:

$$\mathbf{W}^t = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \end{pmatrix} \quad \text{y} \quad \tilde{\mathbf{W}}^t = \begin{pmatrix} -1 & 1 & 1 \\ 1 & -1 & -1 \end{pmatrix}$$

con los datos  $\mathcal{D} = \{((1, 0, 0)^t, (1, 0)^t), ((1, 1, 1)^t, (0, 1)^t)\}$

Según el ejemplo anterior, la NLL de  $\mathbf{W}$  es 0.1269 y la de  $\tilde{\mathbf{W}}$ :

$$\text{NLL}(\tilde{\mathbf{W}}) = -\frac{1}{2}(\log \tilde{\mu}_{11} + \log \tilde{\mu}_{22}) = -\frac{1}{2}(\log \frac{e^{-1}}{e^{-1} + e^1} + \log \frac{e^{-1}}{e^{-1} + e^1}) = 2.1269$$

Por tanto, elegiríamos  $\mathbf{W}$  ya que su NLL es menor que la de  $\tilde{\mathbf{W}}$

## 5. Aprendizaje con descenso por gradiente

- **Descenso por gradiente:** algoritmo iterativo para minimizar una función  $\mathcal{L}(\theta)$  a partir de un valor inicial de los parámetros  $\theta_0$  dado

$$\theta_{i+1} = \theta_i - \eta_i \nabla \mathcal{L}(\theta)|_{\theta_i}$$

- **Factor de aprendizaje:**  $\eta_i > 0$  juega el mismo papel que en Perceptrón; podemos elegir un valor pequeño constante,  $\eta_i = \eta$
- **Dirección de descenso más pronunciada:**  $-\nabla \mathcal{L}(\theta)|_{\theta_i}$  es el neg-gradiente de la función evaluada en  $\theta_i$
- **Convergencia:** si  $\eta$  no es muy grande y la función es convexa (con forma de bol), converge a un mínimo (global)

# Descenso por gradiente en regresión logística

- La **NLL** es una función convexa
- **Gradiente de la NLL:** haremos uso del siguiente resultado

$$\begin{pmatrix} \frac{\partial \text{NLL}}{\partial W_{11}} & \cdots & \frac{\partial \text{NLL}}{\partial W_{1C}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \text{NLL}}{\partial W_{D1}} & \cdots & \frac{\partial \text{NLL}}{\partial W_{DC}} \end{pmatrix} = \frac{\partial \text{NLL}}{\partial \mathbf{W}^t} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n (\boldsymbol{\mu}_n - \mathbf{y}_n)^t$$

- **Descenso por gradiente aplicado a regresión logística:**

$$\mathbf{W}_0 = \mathbf{0}; \quad \mathbf{W}_{i+1} = \mathbf{W}_i - \eta_i \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n (\boldsymbol{\mu}_n - \mathbf{y}_n)^t$$

# Ejemplo

Sea un modelo de regresión logística en notación homogénea para un problema de clasificación en  $C = 2$  clases y datos de dimensión  $D = 2$ , actualiza el valor de  $\mathbf{W}$  aplicando descenso por gradiente con  $\eta = 0.1$ , matriz de pesos iniciales nulos y conjunto de entrenamiento  $\mathcal{D} = \{\mathbf{x}_1 = (1, 0, 0)^t, \mathbf{y}_1 = (1, 0)^t\}$

$$\mathbf{a} = \mathbf{W}^t \mathbf{x} = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$\boldsymbol{\mu} = S(\mathbf{a}) = \left( \frac{e^0}{e^0 + e^0}, \frac{e^0}{e^0 + e^0} \right)^t = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}$$



$$\begin{aligned}
\mathbf{W} &= \mathbf{W} - \eta \mathbf{x}(\boldsymbol{\mu} - \mathbf{y})^t \\
&= \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} - 0.1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} ((0.5, 0.5) - (1, 0)) \\
&= \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} - 0.1 \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} (-0.5, 0.5) \\
&= \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} - \begin{pmatrix} 0.1 \\ 0 \\ 0 \end{pmatrix} (-0.5, 0.5) \\
&= \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} - \begin{pmatrix} -0.05 & 0.05 \\ 0 & 0 \\ 0 & 0 \end{pmatrix} = \begin{pmatrix} 0.05 & -0.05 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}
\end{aligned}$$

## 6. Conclusiones

Hemos visto:

- La distribución categórica y la codificación one-hot
- El modelo probabilístico de clasificación con la función softmax y, en particular, el modelo de regresión logística
- El método de aprendizaje por máxima verosimilitud en regresión logística aplicando descenso por gradiente