



UNIVERSITAT  
POLITÈCNICA  
DE VALÈNCIA

# Regressió Logística

Jorge Civera  
Alfons Juan  
Albert Sanchis

*DSIC*

Departament de Sistemes  
Informàtics i Computació

# Objectius formatius

- Explicar la distribució categòrica
- Representar la codificació one-hot
- Descriure el model probabilístic de classificació amb la funció softmax
- Descriure el model de regressió logística
- Descriure l'aprenentatge per màxima versemblança
- Aplicar descens per gradient en regressió logística

# Índex

<b>1</b>	<b>Distribució categòrica i codificació one-hot</b>	<b>3</b>
<b>2</b>	<b>Model probabilístic de classificació softmax</b>	<b>5</b>
<b>3</b>	<b>Regressió logística</b>	<b>7</b>
<b>4</b>	<b>Aprenentatge per màxima versemblança</b>	<b>9</b>
<b>5</b>	<b>Aprenentatge amb descens per gradient</b>	<b>13</b>
<b>6</b>	<b>Conclusions</b>	<b>17</b>

# 1 Distribució categòrica i codificació one-hot

- **Variable categòrica:** variable aleatòria que pren un valor d'un conjunt finit de categories (no ordenades)
- **Exemples de variables categòriques:** etiqueta de classe, paraula d'un vocabulari, etc.
- **Codificació one-hot:** d'una variable categòrica  $i$  que pren un valor entre  $C$  possibles,  $\{1, \dots, C\}$

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

- **Exemple:** Codificació one-hot de la mostra  $\mathbf{x}_1 = (0, 0)^t$  de la classe  $c_1 = 1$  i  $\mathbf{x}_2 = (1, 1)^t$  de la classe  $c_2 = 2$ :

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

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

# Distribució categòrica i codificació one-hot

- **Distribució categòrica:** distribució de probabilitats entre les  $C$  possibles categories d'una variable categòrica, que ve donada per un vector de paràmetres  $\theta \in [0, 1]^C$  tal que  $\sum_c \theta_c = 1$

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

- **Convenció:**  $0^0 = 1$
- **Exemple:**

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

## 2 Model probabilístic de classificació softmax

- **Normalització probabilística de classificadors:** siga  $G$  qualsevol classificador definit amb funcions discriminants generals  $[g_1, \dots, g_C]$ , es pot definir un equivalent  $G'$  amb funcions discriminants normalitzades probabilísticament  $[g'_1, \dots, g'_C]$

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

Per tant,  $g'_c(\mathbf{x}) = \frac{e^{g_c(\mathbf{x})}}{\sum_{c'} e^{g_{c'}(\mathbf{x})}}$  defineix un classificador equivalent

Aquesta transformació és coneguda com **funció softmax**

- En aquest model probabilístic s'assumeix que els valors  $g_c(\mathbf{x})$  són log-probabilitats no normalitzades denominats **logits**

# Model probabilístic de classificació softmax

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

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

on es compleix

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

### 3 Regressió logística

- **Regressió logística:** model amb softmax i funcions discriminants lineals (en notació homogènia)

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

on

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

- No hi ha diferència amb els classificadors basats en funcions discriminants lineals, a excepció que ara predim les probabilitats de totes les classes



## Exemple

Siga un model de regressió logística en notació homogènia per a un problema de classificació en  $C = 2$  classes i dades representades mitjançant vectors de dimensió  $D = 2$

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

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

Calcula amb quina probabilitat  $\mathbf{x} = (1, 0, 0)^t$  i  $\mathbf{x} = (1, 1, 1)^t$  pertanyen a cada classe:

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

## 4 Aprenentatge per màxima versemblança

- **Objectiu:** establir un criteri per a aprendre  $\mathbf{W}$  a partir d'un conjunt de dades d'entrenament,  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^N$
- **Log-versemblança (condicional):** log-probabilitat de  $\mathcal{D}$  interpretada com a funció de  $\mathbf{W}$

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

- **Aprenentatge per màxima versemblança:** triem una  $\mathbf{W}$  que atorgue màxima probabilitat a  $\mathcal{D}$

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

## Exemple

Calcula la log-versemblança de  $\mathbf{W}^t = \begin{pmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \end{pmatrix}$  amb dues dades  $\mathcal{D} = \{((1, 0, 0)^t, (1, 0)^t), ((1, 1, 1)^t, (0, 1)^t)\}$

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

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

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

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

# Plantejament com a problema de minimització

- **Neg-log-versemblança:** log-versemblança amb el signe canviat i normalitzada pel nombre de dades

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

- **Exemple:** neg-log-versemblança de l'exemple anterior

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

- **Aprenentatge per mínima NLL:** aprenentatge per màxima versemblança plantejat com un problema de minimització

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

## Exemple

Suposem que hem de triar per mínima NLL entre:

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

amb les dades  $\mathcal{D} = \{((1, 0, 0)^t, (1, 0)^t), ((1, 1, 1)^t, (0, 1)^t)\}$

Segons l'exemple anterior, la NLL de  $\mathbf{W}$  és 0.1269 i la de  $\tilde{\mathbf{W}}$ :

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

Per tant, triaríem  $\mathbf{W}$  ja que la seua NLL és menor que la de  $\tilde{\mathbf{W}}$

## 5 Aprenentatge amb descens per gradient

- **Descens per gradient:** algorisme iteratiu per a minimitzar una funció  $\mathcal{L}(\theta)$  a partir d'un valor inicial dels paràmetres  $\theta_0$  donat

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

- **Factor d'aprenentatge:**  $\eta_i > 0$  juga el mateix paper que en Perceptró; podem triar un valor xicotet constant,  $\eta_i = \eta$
- **Direcció de descens més pronunciada:**  $-\nabla \mathcal{L}(\theta)|_{\theta_i}$  és el neg-gradient de la funció avaluada en  $\theta_i$
- **Convergència:** si  $\eta$  no és molt gran i la funció és convexa (amb forma de bol), convergeix a un mínim (global)

# Descens per gradient en regressió logística

- La **NLL** és una funció convexa
- **Gradient de la NLL:** farem ús del següent resultat

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

- **Descens per gradient aplicat a regressió logística:**

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

## Exemple

Siga un model de regressió logística en notació homogènia per a un problema de classificació en  $C = 2$  classes i dades de dimensió  $D = 2$ , actualitza el valor de  $\mathbf{W}$  aplicant descens per gradient amb  $\eta = 0.1$ , matriu de pesos inicials nuls i conjunt d'entrenament  $\mathcal{D} = \{\mathbf{x}_1 = (1, 0, 0)^t, \mathbf{y}_1 = (1, 0)^t\}$

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

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



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

## 6 Conclusions

Hem vist:

- La distribució categòrica i la codificació one-hot
- El model probabilístic de classificació amb la funció softmax i, en particular, el model de regressió logística
- El mètode d'aprenentatge per màxima versemblança en regressió logística aplicant descens per gradient