

## MLP PER MUTUALLY EXCLUSIVE MULTI-CLASS CLASSIFICATION

Fin ora abbiamo risolto il problema di classificare più classi, ma non in modo esclusivo. Il che vuol dire che in uscita potevamo avere più di una classe predetta.

**Esempio:** Classificare gli oggetti di una singola immagine.

### SOFT-MAX

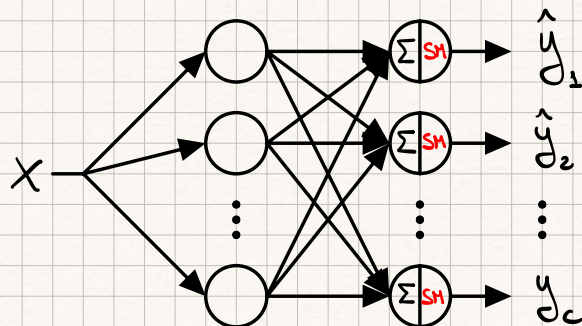
Come classificare in modo esclusivo? La risposta è usare una funzione di attivazione particolare, chiamata **softmax**, in grado di trasformare un vettore di ingresso in una distribuzione di probabilità (la somma degli output della softmax fa 1):

$$\text{Softmax}(\vec{z}) = \begin{bmatrix} \frac{e^{z_1}}{\sum_{j=1}^c e^{z_j}} \\ \vdots \\ \frac{e^{z_c}}{\sum_{j=1}^c e^{z_j}} \end{bmatrix}$$

**Esempio**

$$\begin{bmatrix} 1,3 \\ 5,1 \\ 2,2 \\ 0,7 \\ 1,1 \end{bmatrix} \xrightarrow{\text{softmax}} \begin{bmatrix} 0,02 \\ 0,90 \\ 0,05 \\ 0,01 \\ 0,02 \end{bmatrix} \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix}$$

## Modello



## CATEGORICAL CROSS-ENTROPY

$$\mathcal{L}^{CCS} = \sum_{k=1}^c -y_k \log(\hat{y}_k)$$

È la Cross usata per allenare modelli di classificazione esclusiva.

## DERIVATA SOFTMAX

Calcoliamo le derivate parziali della softmax in due dimensioni  $(x, y)$ :

$$\text{softmax}(x, y) = \left[ \frac{e^x}{e^x + e^y} ; \frac{e^y}{e^x + e^y} \right]$$

$$\text{softmax}_x(x, y) = \left[ \frac{e^x e^y}{(e^x + e^y)^2} ; -\frac{e^x e^y}{(e^x + e^y)^2} \right]$$

$$\text{softmax}_y(x, y) = \left[ -\frac{e^x e^y}{(e^x + e^y)^2} ; \frac{e^x e^y}{(e^x + e^y)^2} \right]$$



Per calcolare la derivata di un rapporto fra due funzioni abbiamo la seguente formula:

$$\left( \frac{f(x)}{g(x)} \right)' = \frac{f'(x)}{g(x)} - \frac{f(x)g'(x)}{g^2(x)}$$

Sostituendo  $f(x) = e^x$  e  $g(x) = e^x + e^y$  abbiamo che

$$\left( \frac{e^x}{e^x + e^y} \right)' = \frac{e^x}{e^x + e^y} - \left( \frac{e^x}{e^x + e^y} \right)^2 = \dots = \frac{e^x e^y}{(e^x + e^y)^2}$$

Considerando la derivata parziale di una sola componente della softmax rispetto alla prima variabile di preattivazione  $z_1$ :

$$\frac{\partial a_1}{\partial z_1} = \frac{e^{z_1}}{e^{z_1} + e^{z_2}} - \left( \frac{e^{z_1}}{e^{z_1} + e^{z_2}} \right)^2 = a_1 - a_1^2 = a_1(1 - a_1)$$

Rispetto alle altre preattivazioni,  $z_2 \dots z_c$  abbiamo

$$\frac{\partial a_1}{\partial z_2} = - \frac{e^{z_1} e^{z_2}}{(e^{z_1} + e^{z_2})^2} = -a_1 a_2$$

## In generale

Possiamo dire che, la  $c$ -esima componente della softmax ha derivata, rispetto alla preattivazione  $z_k$ :

$$\frac{\partial a_c}{\partial z_k} = \begin{cases} a_c(1 - a_c) & \text{se } c=k \\ -a_c a_k & \text{se } c \neq k \end{cases}$$

## Backpropagation

Usando la  $y^{ccs}$ , calcoliamo le derivate:

$$\frac{\partial \mathcal{L}^{ccs}}{\partial W_1} \quad e \quad \frac{\partial \mathcal{L}^{ccs}}{\partial W_2}$$

Considerando i seguenti layer finali:

$$\begin{aligned} z z_2 &= W_2 A_1 & z z_2 &\in \mathbb{R}^{C \times B} \\ z A_2 &= \text{softmax}(z z_2) & z A_2 &\in \mathbb{R}^{C \times B} \end{aligned}$$



Calcoliamo  $\partial \mathcal{L} / \partial W_2$ :

Dato un mini-Batch  $B$ , la loss diventa:

$$\mathcal{L}^{CCS} = \sum_{k=1}^C \sum_{b=1}^B -y_k \log(\hat{y}_k)$$

Considerando  $B=1$  e la classe  $c$ -esima:

$$\frac{\partial \mathcal{L}^{CCS}}{\partial \hat{y}_c} = -\frac{y_c}{\hat{y}_c} \xrightarrow{\hat{y}_c = zA_2^c} -\frac{y}{zA_2^c}$$

Per calcolare  $\partial \mathcal{L} / \partial W_2$  sviluppiamo la seguente chain rule:

$$\frac{\partial \mathcal{L}}{\partial W_2} = \underbrace{\frac{\partial \mathcal{L}}{\partial z z_2^c}}_{A_1^T} \underbrace{\frac{\partial z z_2^c}{\partial W_2}}_{A_2^T}$$

$$\frac{\partial \mathcal{L}}{\partial z z_2^c} = \sum_{k=1}^C \frac{\partial \mathcal{L}}{\partial A_2^k} \frac{\partial z A_2^k}{\partial z z_2^c} = \underbrace{\frac{\partial \mathcal{L}}{\partial A_2^c}}_{\text{Derivata Softmax } k=c} \underbrace{\frac{\partial z A_2^c}{\partial z z_2^c}}_{\text{Derivata della Softmax } k=c} + \sum_{\substack{k=1 \\ k \neq c}}^C \underbrace{\frac{\partial \mathcal{L}}{\partial A_2^k}}_{\text{Derivata Softmax } k \neq c} \underbrace{\frac{\partial z A_2^k}{\partial z z_2^c}}_{\text{Derivata della Softmax } k \neq c} =$$

→ Sommatoria delle derivate di ogni componente della softmax

$$= -\underbrace{\frac{y_c}{zA_2^c}}_{\partial \mathcal{L} / \partial \hat{y}_c} \cdot \underbrace{zA_2^c (1 - zA_2^c)}_{\text{Derivata Softmax } k=c} + \sum_{k \neq c} \underbrace{\left( -\frac{y_k}{zA_2^k} \right)}_{\partial \mathcal{L} / \partial \hat{y}_k} \underbrace{\left( -zA_2^c zA_2^k \right)}_{\text{Derivata della Softmax } k \neq c}$$

$$\begin{aligned}
 &= -y_c (1 - zA_2^c) + \sum_{k \neq c}^C y_k \cdot zA_2^c = \\
 &= -y_c + \underbrace{y_c \cdot zA_2^c + zA_2^c \cdot \sum_{k \neq c}^C y_k}_{\text{Mettiamo in evidenza } zA_2^c} =
 \end{aligned}$$

$$\begin{aligned}
 &= -y_c + zA_2^c \underbrace{\sum_{k=1}^C y_k}_{=1} \text{ per la definizione di softmax}
 \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial z z_2^c} = zA_2^c - y_c$$

Quindi, mettendo insieme:

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial z z_2^c} \frac{\partial z z_2^c}{\partial W_2} = (zA_2^c - y_c) A_1^T$$

Calcoliamo  $\partial \mathcal{L} / \partial W_1$

$$\frac{\partial \mathcal{L}}{\partial W_1} = \underbrace{\frac{\partial \mathcal{L}}{\partial z z_2^c}}_{zA_2^c - y_c} \underbrace{\frac{\partial z z_2^c}{\partial A_1}}_{W_2^T} \underbrace{\frac{\partial A_1}{\partial z z_1}}_{\sigma'(z z_1)} \underbrace{\frac{\partial z z_1}{\partial W_1}}_{A_0^T}$$

**Nota:** la semplificazione  $\partial \mathcal{L} / \partial z z_2 = zA_2^c - y_c$  vale solo usando la coppia  $\mathcal{L}_{\text{CCS}}$  e softmax.