

# LOGISTIC REGRESSION

Fin ora abbiamo usato l'MLP per risolvere problemi di regressione lineare. Possiamo usarlo anche per risolvere problemi di Regressione Logistica.

La Regressione Logistica cerca di risolvere un problema di classificazione. L'obiettivo è calcolare la probabilità:

$$P(y=1 | x)$$

Ovvero la probabilità che, per l'ingresso  $x$  appartenga alla classe positiva ( $y=1$ ). Per risolvere questo compito devo allenare i pesi  $w_1 \dots w_n$  associati alle feature d'ingresso così da decidere se "I dati stanno sotto o sopra l'equazione  $w_1 x_1 + \dots + w_n x_n$ ", e quindi associarli a una delle due classi (Positiva o Negativa).

Prendiamo un ingresso  $x \in \mathbb{R}^R$ , dato a un modello che dà in uscita  $\hat{y} \in [0, 1]$



Il nostro modello stima la probabilità di  $x$  di appartenere alla classe  $C$ :

$$\hat{y}(x) = P(x \in C)$$

Usando come funzione di attivazione la sigmoide:

$$\hat{y}(x, w) = \sigma(w^T \cdot \text{ext}(x))$$

## Threshold

Per decidere l'appartenenza a una classe, data la probabilità di appartenenza, bisogna settare una soglia oltre il quale si decide di assegnare l'ingresso alla classe positiva.

## BINARY CROSS-ENTROPY

Per allenare un modello di tipo logistic regression, si usa come Loss Function la BCS:

$$\mathcal{L}^{BCS} := -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

Questa misura la **dissimilarità** fra l'uscita predetta e quella attuale.

Questa è una loss convessa, per cui possiamo raggiungere l'ottimo:

$$\frac{\partial \mathcal{L}^{BCS}}{\partial W} = \frac{\partial \mathcal{L}^{BCS}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial W}$$

Dove:

$$1) \frac{\partial \mathcal{L}^{BCS}}{\partial \hat{y}} = -y \cdot \frac{1}{\hat{y}} - (1-y) \frac{1}{1-\hat{y}} \cdot (-1) = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$2) \frac{\partial \hat{y}}{\partial W} = \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial W} = \underbrace{\sigma(z)(1-\sigma(z))}_{\sigma(z)=\hat{y}} \cdot \underbrace{\frac{\partial W^T x}{\partial W}}_x = \hat{y}(1-\hat{y}) x$$



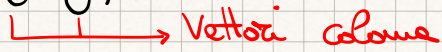
$$-\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}} = \frac{-y(1-\hat{y}) + \hat{y}(1-y)}{\hat{y}(1-\hat{y})} = \frac{-y + \cancel{\hat{y}} + \hat{y} - y\hat{y}}{\hat{y}(1-\hat{y})} = \frac{\hat{y} - y}{\hat{y}(1-\hat{y})}$$

Calcolando ora il gradiente della loss:

$$\frac{\partial \mathcal{L}^{BCE}}{\partial W} = \frac{\hat{y} - y}{\cancel{\hat{y}(1-\hat{y})}} \cancel{\hat{y}(1-\hat{y})} x = (\hat{y} - y) x$$

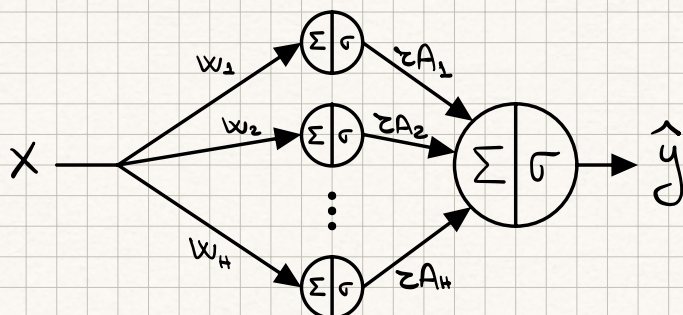
Se usiamo un mini-batch,  $x$  diventa una matrice  $X \in \mathbb{R}^{B \times (D+1)}$ :

$$\frac{\partial \mathcal{L}^{BCE}}{\partial W} = X^T (\hat{y} - y)$$



## PROBLEMA DI CLASSIFICAZIONE MULTICLASSE

L'idea è sempre dare una classe all'ingresso  $X$ :



Dove  $\hat{y}$  è la probabilità dell'ingresso di appartenere a una certa classe.

Nei problemi **multiclasse** abbiamo tanti neuroni nel layer di output quante sono le classi target.

Abbiamo quindi un vettore in uscita che ha i valori di probabilità di appartenenza a tutte le classi di training:

$$\hat{y} = \begin{bmatrix} 0.1 \\ 0.7 \\ 0.9 \\ 0.3 \end{bmatrix}$$

Per esempio se abbiamo un'immagine contenente più oggetti, probabilmente contenerà oggetti appartenenti alla classe 1 e 2 (0.7 e 0.9 di probabilità)

Supponendo un ingresso alla rete  $X \in \mathbb{R}^{3 \times D}$  avremo una uscita  $Y \in \mathbb{R}^{C \times B}$

$$X = \begin{bmatrix} \text{---} \end{bmatrix} \xrightarrow{\text{model}} Y = \begin{bmatrix} \text{---} \end{bmatrix}$$

B      C  
D      B



## FORWARD PROPAGATION

$$\begin{aligned}
 zA_0 &= X^T && \in \mathbb{R}^{D \times B} \\
 A_0 &= \text{ext}(zA_0) && \in \mathbb{R}^{(D+1) \times B} \\
 zZ_1 &= W_1 \cdot A_0 && : W_1 \in \mathbb{R}^{H \times (D+1)} \rightarrow zZ_1 \in \mathbb{R}^{H \times B} \\
 zA_1 &= \sigma(zZ_1) \\
 A_1 &= \text{ext}(zA_1) && : \in \mathbb{R}^{(H+1) \times B} \\
 zZ_2 &= W_2 \cdot A_1 && : W_2 \in \mathbb{R}^{C \times (H+1)} \rightarrow zZ_2 \in \mathbb{R}^{C \times B} \\
 zA_2 &= \sigma(zZ_2) \\
 Y &= zA_2 && : Y \in \mathbb{R}^{C \times B}
 \end{aligned}$$

```

LAB_05 > MLP_BCESIG_forward.m
1 % function [A0,rZ1,A1,rZ2,rA2] = MLP_BCESIG_forward(X, W1, W2)
2 function [rA2,A1,A0,rZ1] = MLP_BCESIG_forward(X, W1, W2)
3 % Compute the feedforward step
4
5 rA0 = X'; % rA0 is the "reduced A0" and it coincides with the transpose of the tall input matrix (the latter is BxD in size)
6 A0 = MLP_extend(rA0); % A0 = Ext(rA0). It is the "extended" version of rA0, obtained by it by adding a row of ones as its new first row
7 rZ1 = W1*A0; % rZ1 = \sum(W1,A0). It is the pre-activation at layer 1 (the hidden one)
8 rA1 = MLP_sigmoid(rZ1); % rA1 = \sigma(rZ1). It is the output of the first layer (the hidden one)
9 A1 = MLP_extend(rA1); % A1 = Ext(rA1). It is the extended version of rA1
10 rZ2 = W2*A1; % rZ2 = \sum(W2,A1). It is the pre-activation at layer 2 (the output one)
11 rA2 = MLP_sigmoid(rZ2); % rA2 is the output of the second layer (the output one). rA2 is a column vector if B=1, otherwise it is a matrix
12
13 end

```

## MATRICI DEI PESI

Per quando riguarda  $W_1$  e  $W_2$

$$W_1 = \begin{bmatrix} w_{11}^1 \\ \vdots \\ w_{1H}^1 \end{bmatrix} \quad \text{dove } \forall w_1^i \in \mathbb{R}^{D+1}$$

$$W_2 = \begin{bmatrix} w_{21}^2 \\ \vdots \\ w_{2C}^2 \end{bmatrix} \quad \text{dove } \forall w_2^i \in \mathbb{R}^{H+1}$$

## LOSS FUNCTION

Abbiamo visto che per problemi con classe binaria si usa la Binary Cross-Entropy:

$$\mathcal{L}^{BCS} = -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

Per problemi che hanno  $C > 2$  usiamo:

$$\mathcal{L}^{tot} = \sum_c \mathcal{L}_c^{BCS}$$

la sommatoria delle BCS di ogni classe. Ogni neurone di uscita rappresenta una classe:

$$\mathcal{L}_c^{BCS} = -y_c \log \hat{y}_c - (1-y_c) \log (1-\hat{y}_c)$$

Alleniamo ogni neurone indipendentemente e cerchiamo di minimizzare la somma di tutte le  $\mathcal{L}_c^{BCS}$ . Per farlo facciamo la **backpropagation**:

$$\frac{\partial \mathcal{L}}{\partial W_1}, \quad \frac{\partial \mathcal{L}}{\partial W_2}$$



1) Derivata su  $W_2$

$$\frac{\partial \mathcal{L}}{\partial W_2} = \underbrace{\frac{\partial \mathcal{L}}{\partial \hat{y}}}_{\text{red}} \underbrace{\frac{\partial \hat{y}}{\partial z_2}}_{\text{blue}} \underbrace{\frac{\partial z_2}{\partial W_2}}_{\text{green}}$$

Possiamo calcolare ora  $\partial \mathcal{L} / \partial z_2$  e  $\partial z_2 / \partial W_2$

$$1) \frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\hat{y} - y}{\hat{y}(1 - \hat{y})} \quad \left. \vphantom{\frac{\partial \mathcal{L}}{\partial \hat{y}}} \right\} \text{Calcolata sopra}$$

$$2) \frac{\partial \hat{y}}{\partial z_2} = \frac{\partial \sigma(z_2)}{\partial z_2} = \sigma(z_2)(1 - \sigma(z_2)) = \hat{y}(1 - \hat{y})$$

$$3) \frac{\partial z_2}{\partial W_2} = A_1^T$$

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\hat{y} - y}{\cancel{\hat{y}(1 - \hat{y})}} \cancel{\hat{y}(1 - \hat{y})} A_1^T = (\hat{y} - y) A_1^T$$

2) Derivata su  $W_1$

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