

Logistic Regression and Gradient Descent

Pattern Recognition and Machine Learning

Aniello Panariello, Lorenzo Bonicelli, Matteo Boschini

October 7th, 2022

University of Modena and Reggio Emilia

lo split degli elementi viene dato.



We are given a training set $\{X_i, Y_i\}_{i=1}^N$, with $X_i \in \mathbb{R}^m$ and $Y_i \in \{0, 1\}$ for each $i = 1, \dots, N$.

- N is the number of training examples;
- each example $X_i = \{x_i^{(1)}, \dots, x_i^{(m)}\}$ is a vector of m features;
- each label Y_i is either 0 or 1.

$f: X \rightarrow Y$ l.c. funzioni sul
training set.

We need to learn a function that maps X to Y such that "it works well on the training set".

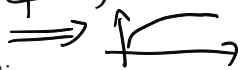
propos \Rightarrow vuoi f che preso x spiti la prob. che $y \in X$

We need to learn the parameters \mathbf{w} of a parametric function that maps X to Y such that “it works well on the training set”.

We need to learn the parameters \mathbf{w} of a parametric function that maps X to Y such that **some error is as low as possible on the training set.**

- Define **MSE** errors
- MINIMIZE ERROR; \rightarrow LOSS FUNCTION
- ω si impone con il GRADIENT DESCENT

Logistic regression classification rule

REGRESSION = PREDECE UN NUMERO
 Det. pl. $f: \dots$ voglio funzioni che ci
 mappi \Rightarrow 

The function for classification has the following form:

$$F(X_i, w) = \sigma(w^T \cdot X_i), \quad \text{where} \quad \sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

1. REGRESSION NON È MODULO DI REGRESSIONE MA DI CLASSIFICAZIONE

esempio in N dim

$$y_i = \sum_{j=1}^N a_j x_{ij} + b = w^T x_i + b$$

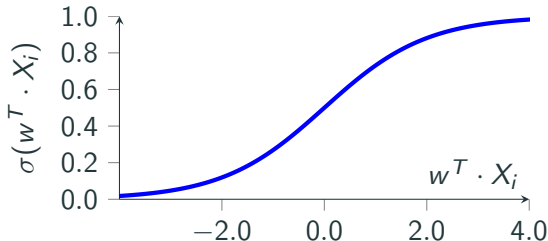
FUNZIONI DI
 REGRESSIONE
 CHE ASSIGNA
 UNO SCALARE $m \in \mathbb{R}$
 ME MAI VIGILANTE
 ASSOCIARE AD UNA
 PROBABILITÀ

The function for classification has the following form:

PARCHEGGIO DI QUESTA N-DIM.
IN UNA FUNZ. SIGMOID CHE
OSSERVISI (0=1)

$$F(X_i, w) = \sigma(w^T \cdot X_i), \quad \text{where} \quad \sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x} \quad \left(\begin{array}{l} \text{FUNCTION} \\ \text{LOGISTICA} \end{array} \right)$$

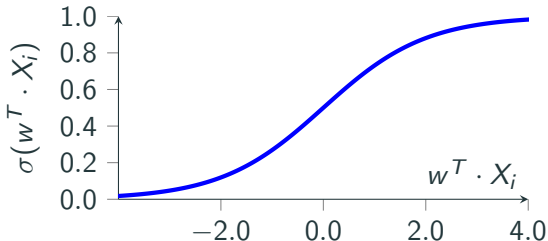
- $\sigma(x)$ is called **sigmoid function**;



The function for classification has the following form:

$$F(X_i, w) = \sigma(w^T \cdot X_i), \quad \text{where} \quad \sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

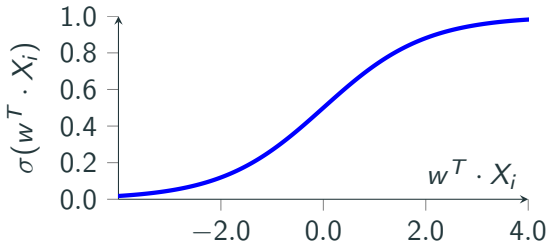
- $\sigma(x)$ is called **sigmoid function**;
- w is a vector in \mathbb{R}^m , and is called **weight vector**;



The function for classification has the following form:

$$F(X_i, w) = \sigma(w^T \cdot X_i), \quad \text{where} \quad \sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{1 + e^x}$$

- $\sigma(x)$ is called **sigmoid function**;
- w is a vector in \mathbb{R}^m , and is called **weight vector**;
- w is initialized randomly, but will improve as training goes.



is label y_i identifica distribuzione label su training set.

During training, we want to **minimize** the following function:

$$\mathcal{L}(w) = -\frac{1}{N} \sum_{i=1}^N [Y_i \log(\underbrace{F(X_i, w)}_{\substack{\text{è la mia predizione} \\ \uparrow \\ \text{per la classe 1}}}) + (1 - Y_i) \log(1 - F(X_i, w))]$$

serve modo x trovare w in modo che $\mathcal{L}(w)$ si minimizza

\Rightarrow GRADIENT DESCENT

$$L(\theta|\mathcal{D}) = \sum_{i=1}^N y_i \log(p_{y=1|x}) + (1-y_i) \log(p_{y=0|x})$$

$$p(y|x) = \frac{p(x, y)}{p(x)} \Rightarrow \text{unobserved}$$

$$W_i \rightarrow X_i, \quad W_i = \{\underbrace{\omega_1 \dots \omega_m}_{\text{wavy line}}\}$$

$$L(w)_{X_i} =$$

$$\frac{\partial L}{\partial w_i}$$

$$\Rightarrow w_i \text{ min } L(w)_{X_i}$$

Gradient descent is an iterative optimization algorithm for finding the minimum of a function. How? Take step proportional to the negative of the gradient of the function at the current point.

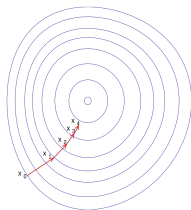


Figure 1: Gradient descent on a series of level sets

¹Credits for this slide: Andrea Palazzi https://github.com/ndrplz/deep_learning_lectures

If we consider a function $f(\theta)$, the **gradient descent update** can be expressed as:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} f(\theta) \quad (1)$$

for each parameter θ_j .

$$\theta_{j+1} = \theta_{j,old} - \alpha \frac{\partial}{\partial \theta_{j,old}} f(\theta)$$

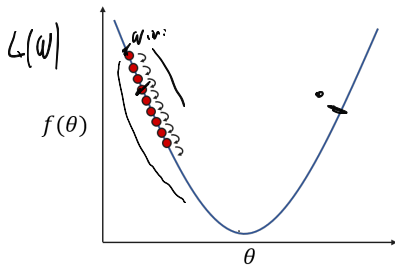
The size of the step is controlled by **learning rate** α .

¹Credits for this slide: Andrea Palazzi https://github.com/ndrplz/deep_learning_lectures

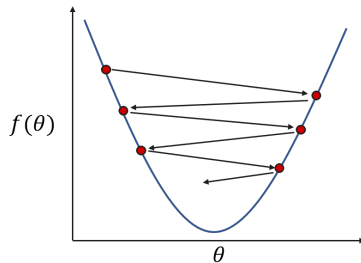
Gradient Descent for 1-d function $f(\theta)$.

¹Credits for this slide: Andrea Palazzi https://github.com/ndrplz/deep_learning_lectures

Choosing the the right **learning rate** α is essential to correctly proceed towards the minimum. A step *too small* could lead to an extremely *slow* convergence. If the step is *too big* the optimizer could *overshoot* the minimum or even *diverge*.



Learning Rate too small



Learning Rate too big

¹Credits for this slide: Andrea Palazzi https://github.com/ndrplz/deep_learning_lectures

Back to our problem. We need to take the derivative of this function w.r.t. w :

$$\begin{aligned}\mathcal{L}(w) &= -\frac{1}{N} \sum_{i=1}^N [Y_i \log(F(X_i, w)) + (1 - Y_i) \log(1 - F(X_i, w))] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i \log \left(\frac{e^{w^T \cdot X_i}}{1 + e^{w^T \cdot X_i}} \right) + (1 - Y_i) \log \left(\frac{1}{1 + e^{w^T \cdot X_i}} \right) \right] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i (w^T \cdot X_i) - Y_i \log \left(1 + e^{w^T \cdot X_i} \right) + (Y_i - 1) \log \left(1 + e^{w^T \cdot X_i} \right) \right] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i (w^T \cdot X_i) - \log \left(1 + e^{w^T \cdot X_i} \right) \right]\end{aligned}$$

Back to our problem. We need to take the derivative of this function w.r.t. w :

$$\mathcal{L}(w) = -\frac{1}{N} \sum_{i=1}^N \left[Y_i (w^T \cdot X_i) - \log \left(1 + e^{w^T \cdot X_i} \right) \right]$$

SOMMA COSTI
SE GRADIENTE È
NO FA LA MEDIA
SE $w = 1$

GRADIENTE X
↓ COSTO
Σ COSTO

$$\begin{aligned} \leftarrow \frac{\delta \mathcal{L}(w)}{\delta w_j} &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i x_i^{(j)} - \frac{e^{w^T \cdot X_i}}{1 + e^{w^T \cdot X_i}} x_i^{(j)} \right] \\ &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i - \frac{e^{w^T \cdot X_i}}{1 + e^{w^T \cdot X_i}} \right] x_i^{(j)} \\ &= -\frac{1}{N} \sum_{i=1}^N \left[Y_i - F(X_i, w) \right] x_i^{(j)} \end{aligned}$$

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

dove m è il
numero di
features

Back to our problem. We need to take the derivative of this function w.r.t. w :

il numero di parametri
 w dipende dalle
features degli N
sample X_i

$$\frac{\delta \mathcal{L}(w)}{\delta w} = \begin{bmatrix} \frac{\delta \mathcal{L}(w)}{\delta w_1} \\ \frac{\delta \mathcal{L}(w)}{\delta w_2} \\ \vdots \\ \frac{\delta \mathcal{L}(w)}{\delta w_m} \end{bmatrix} = \begin{bmatrix} -\frac{1}{N} \sum_{i=1}^N (Y_i - F(X_i, w)) x_i^{(1)} \\ -\frac{1}{N} \sum_{i=1}^N (Y_i - F(X_i, w)) x_i^{(2)} \\ \vdots \\ -\frac{1}{N} \sum_{i=1}^N (Y_i - F(X_i, w)) x_i^{(m)} \end{bmatrix} = -\frac{X^T \cdot (Y - F(X, w))}{N}$$

AVENDO GLI ELEMENTI
X ESEMPIO 0.5
IN $w_i, i = 1 \dots M$

We will update the vector w accordingly:

$$w \leftarrow w - \alpha \frac{\delta \mathcal{L}(w)}{\delta w}$$

Algorithm 1 pseudocode for training

- 1: $X, Y \leftarrow \text{load_training_data}()$
 - 2: set learning rate $\alpha \leadsto 10^{-3} \sim 10^{-4}$
 - 3: initialize w randomly
 - 4: **for** $e = 1$ to number_of_training_steps **do**
 - 5: compute the prediction according to the current weights $F(X, w) =$
 - 6: compute the loss function $\mathcal{L}(w) = \frac{1}{n} \sum_{i=1}^n \ell(w; x_i, y_i)$
 - 7: compute the derivative of the loss function w.r.t. weights $\frac{\delta \mathcal{L}(w)}{\delta w}$
 - 8: update the weight vector $w \leftarrow w - \alpha \frac{\delta \mathcal{L}(w)}{\delta w}$
 - 9: **end for**
-

$$\begin{bmatrix} F(x_1, w) \\ \vdots \\ F(x_n, w) \end{bmatrix}$$



- We want to predict if a character is alive or dead;
- (Some) of our features are:
 - male or female;
 - married or not;
 - number of deaths witnessed;
 - number of dead relatives;
 - ... and many more.