

Regressão Logística

Advanced Institute for Artificial Intelligence – AI2

<https://advancedinstitute.ai>

Existem problemas para os quais desejamos estimar a **classe** (rótulo) de uma determinada amostra com base no seu conjunto de dados de entrada, isto é, as **características** do seu problema.

Problemas de classificação são similares aos de regressão pois desejamos aprender uma função que, dada uma entrada, estime um valor de saída. A diferença é que nossa saída será mapeada para uma **probabilidade** (*soft classification*) ou **rótulo** (*hard classification*) da amostra.

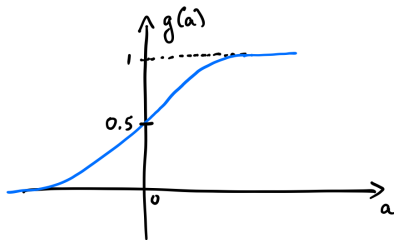
A técnica de Regressão Logística é uma das mais utilizadas na literatura, principalmente por ser simples e dar bons resultados em diversas situações. Ela tem esse nome pelo fato de estimar a probabilidade de uma dada amostra pertencer à uma classe em específico. Portanto, é uma técnica de classificação do tipo *soft*.

Definição do problema: seja um conjunto de dados $\mathcal{X} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_z, y_z)\}$ tal que $\mathbf{x}_i \in \mathbb{R}^{n+1}$ corresponde ao dado de entrada e $y_i \in [0, 1]$ denota o seu respectivo valor de saída. Temos, ainda, que \mathcal{X} pode ser **particionado** da seguinte forma: $\mathcal{X} = \mathcal{X}^1 \cup \mathcal{X}^2$, em que \mathcal{X}^1 e \mathcal{X}^2 denotam os conjuntos de dados de **treinamento** e **teste**, respectivamente. Nosso objetivo é, dado o conjunto de treinamento, aprender uma função $h : \mathbb{R}^{n+1} \rightarrow [0, 1]$ que consiga estimar a probabilidade de uma amostra pertencer à uma dada classe. Por que usar uma função logística (sigmoide)?

A Função Logística possui algumas propriedades interessantes com a seguinte formulação:

$$g(a) = \frac{1}{1 + e^{-a}}, \quad (1)$$

tal que $g(a) \in [0, 1]$.



Como funciona o Regressor Logístico? Dado que queremos obter uma saída $h_w(x) \in [0, 1]$, basta modificarmos nossa entrada na Equação 1 como segue:

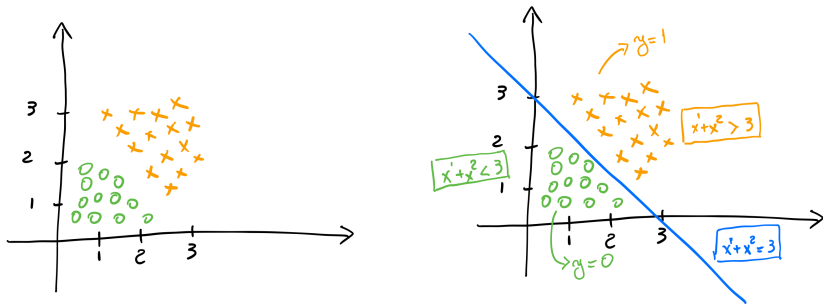
$$\begin{aligned} h_w(x) &= g(w^T x) \\ &= \frac{1}{1 + e^{-w^T x}}. \end{aligned} \quad (2)$$

Agora, o termo $w^T x$ é chamado de **função base**.

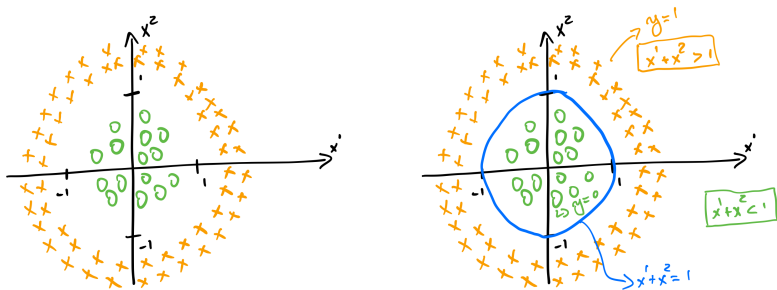


Na prática, queremos que $h_w(x) \geq 0.5$ quando $w^T x \geq 0$. De maneira análoga, temos que $h_w(x) < 0.5$ quando $w^T x < 0$.

Suponha que tenhamos a seguinte situação, em que $h_w(x) = g(w^T x)$ e $w = [-3 \ 1 \ 1]$. Queremos atribuir $y = 1$ se $w^T x > 0$, ou seja, se $-3 + x^1 + x^2 > 0 \Rightarrow x^1 + x^2 > 3$. De maneira análoga, queremos atribuir $y = 0$ se $w^T x < 0$, ou seja, se $-3 + x^1 + x^2 < 0 \Rightarrow x^1 + x^2 < 3$.



Suponha agora outra situação, em que $h_{\mathbf{w}}(\mathbf{x}) = g(w_0 + w_1x^1 + w_2x^2 + w_3(x^1)^2 + w_4(x^2)^2)$ e $\mathbf{w} = [-1 \ 0 \ 0 \ 1 \ 1]$. Queremos atribuir $y = 1$ se $-1 + (x^1)^2 + (x^2)^2 > 0 \Rightarrow (x^1)^2 + (x^2)^2 > 1$. De maneira análoga, queremos atribuir $y = 0$ se $1 + (x^1)^2 + (x^2)^2 < 0 \Rightarrow (x^1)^2 + (x^2)^2 < 1$.



Desta forma, podemos notar que, dependendo da função base escolhida, diferentes superfícies de separação podem ser obtidas. Usualmente, quanto maior o grau do polinômio, mais complexa será a sua superfície.

Assim sendo, temos as seguintes informações até o momento:

- Função hipótese: $h_{\mathbf{w}}(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}}$.
- Função base: $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ (padrão).
- Função de custo: ?

Por que não é interessante utilizar MSE como função de custo para o Regressor Logístico?

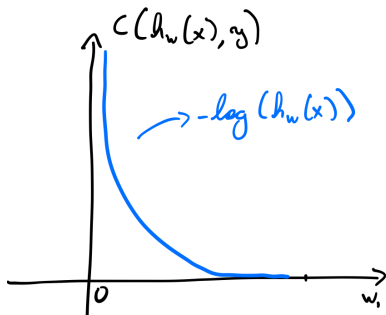
Temos dois problemas principais com MSE quando usamos a técnica de Regressão Logística:

- ❶ Suponha que o rótulo verdadeiro de uma amostra $x \in \mathcal{X}$ qualquer seja $y = 1$, e nosso classificador tenha resultado como saída $h_{\mathbf{w}}(x) = 0$. Neste caso (para $m = 1$), temos que $J(\mathbf{w}) = \frac{1}{2}(1 - 0)^2 = 0.5$. Note que essa é uma penalização muito pequena para um **erro** de classificação. Assim sendo, MSE não penaliza muito fortemente erros de classificação e pode levar a um aprendizado insuficiente.
- ❷ A função de custo MSE **não é convexa** para o Regressor Logístico (provado matematicamente).

Vamos, então, definir uma nova função de custo que seja convexa. Nossa função será dividida, inicialmente, em duas partes:

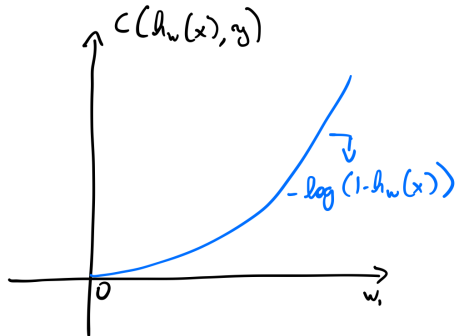
$$C(h_{\mathbf{w}}(\mathbf{x}), y) = \begin{cases} -\log(h_{\mathbf{w}}(\mathbf{x})) & \text{se } y = 1 \\ -\log(1 - h_{\mathbf{w}}(\mathbf{x})) & \text{se } y = 0. \end{cases} \quad (3)$$

Vamos, agora, analisar as duas situações separadamente. Como primeiro caso, temos a situação em que $y = 1$.



Neste caso, quando temos um erro, ou seja, $h_{\mathbf{w}}(\mathbf{x}) = 0$, então $C(h_{\mathbf{w}}(\mathbf{x}), y) = -\log(0) = \infty$. No caso de um acerto, ou seja, $h_{\mathbf{w}}(\mathbf{x}) = 1$, temos que $C(h_{\mathbf{w}}(\mathbf{x}), y) = -\log(1) = 0$.

A próxima situação ocorre quando $y = 0$.



Neste caso, quando temos um erro, ou seja, $h_w(x) = 1$, então $C(h_w(x), y) = -\log(1 - 1) = -\log(0) = \infty$. No caso de um acerto, ou seja, $h_w(x) = 0$, temos que $C(h_w(x), y) = -\log(1 - 0) = -\log(1) = 0$.

Podemos escrever a Equação 3 agrupando as duas situações descritas anteriormente:

$$C(h_{\mathbf{w}}(\mathbf{x}), y) = -y \log(h_{\mathbf{w}}(\mathbf{x})) - (1 - y) \log(1 - h_{\mathbf{w}}(\mathbf{x})). \quad (4)$$

Quando $y = 1$, termo vermelho é zerado e a equação acima torna-se a Equação 3 para este caso. Por outro lado, quando $y = 0$, termo azul é zerado e a equação acima torna-se a Equação 3 para este caso.

A nossa função de custo final para o Regressor Logístico é dada por:

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{m} \sum_{i=1}^m C(h_{\mathbf{w}}(\mathbf{x}_i), y_i) \\ &= \frac{1}{m} \sum_{i=1}^m [-y_i \log(h_{\mathbf{w}}(\mathbf{x}_i)) - (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))]. \end{aligned} \quad (5)$$

A função acima é também conhecida por **entropia cruzada binária**.

Basta, agora, aprendermos o conjunto de parâmetros \mathbf{w} de maneira semelhante à Regressão Linear, ou seja, podemos fazer uso da técnica de gradiente descendente. A derivada da função de custo apresentada na Equação 5 é dada como segue:

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m [(h_{\mathbf{w}}(\mathbf{x}_i) - y_i)x_i^j]. \quad (6)$$

Muito embora a formulação acima seja idêntica à Regressão Linear, devemos lembrar que a função hipótese ($h_{\mathbf{w}}(\mathbf{x})$) é diferente.

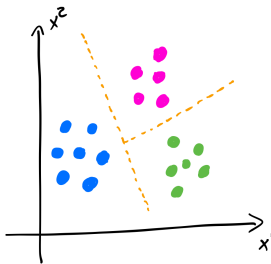
Desta forma, o algoritmo do gradiente descendente pode ser sumarizado da seguinte forma:

- ➊ Atribua valores aleatórios para w .
- ➋ Avalie a função de custo $J(w)$.
- ➌ Caso o **critério de parada tenha sido atingido**, vá para o passo 6.
- ➍ $w_j^{(t+1)} = w_j^{(t)} - \alpha \frac{1}{m} \sum_{i=1}^m [(h_w(x_i) - y_i)^2 x_i^j]$.
- ➎ Retorne ao Passo 2.
- ➏ Fim do algoritmo.

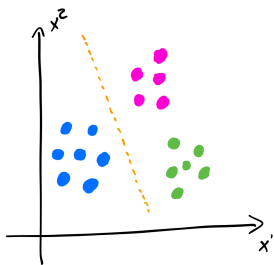
Note que os pesos em w precisam ser atualizados simultaneamente!

No entanto, note que o Regressor Logístico é, naturalmente, um **classificador binário**, isto é, $y_i \in \{0, 1\}$, $\forall i = 1, 2, \dots, z$. Desta forma, como podemos atuar em problemas que possuem **múltiplas classes**?

Agora, nossos rótulos podem ser representados da seguinte forma: $y_i \in \{0, 1, \dots, c - 1\}$, em que c corresponde ao número de classes. Suponha a seguinte situação em que $c = 3$.

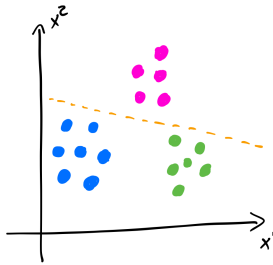


Existem algumas abordagens para tratar esse problema, tais como OVO (*one-versus-one*) e OVA (*one-versus-all*). Vamos considerar a abordagem OVA. Neste caso, para um problema de classificação com c classes, temos que criar c classificadores. Vejamos o exemplo abaixo em que $c = 3$.



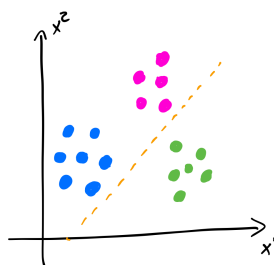
Azul versus todos

$$h_w^1(x)$$



Rosa versus todos

$$h_w^2(x)$$



Verde versus todos

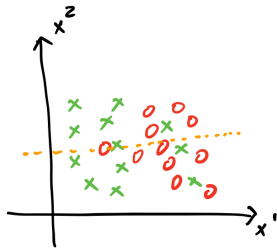
$$h_w^3(x)$$

Basta, então, treinarmos cada classificador $h_w^i(x)$, $\forall i = 0, 1, \dots, c - 1$ em \mathcal{X}^1 . Dada uma amostra $x \in \mathcal{X}^2$, seu rótulo y será aquele que obedece à seguinte regra de decisão:

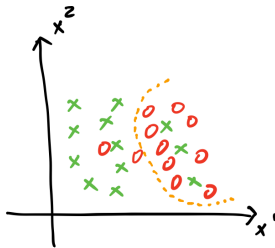
$$y = \arg \max_i \{h_w^i(x)\}. \quad (7)$$

Regularização

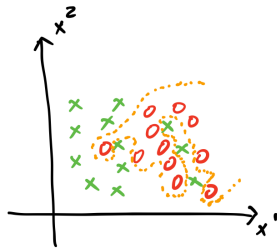
Como vimos anteriormente, a ideia da regularização é evitar com que a técnica fique muito **especializada** (viciada) no conjunto de treinamento e, portanto, não consiga generalizar muito bem no conjunto de teste. Vejamos as três principais situações que podem ocorrer em classificadores de padrões.



Subtreinamento



Bom treinamento



Supertreinamento

Com isso, podemos modificar a Equação 5 da função de custo como segue:

$$J(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m [-y_i \log(h_{\mathbf{w}}(\mathbf{x}_i)) - (1 - y_i) \log(1 - h_{\mathbf{w}}(\mathbf{x}_i))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2, \quad (8)$$

em que λ corresponde à taxa de regularização. Já as derivadas parciais da nova função de custo equivalem à derivadas da MSE no caso da regressão linear, como segue:

$$\frac{\partial J(\mathbf{w})}{\partial w_0} = \frac{1}{m} \sum_{i=1}^m (h_{\mathbf{w}}(x_i) - y_i), \quad (9)$$

e

$$\frac{\partial J(\mathbf{w})}{\partial w_j} = \frac{1}{m} \sum_{i=1}^m [(h_{\mathbf{w}}(x_i) - y_i)x_i^j] + \frac{\lambda}{m} w_j. \quad (10)$$

O algoritmo do gradiente descendente permanece o mesmo.