

# Mathematical model for logistic regression

Professor Montaz Ali

25th January 2021

Logistic Regression (LR) incorporates two models: Linear Regression (which is a linear model) and a non-linear transformation. Here we consider a linear model (LM) as being linear in its parameters and not in its variables:

$$L_m(x) = \sum_{i=0}^n a_i x_i = a^T x \quad (1)$$

where,  $x = (x_0, x_1, \dots, x_n)^T$ ,  $a_0 = b$  and  $x_0 = 1$ . The above LM is the linear combination of its inputs say  $x_i$ ,  $i = 1 \dots n$ . Inputs come from a known data set of size  $N$ . The  $j$ -th input will be denoted by  $x^j$ ,  $j = 1, 2, \dots, N$ . The process of obtaining the optimal parameters of  $L_m(x)$  is known as linear regression. Although we will not use linear classification here but it is needed to put things in perspective. Linear classification has two outcome 'yes' 'no' (1 or -1), i.e., binary classification. Two outcomes can be obtained from  $L_m(x)$  by its sign i.e  $\text{sign}(L_m(x)) = \text{sign}(a^T x) = f(x)$ . Here sign is acting as a non-linear transformation to convert  $a^T x$  to  $f(x)$ , where  $f(x) \in \{1, -1\}$ .

We can consider another non-linear transformation  $F(\cdot)$  such that  $F(a^T x) = f(x)$  where  $f(x) \in [0, 1]$ . Such a transformation is known as logistic function and is given by

$$F(x) = \frac{e^x}{1 + e^x}. \quad (2)$$

It now follows that

$$f(x) = \frac{e^{a^T x}}{1 + e^{a^T x}} \quad (3)$$

is a probability which is a generalization of binary classification and has many applications but we will use this to predict the heart attacks.

When the data is collected individuals, the person may or may not have a heart attack. We use a parameter  $\omega$  to denote this with  $\omega = 1$  denoting the heart attack and  $\omega = -1$  otherwise. When we build a model to predict  $\omega$ , our noise target will be as follows

$$\Pr(\omega/x) = \begin{cases} f(x) & \text{if } \omega = 1 \\ 1 - f(x) & \text{if } \omega = -1 \end{cases}.$$

It follows that  $\Pr(\omega/x) = F(\omega a^T x)$  since  $F(-a^T x) = 1 - F(a^T x)$ . We will use our data set

$$D = \{x^1, x^2, \dots, x^N\}$$

to construct  $L_m(x)$  such that  $\Pr(\omega/x)$  predicts our corresponding data  $(\omega/x)$  as closely as possible.

Consider the following likelihood function

$$L_f(a) = \prod_{i=1}^N \Pr(\omega^i/x^i) = \prod_{i=1}^N F(\omega^i a^T x^i) \quad (4)$$

We want to maximize (4) or equivalently to minimize some error measure. Since all are positive quantity in (4), we can write:

$$\max_a L_f(a) = \max_a \ln \prod_{i=1}^N F(\omega^i a^T x^i) = \sum_{i=1}^N \ln (F(\omega^i a^T x^i)).$$

We can then minimize the following

$$\min_a \bar{L}_f(a) = \min_a - \sum_{i=1}^N \ln (F(\omega^i a^T x^i)) = \min_a \sum_{i=1}^N \ln [1 + \exp(-\omega^i a^T x^i)] \quad (5)$$

which is a non-linear function and has to be optimized iteratively using gradient descent. The gradient of the function is given by

$$\nabla \bar{L}_f(a) = \sum_{i=1}^N \frac{-\omega^i x^i}{1 + e^{\omega^i a^T x^i}}. \quad (6)$$

The procedure is as follows:

1. Initialize  $a^k$ , for  $k=0$ , Calculate  $\nabla \bar{L}_f(a^k)$ ,
2. Set  $\alpha = 1$
3. Find  $a^{k+1} = a^k - \alpha \times \nabla \bar{L}_f(a^k)$ ,
4. Compare if  $L_f(a^{k+1}) < L_f(a^k)$  then set  $k = k + 1$  and go to 5 else set  $\alpha = \frac{1}{2} \times \alpha$  and go to 3.
5. Calculate  $\nabla \bar{L}_f(a^{k+1})$ . Stop if  $\|\nabla \bar{L}_f(a^{k+1})\|$  is small else go to 2.