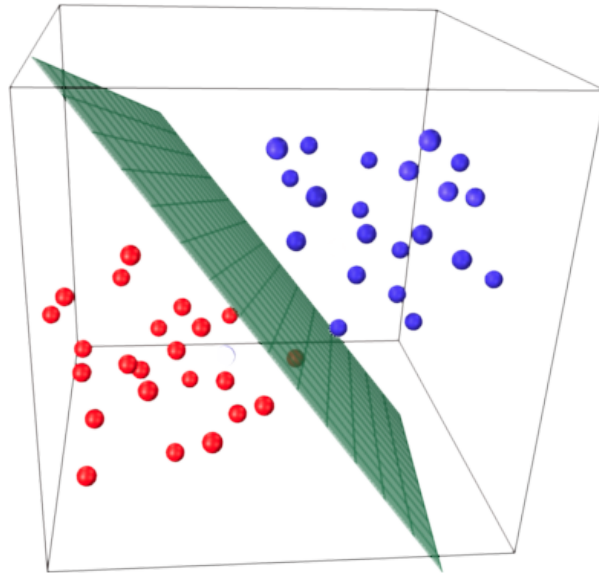


The main focus of logistic regression is to separate the two classes using a hyperplane and then get values as 1 or 0 approximately when put into sigmoid function.



The following is the hyperplane which separates the classes giving positive and negative values respectively. The plane gives positive for one side and negative for other as its a hyperplane of that particular space. The equation of the comes out to be:

$$W_0 + W_1x_1 + \dots + W_nx_n = 0$$

where  $W_0, W_1, \dots$  are weights which has to be trained using the gradient decent method, and  $x_1, x_2, \dots$  are features of the dependent variable.

Then we can represent the features and weights as vectors for easy representation and manipulation. Let:

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ \vdots \end{bmatrix} \text{ and } W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ \vdots \end{bmatrix} \text{ then } z(x) = W^T X + W_0$$

then we get the hypothesis function as:

$$h(x) = \sigma(z(x)) = \frac{1}{1 + e^{-(W^T X + W_0)}} = \frac{1}{1 + e^{-(W_0 + Wx_1 + \dots + W_nx_n)}}$$

from bernoulli's distribution equation we get:

$$B(y_p) = (y_p)^y (1 - y_p)^{1-y} \quad \text{where, } y_p \text{ is predicted value}$$

$$y \text{ is the actual value}$$

we define cost function as  $-\log(B(x))$  for one sample, as it penalises more even for small deviation from the actual value(can be inferred from the graph). Then for entire set it is,

$$\text{cost}(y_p) = -\frac{1}{m} \sum_{i=1}^m [(y) \log(y_p) + (1-y) \log(1-y_p)] \quad (1)$$

then,

$$\text{cost}(y_p) = \begin{cases} y = 1 & \sum \log(1 + e^{-z(x)}) \\ y = 0 & \sum \log(1 + e^{z(x)}) \end{cases}$$

then differentiating we get,

$$\frac{\partial \text{cost}(y_p)}{\partial W_0} = \frac{1}{m} \begin{cases} y = 1 & \sum \frac{-e^{-z(x)}}{(1 + e^{-z(x)})} = \sigma(z(x)) - 1 \\ y = 0 & \sum \frac{e^{z(x)}}{(1 + e^{z(x)})} = \sigma(z(x)) - 0 \end{cases}$$

$$\frac{\partial \text{cost}(y_p)}{\partial W_j} = \frac{1}{m} \begin{cases} y = 1 & \sum \frac{-x_j e^{-z(x)}}{(1 + e^{-z(x)})} = x_j(\sigma(z(x)) - 1) \\ y = 0 & \sum \frac{x_j e^{z(x)}}{(1 + e^{z(x)})} = x_j(\sigma(z(x)) - 0) \end{cases}$$

then we can write it just as ,

$$\frac{\partial \text{cost}(y_p)}{\partial W_0} = \frac{1}{m} \sum_{i=1}^m [h(X) - y_i]$$

$$\frac{\partial \text{cost}(y_p)}{\partial W_j} = \frac{1}{m} \sum_{i=1}^m x_j^i [h(X) - y_i]$$

In vector form the approximate algorithm would be,

$$W_0 = W_0 - \frac{\alpha}{m} \sum_{i=1}^m [\sigma(W^T X + W_0) - y_i]$$

$$W_j = W_j - \frac{\alpha}{m} \sum_{i=1}^m x_j^i [\sigma(W^T X + W_0) - y_i]$$

where  $y_i \in \{0, 1\}$

where  $j = 1, 2, \dots, n$

Alpha is a constant which determine the rate of covergence of the process, it shouldn't be too small or too big. If its too small the convergence rate will be very slow and on other hand if its too big then it may never hit the minimum and instead deviate far away.

These equations should be iterated several times to get the desired minimum-error weights. .