

Binary Classification:

In this our dependent variable (y) is a discrete variable where it takes two values i.e. $y \in \{0,1\}$ where 0 represents negative class and 1 represents a positive class.

Examples:

Email : Spam or not ?

Online Transactions : Fraud or not ?

Tumor : malignant or Benign ?

Threshold classifier output : We use linear regression and map all predictions greater than 0.5 as a 1 and all less than 0.5 as a 0. However, this method doesn't work well because classification is not actually a linear function.

What if we use linear reg for classification problem ?

If we use linear regression model then our predicted output might be <0 or sometimes it might be >1 , both the cases are not suitable for classification problem.

We can use classification algorithms such as **logistic Regression** for the binary classification.

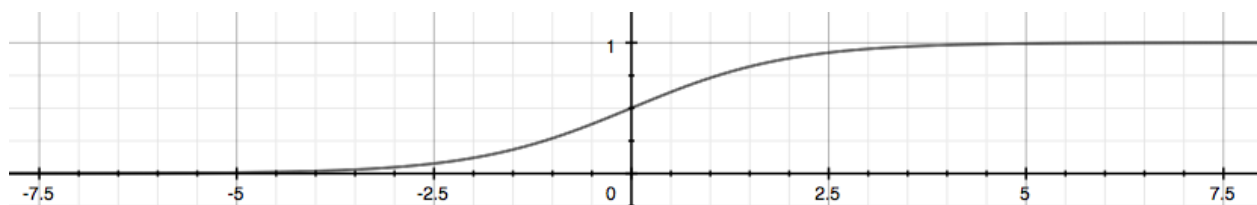
Logistic Regression :

Our goal is to have our hypothesis function between 0 to 1, i.e. $0 \leq h_{\theta}(x) \leq 1$.

θ

For Logistic regression we modify our hypothesis function as $h_{\theta}(x) = g(h_{\theta}(x))$, i.e. $h_{\theta}(\theta^T X) = \frac{1}{1+e^{-\theta^T X}}$

where $g(z)$ is a logistic/sigmoid function, where $g(z) = \frac{1}{1+e^{-z}}$, the graph looks like below.



Interpretation of Hypothesis Output:

$h_{\theta}(x)$ = estimated probability that $y=1$ on input x .

Example of cancer prediction : $h_{\theta}(x) = 0.7 \rightarrow$ it tells that 70% chance of tumor being malignant tumor.

Another way of notation is $h_{\theta}(x) = P(y=1/x; \theta) \rightarrow$ "probability that $y=1$, given x , parameterized by θ "

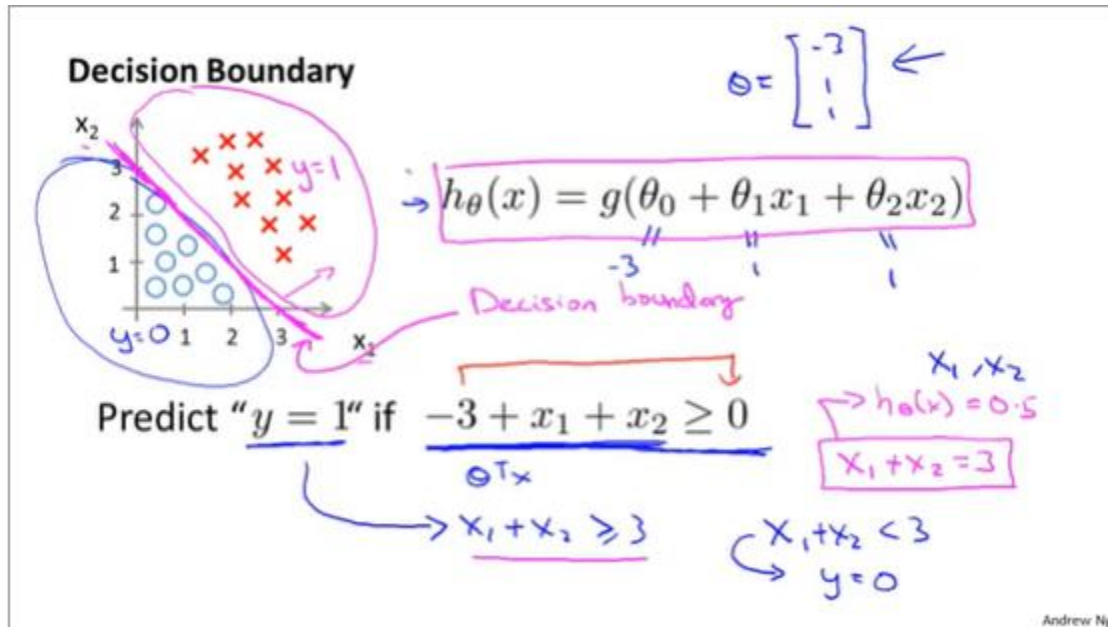
Also, $P(y=1/x; \theta) + P(y=0/x; \theta) = 1$.

Decision Boundary :

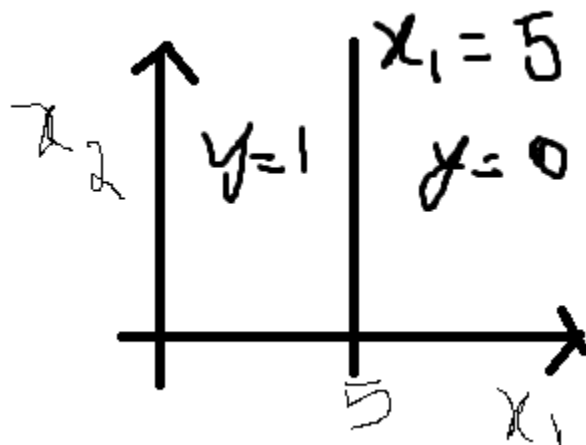
From sigmoid graph, we can say that $g(z) \geq 0.5$ when $z \geq 0$, suppose if $z=0 \rightarrow g(z) = 1/2$

That is

- for $y=1$, $\rightarrow g(\theta^T X) \geq 0.5$ i.e. $(\theta^T X) \geq 0$. Similarly
- for $y=0 \rightarrow g(\theta^T X) < 0.5$ i.e. $(\theta^T X) < 0$.

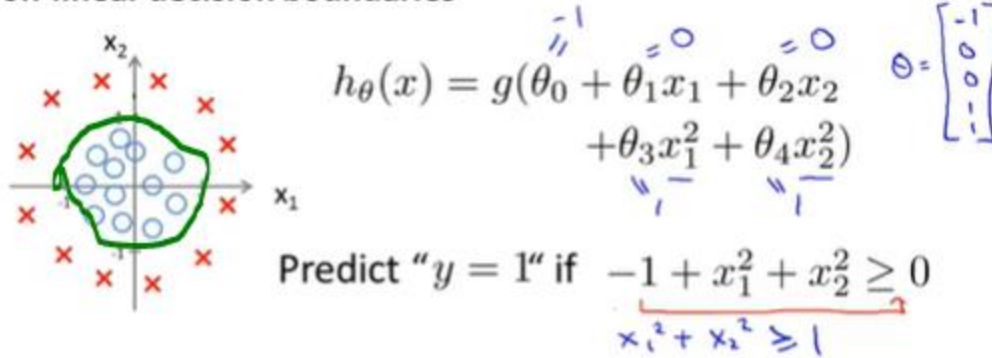


- Consider logistic regression with two features x_1 and x_2 . Suppose $\theta_0 = 5$, $\theta_1 = -1$ and $\theta_2 = 0$, so that $h_\theta(x) = g(5 - x_1)$. Which of these shows the decision boundary of $h_\theta(x)$?
- For $y=1$, $\theta^T X$ i.e. $5 - x_1 \geq 0 \rightarrow x_1 \leq 5$ is the required decision Boundary equation.



- The decision boundaries might not always be linear, it can also be non-linear as well.

Non-linear decision boundaries



Cost Function for Logistic Regression:

Training set : $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$, m examples and we have our $X \in \begin{bmatrix} x^0 \\ \vdots \\ x^n \end{bmatrix}$ i.e. a (n+1) dimensional vector, $x^0 = 1$ and $y \in \{0, 1\}$.

Our hypothesis/Objective function, $h_{\theta}(\theta^T X) = \frac{1}{1 + e^{-\theta^T X}}$

Now how to choose parameters θ ?

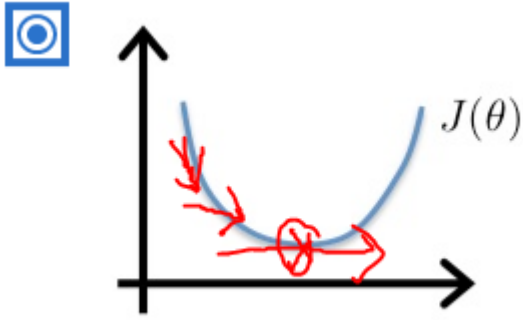
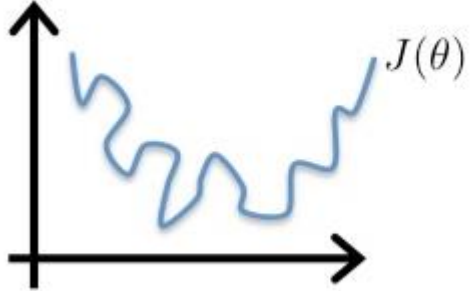
Cost function :

$$\text{For Linear Regression cost function} \rightarrow J(\theta) = \frac{1}{2m} * \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$

$$= \frac{1}{m} * \sum_{i=1}^m \text{cost}(h_{\theta}(x^i), y^i)$$

Where $\rightarrow \text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$, for logistic regression this cost function will be a non-convex (it will have a lot of local minimums).

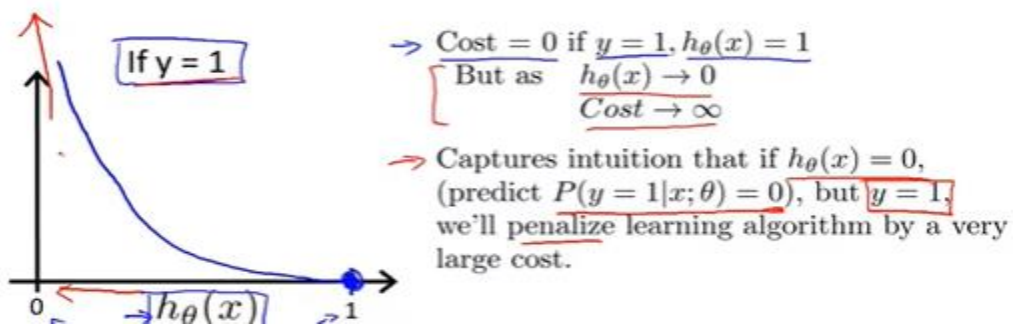
Because of squared cost function and the sigmoid function's non-linearity the graph will be non-convex, so it will be impossible to find global minimum.

Convex	Non-Convex
	

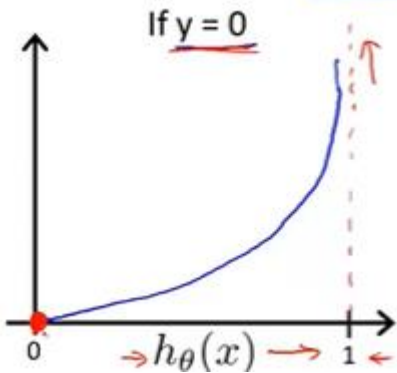
So, we use a different cost function for logistic regression.

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & , \text{ if } y = 1 \\ -\log(1 - h_{\theta}(x)) & , \text{ if } y = 0 \end{cases}$$

Case 1: If $y = 1$



Case 2 : If $y = 0$.

	<p>If our hypothesis predicts as 1 but the value of y is zero, then we get large cost value to be penalized that shows that our hypothesis is predicting the wrong class label.</p> <p>If hypothesis = 0 and $y = 0$ then the cost value is zero.</p>
---	---

Why $\log 1 = 0$?

- A logarithm is defined as the exponent or power to which a base must be raised to get some new number. It is a convenient approach to express large numbers.

Example : $10^2 = 100$ which can be written as $2 = \log_{10} 100$, Common or Briggsian logarithm is the logarithm with base 10.

So, we know anything to the power of zero is one $\rightarrow 10^0 = 1$
 $= \log 10^0$
 $= 0 \log_{10} 10$
 $= 0. \text{ (since } \log_{10} 10 = 1 \text{)}$

Why anything to power 0 is one i.e. $x^0 = 1$?

- We know $a^m / a^n = a^{m-n}$
- if $m=n$:
 - $\rightarrow a^m / a^m = a^{m-m}$
 - $\rightarrow 1 = a^0$
 - $\rightarrow a^0 = 1$

also remember $0^0 \neq 1$

Simplified cost function and gradient descent:

$$J(\theta) = \frac{1}{2m} * \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2$$
$$= \frac{1}{m} * \sum_{i=1}^m \text{cost}(h_{\theta}(x^i), y^i)$$

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)), & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)), & \text{if } y = 0 \end{cases}$$

Note : $y=0$ or 1 always

$$\rightarrow \text{Cost} = -y \log(h_{\theta}(x)) - (1 - y) \log(1 - h_{\theta}(x))$$

Logistic regression Cost Function : (based on maximum likelihood)

$$J(\theta) = \frac{1}{m} * \sum_{i=1}^m \text{cost}(h_{\theta}(x^i), y^i)$$
$$= -\frac{1}{m} * \left[\sum_{i=1}^m y \log(h_{\theta}(x)) + (1 - y) \log(1 - h_{\theta}(x)) \right]$$

To fit parameter θ : that minimize $J(\theta)$

Gradient descent :

Goal to minimize $J(\theta)$

Repeat {

$$\theta_j = \theta_j - \alpha * \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x^i) - y^i) x_j^i)$$

}

$$J(\theta) = -\frac{1}{m} * \left[\sum_{i=1}^m y \log(h_{\theta}(x)) + (1 - y) \log(1 - h_{\theta}(x)) \right]$$

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_j} &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) \frac{\partial}{\partial \theta_j} g(\theta^T x) \\ &= \left(y \frac{1}{g(\theta^T x)} - (1 - y) \frac{1}{1 - g(\theta^T x)} \right) g(\theta^T x) (1 - g(\theta^T x)) \frac{\partial}{\partial \theta_j} \theta^T x \\ &= (y(1 - g(\theta^T x)) - (1 - y)g(\theta^T x)) x_j \\ &= (y - h_{\theta}(x)) x_j \end{aligned}$$

Vectorized Implementation :

$$\rightarrow h_{\theta}(x) = g(\theta^T x)$$

$$\rightarrow \text{cost} = 1/m * [-y^T * \log(h_{\theta}(x)) - (1 - y)^T * \log(1 - h_{\theta}(x))]$$

$$\rightarrow \theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - y)$$

Out Scope for now : Usage of 1x1 convolution ?

We can shrink the number of channels, based on the number of 1x1 convolutional filters.