# [Summary]Logistic Regression 1

Don't be confused by the name "Logistic Regression"; it is named that way for historical reasons and is actually an approach to **classification problems, not regression problems**.

## Binary Classification

Instead of our output vector y being a continuous range of values, it will only be 0 or 1.

$$y \in \{\, 0,\ 1 \,\}$$

Where 0 is usually taken as the "negative class" and 1 as the "positive class", but you are free to assign any representation to it.

We're only doing two classes for now, called a "Binary Classification Problem."
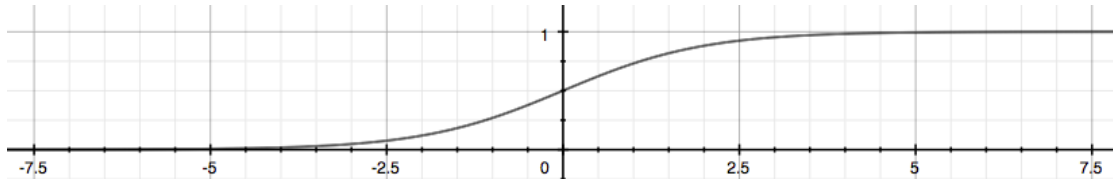
### Hypothesis Representation

Our hypothesis should satisfy:

$$0 \leq h_\theta(x) \leq 1$$

Our new form uses the "Sigmoid Function," also called the "Logistic Function":

$$h_\theta(x) = g(\theta^T x)$$
$$z = \theta^T x$$
$$g(z) = \frac{1}{1 + e^{-z}}$$

$h_\theta$ will give us the **probability** that our output is 1. For example, $h_\theta(x) = 0.7$ gives us the probability of 70% that our output is 1.

$$h_\theta(x) = P(y = 1|x; \theta) = 1 - P(y = 0|x; \theta)$$
$$P(y = 0|x; \theta) + P(y = 1|x; \theta) = 1$$

## Decision Boundary

In order to get our discrete 0 or 1 classification, we can translate the output of the hypothesis function as follow:

$$h_\theta(x) \geq 0.5 \rightarrow y = 1$$
$$h_\theta(x) < 0.5 \rightarrow y = 0$$

The way our logistic function g behaves is that when its input is greater than or equal to zero, its output is greater than or equal to 0.5:

$$g(z) \geq 0.5$$
$$when\ z \geq 0$$

Remember:

$$z = 0,\ e^0 = 1 \Rightarrow g(z) = \frac{1}{2}$$
$$z \rightarrow \infty,\ e^{-\infty} \rightarrow 0 \Rightarrow g(z) = 1$$
$$z \rightarrow -\infty,\ e^{\infty} \rightarrow \infty \Rightarrow g(z) = 0$$

So if our input to g is $\theta^T X$, then that means:

$$h_\theta(x) = g(\theta^T x) \geq 0.5$$
$$when\ \theta^T x \geq 0$$

From these statements we can now say:

$$\theta^T x \geq 0 \Rightarrow y = 1$$
$$\theta^T x < 0 \Rightarrow y = 0$$

The **decision boundary** is the line that separates the area where y = 0 and where y = 1. It is created by our hypothesis function.

The input to the sigmoid function g(z)(e.g. $\theta^T X$) doesn't need to be linear, and could be a funcation that describes a circle (e.g. $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$) or any shape to fit our data.
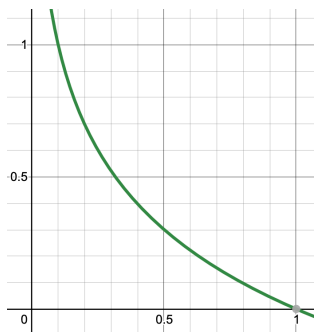
## Cost Function

We cannot use the same cost function that we use for linear regression because the Logistic Function will cause the output to be wavy, causing many local optima. In other words, it will not be a convex function.

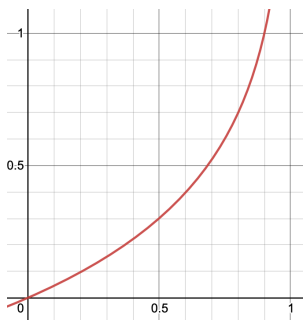Instead, our cost function for logistic regression looks like:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} Cost(h_\theta(x^i),\ y^i)$$
$$Cost(h_\theta(x),\ y) = -\log(h_\theta(x))\ if\ y = 1$$
$$Cost(h_\theta(x),\ y) = -\log(1 - h_\theta(x))\ if\ y = 0$$

-1

-0.5

0    0.5    1

*If y =1*

-1

0.5

0    0.5    1

*If y = 0*

The more our hypothesis is off from y, the larger the cost function output. If our hypothesis is equal to y, then our cost is 0:

$$Cost(h_\theta(x), y) = 0 \; if \; h_\theta(x) = y$$
$$Cost(h_\theta(x), \; y) \to \infty \; if \; y = 0 \; and \; h_\theta(x) \to 1$$
$$Cost(h_\theta(x), \; y) \to \infty \; if \; y = 1 \; and \; h_\theta(x) \to 0$$

## Simplified Cost Function and Gradient Descent

We can compress our cost function's two conditional cases into one case:

$$Cost(h_\theta(x), \; y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$$

We can fully write out our entire cost function as follows:

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

A vectorized implementation is:

$$h = g(X\theta)$$

$$J(\theta) = \frac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$

## Gradient Descent

Remember that the general form of gradient descent is:

Repear {

$$\theta_j := \theta_j - \alpha \frac{\delta}{\delta\theta_j} J(\theta)$$

}

We can work out the derivative part using calculus to get:

Repeat {

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

}

Notice that this algorithm is identical to the one we used in linear regression. We still have to simultaneously update all values in theta.

## Partial derivative of $J(\theta)$

First calculate derivative of sigmoid function (it will be useful while finding partial derivative of J(θ)):

$$\sigma(x)' = \left(\frac{1}{1+e^{-x}}\right)'$$

$$= \frac{-(1+e^{-x})'}{(1+e^{-x})^2}$$

$$= \frac{-1' - (e^{-x})'}{(1+e^{-x})^2}$$

$$= \frac{0 - (-x)'(e^{-x})}{(1+e^{-x})^2}$$

$$= \frac{-(-1)(e^{-x})}{(1+e^{-x})^2}$$

$$= \frac{e^{-x}}{(1+e^{-x})^2}$$

$$= \left(\frac{1}{1+e^{-x}}\right)\left(\frac{e^{-x}}{1+e^{-x}}\right)$$

$$= \sigma(x)\left(\frac{+1 - 1 + e^{-x}}{1+e^{-x}}\right)$$

$$= \sigma(x)\left(\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}}\right)$$

$$= \sigma(x)(1 - \sigma(x))$$

Now we are ready to find out resulting partial derivative:

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

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

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

$$= -\frac{1}{m}\sum_{i=1}^{m}\left[\frac{y^i \frac{\delta}{\delta\theta_j}\sigma(\theta^T x^i)}{h_\theta(x^i)} + \frac{(1-y^i)\frac{\delta}{\delta\theta_j}\left(1 - \sigma(\theta^T x^i)\right)}{1 - h_\theta(x^i)}\right]$$

$$= -\frac{1}{m}\sum_{i=1}^{m}\left[\frac{y^i\sigma(\theta^T x^i)\left(1 - \sigma(\theta^T x^i)\right)\frac{\delta}{\delta\theta_j}\theta^T x^i}{h_\theta(x^i)} + \frac{-(1-y^i)\sigma(\theta^T x^i)\left(1 - \sigma(\theta^T x^i)\right)\frac{\delta}{\delta\theta_j}\theta^T x^i}{1 - h_\theta(x^i)}\right]$$

$$= -\frac{1}{m}\sum_{i=1}^{m}\left[\frac{y^i h_\theta(x^i)\left(1 - h_\theta(x^i)\right)\frac{\delta}{\delta\theta_j}\theta^T x^i}{h_\theta(x^i)} - \frac{(1-y^i)h_\theta(x^i)\left(1 - h_\theta(x^i)\right)\frac{\delta}{\delta\theta_j}\theta^T x^i}{1 - h_\theta(x^i)}\right]$$

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

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

$$= -\frac{1}{m}\sum_{i=1}^{m}[y^i - y^i h_\theta(x^i) - h_\theta(x^i) + y^i h_\theta(x^i)]x_j^i$$

$$= -\frac{1}{m}\sum_{i=1}^{m}[y^i - h_\theta(x^i)]x_j^i$$

$$= \frac{1}{m}\sum_{i=1}^{m}[h_\theta(x^i) - y^i]x_j^i$$