

Email - spam / not spam

Online transaction - Fraudulent / legal

Tumour - malignant / benign

Logistic Regression Model

Want $0 \leq h_{\theta}(x) \leq 1$

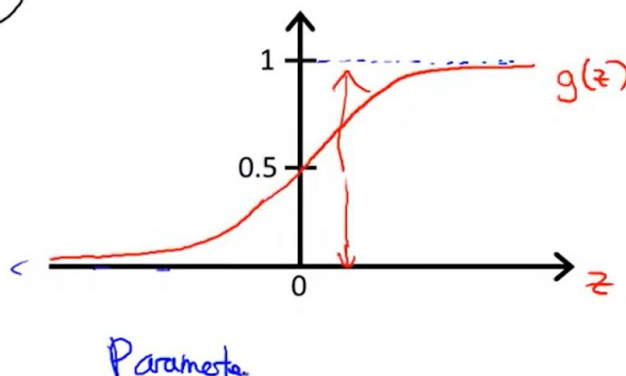
$$h_{\theta}(x) = g(\theta^T x)$$

$$\rightarrow g(z) = \frac{1}{1 + e^{-z}}$$

$\theta^T x$

Sigmoid function
Logistic function

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



Interpretation of Hypothesis Output

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

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$$h_{\theta}(x) = 0.7$$

$$y = 1$$

Tell patient that 70% chance of tumor being malignant

$$h_{\theta}(x) = P(y=1|x;\theta)$$

$$y = 0 \text{ or } 1$$

"probability that $y = 1$, given x ,
parameterized by θ "

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

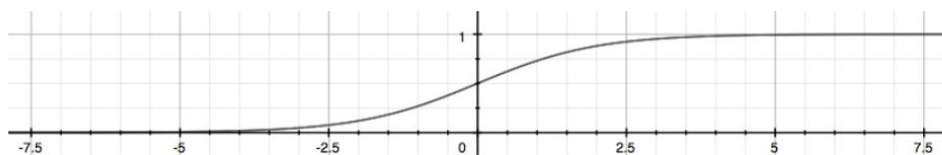
Our new form uses the "Sigmoid Function," also called the "Logistic Function":

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

$$z = \theta^T x$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

The following image shows us what the sigmoid function looks like:

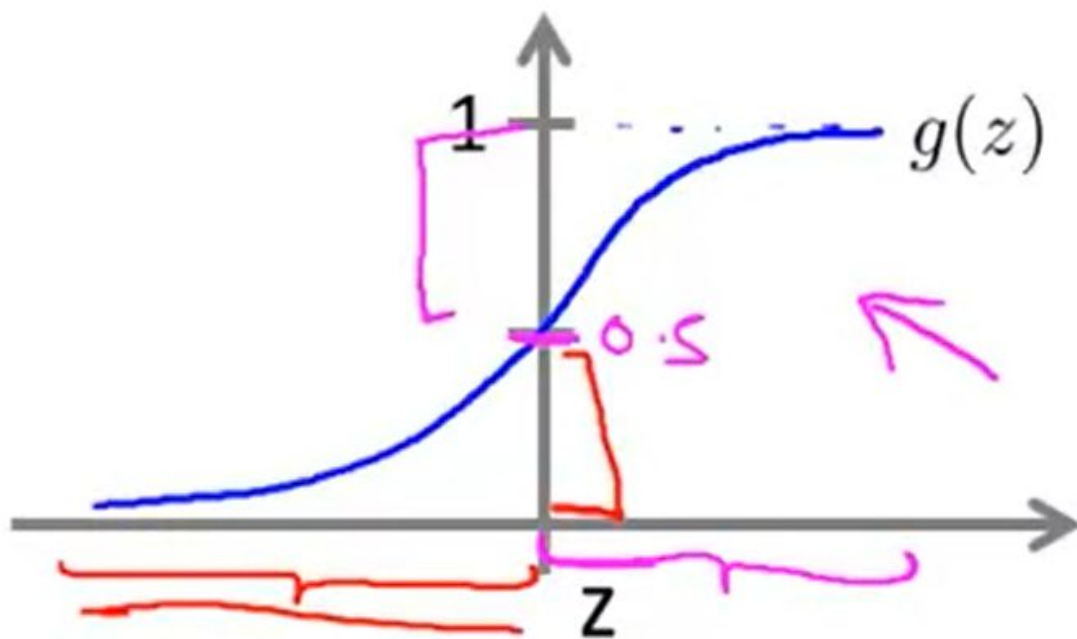


The function $g(z)$, shown here, maps any real number to the $(0, 1)$ interval, making it useful for transforming an arbitrary-valued function into a function better suited for classification.

$h_{\theta}(x)$ will give us the **probability** that our output is 1. For example, $h_{\theta}(x) = 0.7$ gives us a probability of 70% that our output is 1. Our probability that our prediction is 0 is just the complement of our probability that it is 1 (e.g. if probability that it is 1 is 70%, then the probability that it is 0 is 30%).

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

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



$$g(z) \geq 0.5$$

when $z \geq 0$

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

whenever $\theta^T x \geq 0$

\uparrow
 z

$z) < 0.5$

$Y = 1$ when $\theta^T X$ is greater than or equal to zero

$Y = 0$ when $\theta^T X$ is lesser than zero

Logistic regression

$$\rightarrow h_{\theta}(x) = g(\theta^T x) = P(y=1|x;\theta)$$

$$\rightarrow g(z) = \frac{1}{1+e^{-z}}$$

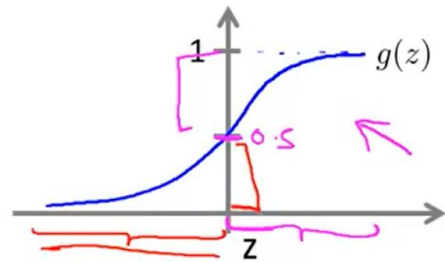
Suppose predict " $y = 1$ " if $h_{\theta}(x) \geq 0.5$

$$\rightarrow \theta^T x \geq 0$$

predict " $y = 0$ " if $h_{\theta}(x) < 0.5$

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

$$\rightarrow \theta^T x < 0$$



$$g(z) \geq 0.5$$

when $z \geq 0$

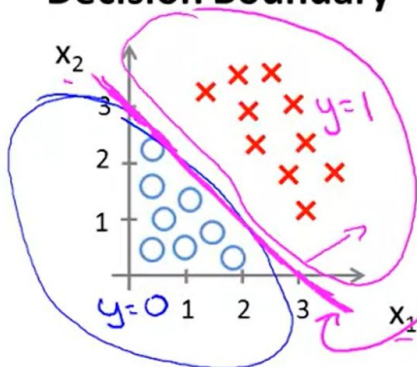
$$h_{\theta}(x) = g(\theta^T x) \geq 0.5$$

whenever $\theta^T x \geq 0$

\uparrow
 z

$$g(z) < 0.5$$

Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

$$\rightarrow h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Decision boundary

Predict " $y = 1$ " if $-3 + x_1 + x_2 \geq 0$

$$\theta^T x$$

$$\rightarrow x_1 + x_2 \geq 3$$

x_1, x_2

$$\rightarrow h_{\theta}(x) = 0.5$$

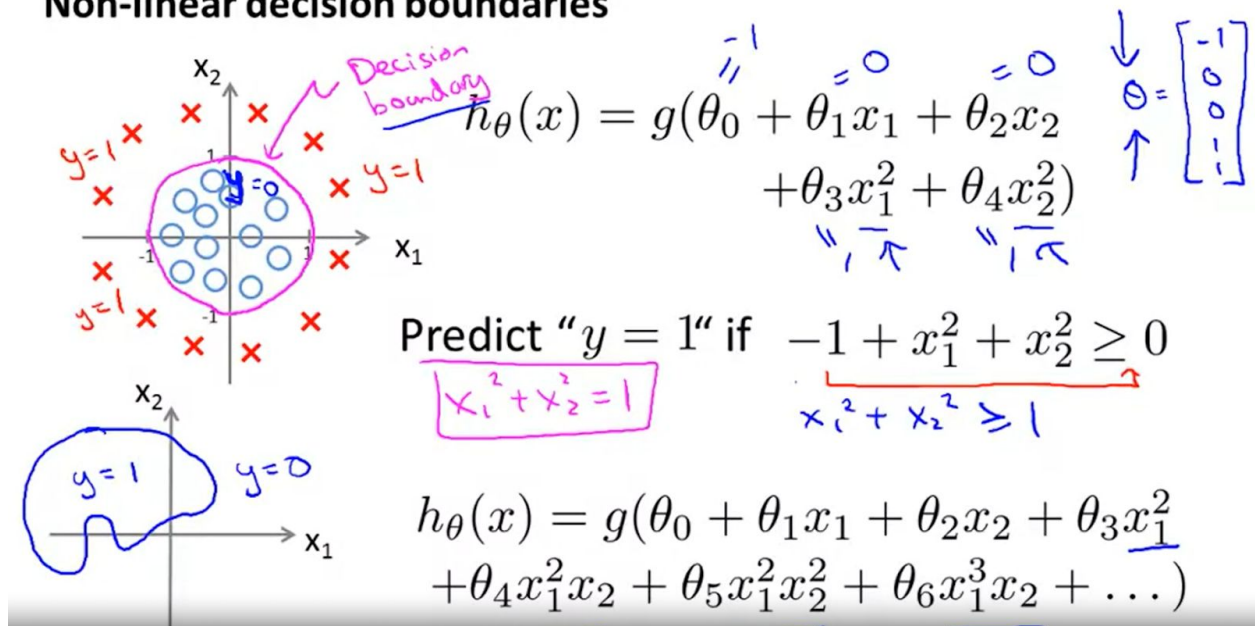
$$x_1 + x_2 = 3$$

$$\rightarrow x_1 + x_2 < 3$$

$$y = 0$$

Decision boundary is a property of the hypotheses(parameter theta) and not of the training set.
By using higher order polynomials we can get complex decision boundaries

Non-linear decision boundaries



Cost function

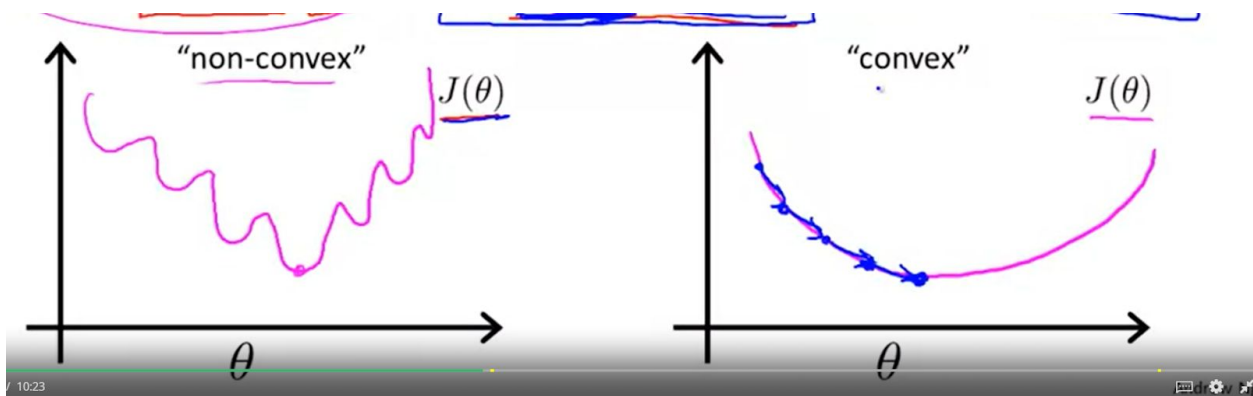
→ Linear regression: $J(\theta) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

$\text{cost}(h_{\theta}(x^{(i)}), y)$

→ $\text{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$

The sigmoid function is a non linear function so it hinders the convex curve which we get normally from linear regression.

If we use normal cost function for logistic regression we would encounter the non convex function. That is, we might not reach global minima. So we use a different cost function.



Cost function

→ ~~Linear~~ regression:
logistic

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

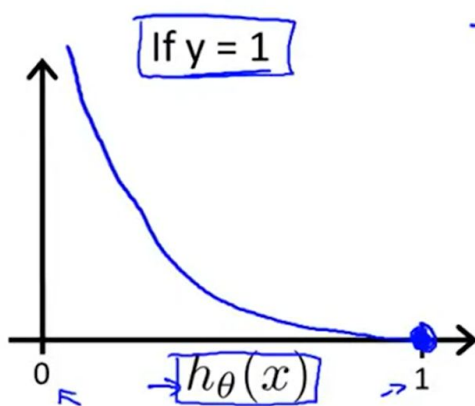
$\rightarrow \text{cost}(h_{\theta}(x^{(i)}), y)$

$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

$\frac{1}{1 + e^{-\sigma_{\theta}^T x}}$

Logistic regression cost function

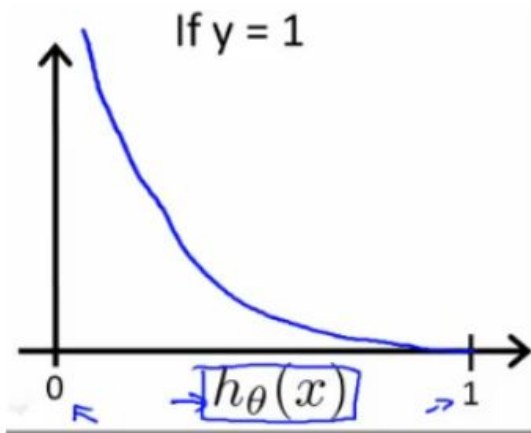
$$\text{Cost}(\underline{h_{\theta}(x)}, y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



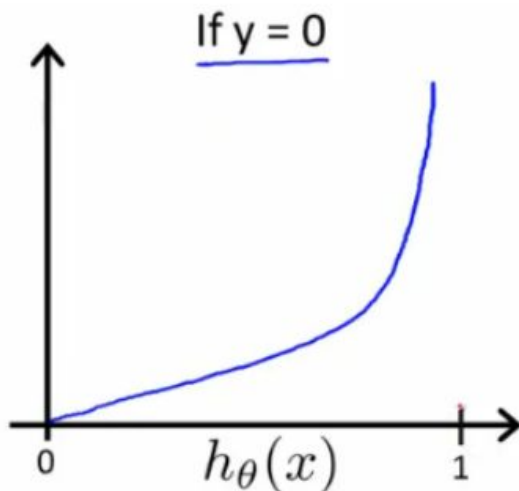
→ Cost = 0 if $y = 1, h_\theta(x) = 1$
 But as $h_\theta(x) \rightarrow 0$
 $Cost \rightarrow \infty$

Captures intuition that if $h_\theta(x) = 0$,
 (predict $P(y = 1|x; \theta) = 0$), but $y = 1$,
 we'll penalize learning algorithm by a very
 large cost.

When $y = 1$, we get the following plot for $J(\theta)$ vs $h_\theta(x)$:



Similarly, when $y = 0$, we get the following plot for $J(\theta)$ vs $h_\theta(x)$:



If our correct answer 'y' is 0, then the cost function will be 0 if our hypothesis function also outputs 0. If our hypothesis approaches 1, then the cost function will approach infinity.

If our correct answer 'y' is 1, then the cost function will be 0 if our hypothesis function outputs 1. If our hypothesis approaches 0, then the cost function will approach infinity.

Note that writing the cost function in this way guarantees that $J(\theta)$ is convex for logistic regression.

Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

Note: $y = 0$ or 1 always

This can be rewritten as

Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

Note: $y = 0$ or 1 always

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

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

If $y=0$: $\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x))$

Logistic regression cost function

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

To fit parameters θ :

$$\min_{\theta} J(\theta) \quad \text{Get } \underline{\theta}$$

To make a prediction given new x :

$$\text{Output } \underline{h_{\theta}(x)} = \frac{1}{1 + e^{-\theta^T x}}$$

$$\underline{p(y=1 | x; \theta)}$$

Gradient Descent

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

Want $\min_{\theta} J(\theta)$:

Repeat {

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

}

(simultaneously update all θ_j)

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

Gradient Descent

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

Want $\min_{\theta} J(\theta)$:

Repeat {

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

(simultaneously update all θ_j)

Gradient descent looks exactly similar for linear regression and logistic regression. Just the definition for the hypothesis has changed.

Gradient Descent

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

Want $\min_{\theta} J(\theta)$:

Repeat {

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

(simultaneously update all θ_j)

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$

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

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

Algorithm looks identical to linear regression!

Optimization algorithm

Given θ , we have code that can compute

$$\begin{aligned} & - J(\theta) \leftarrow \\ & - \frac{\partial}{\partial \theta_j} J(\theta) \leftarrow \quad (\text{for } j = 0, 1, \dots, n) \end{aligned}$$

Optimization algorithms:

- - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

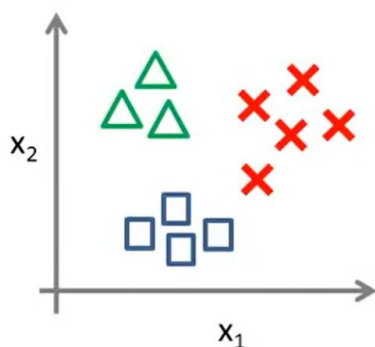
Disadvantages:

- More complex

IF there are multiple variables unlike 2 in binary classification

Multi class classification

One-vs-all (one-vs-rest):

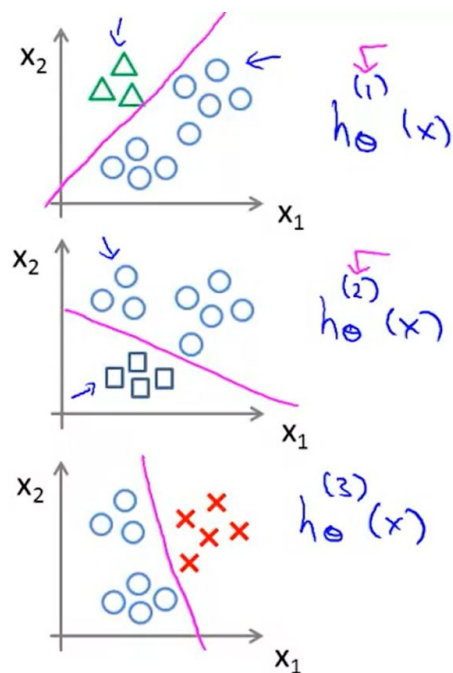


Class 1: \triangle ←

Class 2: \square ←

Class 3: \times ←

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3)$$



Advanced optimization

Press Esc to exit full screen

```
function [jVal, gradient] = costFunction(theta)
    jVal = [code to compute  $J(\theta)$ ];

    gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];

    gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];

    gradient(3) = [code to compute  $\frac{\partial}{\partial \theta_2} J(\theta)$ ];
    ⋮
    gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];
```