

# Machine Learning Stanford course - Summary

Gi Wah Dávalos Loo 

April 13, 2018

## 1 Introduction

### 1.1 Definition

Modern definition by Tom Mitchell:

”A computer program is said to learn from experience **E** with respect to some class of tasks **T** and performance measure **P**, if its performance at tasks in **T**, as measured by **P**, improves with experience **E**.”

### 1.2 Types

In general, any machine learning problem can be assigned to one of two broad classifications: Supervised and Unsupervised learning.

#### 1.2.1 Supervised Learning

We are given a data set and already know what our correct output should look like, having the idea that there is a relationship between the input and the output. Supervised learning problems are categorized into:

1. *Regression*: tries to predict results within a continuous output, meaning that we are trying to map input variables to some continuous function.
2. *Classification*: tries to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories.

Regression	Classification
Linear Regression	Logistic Regression

### 1.2.2 Unsupervised Learning

Allows us to approach problems with little or no idea what our results should look like. We can derive structure from data where we don't necessarily know the effect of the variables.

We can derive this structure by clustering the data based on relationships among the variables in the data.

With unsupervised learning there is no feedback based on the prediction results.

## 2 Linear Regression

The main idea of linear regression is to predict some continuous value from input ( $x_i$ ) using a function that has a linear behaviour ( $\theta_0x_0 + \theta_1x_1 + \dots + \theta_nx_n = \theta^T x$ ). The **goal** is to find the most accurate coefficients for the linear function we must reduce the error as much as possible, in other words, we gotta minimize the cost function.

### 2.1 Types

#### 2.1.1 Univariable

When we have only 1 feature. The linear function for univariable regression should look like this:

$$h_{\theta}(x) = \theta_0 + \theta_1x \quad (1)$$

#### 2.1.2 Multivariable

When we have more than 1 feature. The multivariable linear regression function should look like this:

$$h_{\theta}(x) = \theta_0x_0 + \theta_1x_1 + \theta_2x_2 + \dots + \theta_nx_n \quad (2)$$

### 2.2 Cost Function

We can measure the accuracy of our hypotheses function by using a cost function. This takes an average difference (actually a fancier version of an average) of all the results of the hypotheses with inputs from  $x$ 's and the actual output  $y$ 's.

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

## 2.3 Notation (Matrixes and Vector approaches)

In order to make operations more code/cpu "friendly" we represent our set of data and variables on matrixes and vectors.

$x_j^{(i)}$  = value of feature  $j$  in the  $i^{th}$  training example

$x^{(i)}$  = the input (features) of the  $i^{th}$  training example

$m$  = The number of training examples

$n$  = The number of features

$$\theta_{(n \times 1)} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \quad x_{(n \times 1)} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \quad X_{(m \times n)} = \begin{bmatrix} x_0^1 & x_1^1 & x_2^1 \dots x_n^1 \\ x_0^2 & x_1^2 & x_2^2 \dots x_n^2 \\ x_0^3 & x_1^3 & x_2^3 \dots x_n^3 \\ \vdots & \vdots & \vdots \\ x_0^m & x_1^m & x_2^m \dots x_n^m \end{bmatrix} \quad y_{(m \times 1)} = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$\text{definition format: } h_{(\theta)} = [\theta_0 \quad \theta_1 \dots \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} = \theta^T x$$

$$\text{solution format: } h_{(\theta)} = \begin{bmatrix} x_0^1 & x_1^1 & x_2^1 \dots x_n^1 \\ x_0^2 & x_1^2 & x_2^2 \dots x_n^2 \\ x_0^3 & x_1^3 & x_2^3 \dots x_n^3 \\ \vdots & \vdots & \vdots \\ x_0^m & x_1^m & x_2^m \dots x_n^m \end{bmatrix} \begin{bmatrix} \theta_0 & \theta_1 \dots \theta_n \end{bmatrix} = X\theta \quad (3)$$

$$\text{cost function: } J(\theta) = \frac{1}{2m} (X\theta - y)^2 \quad (4)$$

**Note 1:**  $x_0$  is a "hack" and it's value is always 1. It's for making both  $(\theta$  and  $x)$  same size in order to multiply them.

**Note 2:** The univariable and multivariable  $h_{(\theta)}$  can be solved using  $X\theta$

## 2.4 How to obtain $\theta_j$

The idea is to minimize the **Cost Dunction** (minimize the error). In order to do that, we need to build an equation by comparing the derivative of J to zero.

$$\min(J(\theta)) = \frac{\partial J(\theta)}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})x_j^{(i)} = 0$$

We got 2 ways to approach this problem:

1. *Gradient Descent*
2. *Normal Equation*

### 2.4.1 Gradient Descent

It is an iterative algorithm that, with a fixed learning rate ( $\alpha$ ), updates all values of  $\theta$  simultaneously and decreases the Cost Function. The learning rate ( $\alpha$ ) is adjusted by try/error. The initial values to test should be on a log-scale, at multiplicative steps of about 3 times the previous value (i.e., 0.3, 0.1, 0.03, 0.01 and so on).

repeat until converge {

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

}

Since it's an iterative algorithm we want the steps to be really small in order to make the computation time a lot faster. We care about the relationship between entries, not the actual values; that's why we tend to make the mean zero. For that we got 2 methods:

1. *Feature Scaling*: Involves dividing the input values by the range (i.e. the maximum value minus the minimum value) of the input variable, resulting in a new range of just 1.
2. *Mean normalization*: Involves subtracting the average value for an input variable from the values for that input variable resulting in a new average value for the input variable of just zero.

$$x_i = \frac{x_i - \mu_i}{s_i}$$

Where  $\mu_i$  is the average of all the values for feature ( $i$ ) and  $s_i$  is the range of values (max - min), or  $s_i$  is the standard deviation.

Note that dividing by the range, or dividing by the standard deviation, give different results.

### 2.4.2 Normal Equation

This is the analytical way to tackle this problem. We actually solve the equation  $\frac{\partial J(\theta)}{\partial \theta_j} = 0$ . There is no need to do feature scaling with the normal equation.

$$\theta = (X^T X)^{-1} X^T y$$

### 2.4.3 Comparison

The main differences between this two approaches:

Gradient Descent	Normal Equation
Need to choose alpha	No need to choose alpha
Needs many iterations	No need to iterate
$O(kn^2)$	$O(n^3)$ , need to calculate inverse of $X^T X$
Works well when $n$ is large	Slow if $n$ is very large

## 3 Logistic Regression

The main idea of logistic regression is to predict from a set of discrete binary values (ex: have cancer or not  $[0, 1]$ ). In order to accomplish that, we use a linear function (Decision boundary, can be non-linear) to map all predictions greater than 0.5 as 1 and all less than 0.5 as a 0. This **prediction** is done by a **sigmoid** ( $g(x)$ ) function that normalizes the output of the **Decision boundary** ( $\theta^T x$ ) to just 1 or 0.

So in conclusion, there is a **hypotheses** ( $h_\theta(x)$ ) function that is the sigmoid ( $g(x)$ ) function applied to the Decision boundary; and the **goal** is to find the most accurate parameters for the Decision boundary ( $\theta^T x$ ).

### 3.1 Components

#### 3.1.1 Sigmoid function

Also called *Logistic function*, 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. The function has this form:

$$g(z) = \frac{1}{1 + e^{-z}} \quad (5)$$

And looks like:

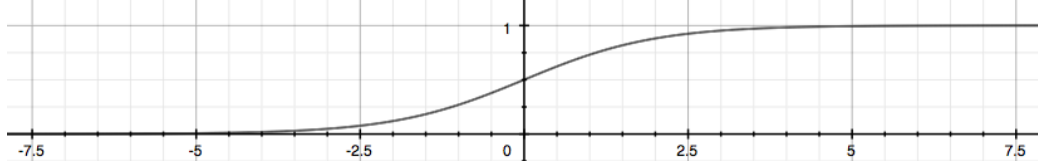


Figure 1: Sigmoid function

### 3.1.2 Hypotheses function

This hypotheses function will predict, using some input data ( $x_i$ ) and a set of precomputed coefficients ( $\theta$ ), a number between 0 and 1.

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

In order to achieve this, we use the sigmoid function to map our output to 0 - 1 interval on the decision boundary. The function has this form:

$$h_{\theta}(x) = g(\theta^T x) \quad (6)$$

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

The result of this function can be interpreted as *the probability that the 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%).

### 3.1.3 Decision boundary

This function describes the line (can be polynomial) that determines the limit of the features in order to predict wether 0 or 1. The representation for the decision boundary is:

$$\theta^T x$$

In order to visualize it, we need to separate the ones that are  $y = 0$  and the ones that are  $y = 1$ , then plot all features. But all axis must be the features ( $x^{(j)}$ ) displayed separately the  $0$ 's from the  $1$ 's  
And it can have many forms, also polynomial.

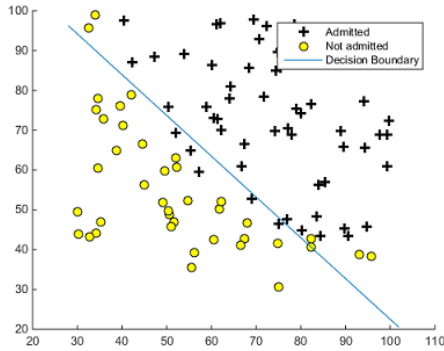


Figure 2: Linear Decision boundary

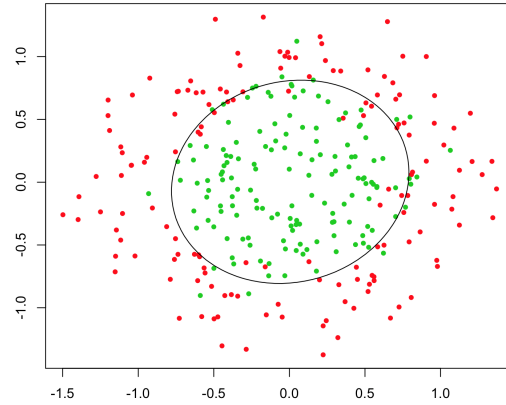


Figure 3: Cuadratic Decision boundary

We can interpret the output of the hypotheses function as:

$$\begin{aligned} h_{\theta}(x) &\geq 0.5 \rightarrow y = 1 \\ h_{\theta}(x) &< 0.5 \rightarrow y = 0 \end{aligned}$$

If we observe Figure 1 we can observe that in order to get  $g(z) \geq 0.5$  we need  $z$  to be greater than 0, in other words:

$$\begin{aligned} h_{\theta}(x) = g(\theta^T x) &\geq 0.5 \rightarrow y = 1 \\ \text{when } \theta^T x &\geq 0 \end{aligned}$$

Remember:

$$\begin{aligned} z = 0, e^0 = 1 &\Rightarrow g(z) = 1/2 \\ z \rightarrow \infty, e^{\infty} &\rightarrow 0 \Rightarrow g(z) = 1 \\ z \rightarrow -\infty, e^{-\infty} &\rightarrow 0 \Rightarrow g(z) = 0 \end{aligned}$$