

Week 3 Lecture Notes

ML: Logistic Regression

Now we are switching from regression problems to **classification problems**. 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.

math(3)

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

We're using binary classifiers for now called a "Binary Classification Problem."

Our method is to use linear regression and hope all predictions greater than 0.5 are 1 and all less than 0.5 are 0. This method doesn't work well because classification is not actually a linear function.

Hypothesis Representation

Our hypothesis should satisfy:

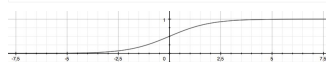
$$0 \leq h_{\theta}(x) \leq 1$$

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

$$h_{\theta}(x) = \text{sigmoid}(z)$$

$$z = \theta^T x$$

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



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. The mapping with interactive plot of sigmoid function: <https://www.khanacademy.org/machine-learning/a/interactive-plot-of-sigmoid-function/a/interactive-plot-of-sigmoid-function>

We start with our old hypothesis (linear regression), except that we want to restrict the range to 0 and 1. This is accomplished by plugging $\theta^T x$ into the logistic function.

h_{θ} 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.

$$h_{\theta}(x) = \text{Pr}(y = 1|x, \theta) = 1 - \text{Pr}(y = 0|x, \theta)$$

$$\text{Pr}(y = 0|x, \theta) = \text{Pr}(y = 1|x, \theta) = 1$$

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%).

Decision Boundary

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

$$h_{\theta}(x) \geq 0.5 \rightarrow y = 1$$

$$h_{\theta}(x) < 0.5 \rightarrow y = 0$$

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

$$g(z) \geq 0.5$$

$$\text{when } z \geq 0$$

Remember:

$$z = \theta_0 x^0 + \dots + \theta_n x^n = z$$

$$z = \theta_0 x^0 + \dots + \theta_n x^n = z$$

$$z = \theta_0 x^0 + \dots + \theta_n x^n = z$$

So if our input is given $\theta^T x$, then that means:

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

$$\text{when } \theta^T x \geq 0$$

From these statements we can see that:

$$\theta^T x \geq 0 \rightarrow y = 1$$

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

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

Example:

$$\theta = -1$$

$$\theta = 0$$

$$y = 1 \text{ if } z = (-1)x_0 + 0x_1 \geq 0$$

$$z = -x_0 \geq 0$$

$$-x_0 \geq 0$$

$$x_0 \leq 0$$

In this case, our decision boundary is a straight vertical line placed on the graph where $x_0 = 0$, and everything to the left of that divides $y = 1$, while everything to the right divides $y = 0$.

Again, the input to the sigmoid function g (e.g. $\theta^T x$) doesn't need to be linear, and could be a function that describes a circle (e.g. $z = A_0 + A_1 x_1^2 + A_2 x_2^2$) or any shape in \mathbb{R}^n our data.

Cost Function

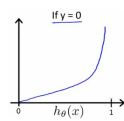
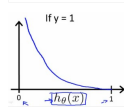
We cannot use the same cost function that we use for linear regression because the Logistic Function will cause the outputs to be many, causing many local optima. In other words, it will not be a convex function.

Instead, our cost function for logistic regression looks like:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

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

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



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.

$$\text{Cost}(h_{\theta}(x), y) = 0 \text{ if } h_{\theta}(x) = y$$

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

$$\text{Cost}(h_{\theta}(x), y) = -\log(1 - h_{\theta}(x)) \text{ if } y = 0 \text{ and } h_{\theta}(x) < 1$$

If our current 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 current 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.

Simplified Cost Function and Gradient Descent

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

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

Notice that when y is equal to 1, then the second term $(1 - y) \log(1 - h_{\theta}(x))$ will be zero and will not affect the result. If y is equal to 0, then the first term $-y \log(h_{\theta}(x))$ will be zero and will not affect the result.

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

$$J(\theta) = -\frac{1}{2} \sum_{i=1}^n [y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))]$$

A vectorized implementation is:

$$J = \text{cost}(\theta)$$

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

Gradient Descent

Remember that the general form of gradient descent is:

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

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

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

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

A vectorized implementation is:

$$\theta := \theta - \alpha X^T (y - \hat{y})$$

Partial derivative of $J(\theta)$

First calculate derivative of sigmoid function (it will be useful while finding partial derivative of $J(\theta)$)

$$g'(z) = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \left(\frac{1 + e^{-z} - 1}{1 + e^{-z}} \right) = \frac{1}{1 + e^{-z}} \left(\frac{e^{-z}}{1 + e^{-z}} \right) = \frac{e^{-z}}{(1 + e^{-z})^2}$$

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

Now we are ready to find out our partial derivative:

[illegible]

The vectorized versions

$$\nabla J(\theta) = \frac{1}{n} \cdot X^T \cdot (g(X \cdot \theta) - \bar{y})$$

Advanced Optimization

"Conjugate gradient", "BFGS", and "LBFGS" are more sophisticated, faster ways to optimize 0-1ss; can be used instead of gradient descent. A. Ng suggests not to write these more sophisticated algorithms yourself (unless you are an expert in numerical computing) but use the libraries issued as they're already tested and highly optimized. Octave provides them.

We first need to provide a function that evaluates the following two functions for a given input value θ :

$$\frac{\partial}{\partial \theta} J(\theta)$$

We can write a single function that returns both of these

```
1 function [J(theta), gradient] = costFunction(theta)
2     J(theta) = 0; % cost to compute J(theta) ...
3     gradient = []; % code to compute derivative of J(theta) ...
4 end
```

Then we can use Octave's fminunc optimization algorithm along with the optimoptions function that creates an options structure containing the options we want to use in fminunc. Using the value for MaxIter should be an integer, not a floating point, value as the default is 50.

A simpler approach

A "bias" feature is simply a way to move the "best fit" learned vector to better fit the data. For example, consider a learning problem with a single feature X_1 . The formula without the bias feature is just $\text{bias}(X_1) = y$. This is graphed as a line that always passes through the origin, with slope y/x . The b_0 term allows the line to pass through a different point on the y -axis. The self-bias strategy gives a better fit, but all bias fits pass go through the origin (is it right?)

Jon Orton