

## Logistic Regression

Based on a chapter by Chris Piech

**Logistic regression** is a classification algorithm<sup>1</sup> that works by trying to learn a function that approximates  $P(Y|X)$ . It makes the central assumption that  $P(Y|X)$  can be approximated as a sigmoid function applied to a linear combination of input features. It is particularly important to learn because logistic regression is the basic building block of artificial neural networks.

Mathematically, for a single training data point  $(\mathbf{x}, y)$ , logistic regression assumes:

$$P(Y = 1 | \mathbf{X} = \mathbf{x}) = \sigma(z) \text{ where } z = \theta_0 + \sum_{i=1}^m \theta_i x_i$$

This assumption is often written in the equivalent forms:

$$\begin{aligned} P(Y = 1 | \mathbf{X} = \mathbf{x}) &= \sigma(\theta^T \mathbf{x}) && \text{where we always set } x_0 \text{ to be } 1 \\ P(Y = 0 | \mathbf{X} = \mathbf{x}) &= 1 - \sigma(\theta^T \mathbf{x}) && \text{by total law of probability} \end{aligned}$$

Using these equations for probability of  $Y | X$  we can create an algorithm that selects values of theta that maximize that probability for all data. I am first going to state the log probability function and partial derivatives with respect to theta. Then later we will (a) show an algorithm that can chose optimal values of theta and (b) show how the equations were derived.

An important thing to realize is that: given the best values for the parameters  $(\theta)$ , logistic regression often can do a great job of estimating the probability of different class labels. However, given bad, or even random, values of  $\theta$  it does a poor job. The amount of “intelligence” that your logistic regression machine learning algorithm has depends on how good its values of  $\theta$  are.

### Notation

Before we get started I want to make sure that we are all on the same page with respect to notation. In logistic regression,  $\theta$  is a vector of parameters of length  $m$  and we are going to learn the values of those parameters based off of  $n$  training examples. The number of parameters should be equal to the number of features of each data point (see section 1).

Two pieces of notation that we use often in logistic regression that you may not be familiar with:

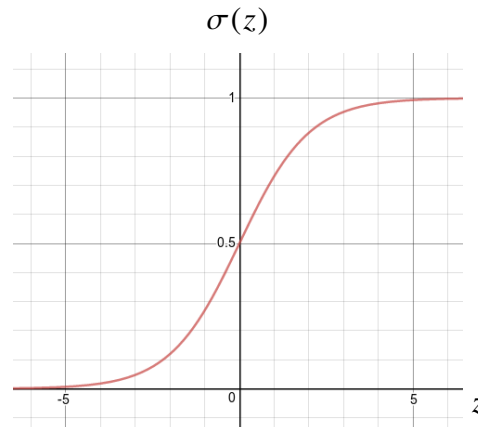
$$\begin{aligned} \theta^T \mathbf{x} &= \sum_{i=1}^m \theta_i x_i = \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_m x_m \\ \sigma(z) &= \frac{1}{1 + e^{-z}} \end{aligned}$$

The superscript  $T$  in  $\theta^T \mathbf{x}$  represents a matrix transpose; the operation  $\theta^T \mathbf{x}$  is equivalent to taking the *dot product* of the vectors  $\theta$  and  $\mathbf{x}$ , or simply a weighted sum of the components of  $\mathbf{x}$  (with  $\theta$  containing the weights).

---

<sup>1</sup>Yes, this is a terribly confusing name, given that **regression** refers to tasks that require predicting continuous values. Perhaps *logistic classification* would have been better.

The function  $\sigma(z) = \frac{1}{1+e^{-z}}$  is called the **logistic function** (or *sigmoid function*); it looks like this:



An important quality of this function is that it maps all real numbers to the range (0, 1). In logistic regression,  $\sigma(z)$  turns an arbitrary “score”  $z$  into a number between 0 and 1 that is interpreted as a probability. Positive numbers become high probabilities; negative numbers become low ones.

## Log Likelihood

In order to choose values for the parameters of logistic regression, we use maximum likelihood estimation (MLE). As such we are going to have two steps: (1) write the log-likelihood function and (2) find the values of  $\theta$  that maximize the log-likelihood function.

The labels that we are predicting are binary, and the output of our logistic regression function is supposed to be the probability that the label is one. This means that we can (and should) interpret each label as a Bernoulli random variable:  $Y \sim \text{Ber}(p)$  where  $p = \sigma(\theta^T \mathbf{x})$ .

To start, here is a super slick way of writing the probability of one data point (recall this is the equation form of the probability mass function of a Bernoulli):

$$P(Y = y | X = \mathbf{x}) = \sigma(\theta^T \mathbf{x})^y \cdot [1 - \sigma(\theta^T \mathbf{x})]^{(1-y)}$$

Now that we know the probability mass function, we can write the likelihood of all the data:

$$L(\theta) = \prod_{i=1}^n P(Y = y^{(i)} | X = \mathbf{x}^{(i)}) \quad \text{the likelihood of independent training labels}$$

$$= \prod_{i=1}^n \sigma(\theta^T \mathbf{x}^{(i)})^{y^{(i)}} \cdot [1 - \sigma(\theta^T \mathbf{x}^{(i)})]^{(1-y^{(i)})} \quad \text{substituting the likelihood of a Bernoulli}$$

And if you take the log of this function, you get the log likelihood for logistic regression. The log likelihood equation is:

$$LL(\theta) = \sum_{i=1}^n y^{(i)} \log \sigma(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \log[1 - \sigma(\theta^T \mathbf{x}^{(i)})]$$

Recall that in MLE the only remaining step is to choose parameters ( $\theta$ ) that maximize log likelihood.

## Gradient of Log Likelihood

Now that we have a function for log-likelihood, we simply need to choose the values of theta that maximize it. Unfortunately, if we try just setting the derivative equal to zero, we'll quickly get frustrated: there's no closed form for the maximum. However, we can find the best values of theta by using an optimization algorithm. The optimization algorithm we will use requires the partial derivative of log likelihood with respect to each parameter. First I am going to give you the partial derivative (so you can see how it is used); we'll derive it a bit later:

$$\frac{\partial LL(\theta)}{\partial \theta_j} = \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)}$$

## Gradient Ascent Optimization

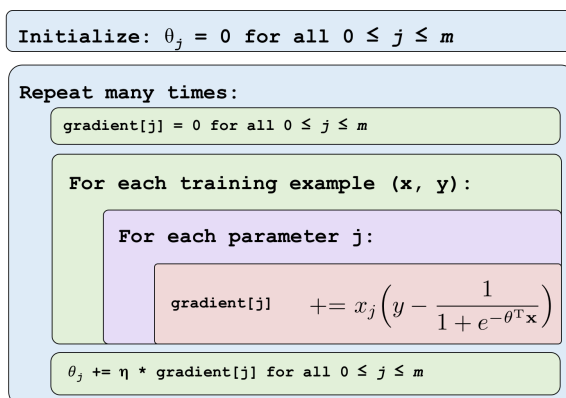
Our goal is to choose parameters ( $\theta$ ) that maximize likelihood, and we know the partial derivative of log likelihood with respect to each parameter. We are ready for our optimization algorithm.

In the case of logistic regression we can't solve for  $\theta$  mathematically. Instead we use a computer to choose  $\theta$ . To do so we employ an algorithm called gradient ascent (a classic in optimization theory). The idea behind gradient ascent is that gradients point "uphill". If you continuously take small steps in the direction of your gradient, you will eventually make it to a local maximum. In the case of logistic regression you can prove that the result will always be a global maximum.

The update to our parameters that results in each small step can be calculated as:

$$\begin{aligned} \theta_j^{\text{new}} &= \theta_j^{\text{old}} + \eta \cdot \frac{\partial LL(\theta^{\text{old}})}{\partial \theta_j^{\text{old}}} \\ &= \theta_j^{\text{old}} + \eta \cdot \sum_{i=1}^n [y^{(i)} - \sigma(\theta^T \mathbf{x}^{(i)})] x_j^{(i)} \end{aligned}$$

Where  $\eta$  is the magnitude of the step size that we take. If you keep updating  $\theta$  using the equation above, you will converge on the best values of  $\theta$ . You now have an intelligent model. Here is the gradient ascent algorithm for logistic regression in pseudo-code:



It is also common to have a parameter  $\theta_0$  that is added as a constant to the  $\theta^T \mathbf{x}$  inside the sigmoid. Rather than computing special derivatives for  $\theta_0$ , we can simply define an additional feature  $x_0$  that always takes the value 1. Taking a weighted average then results in adding  $\theta_0$ , the weight for  $x_0$ .

## Derivations

In this section we provide the mathematical derivations for the gradient of log-likelihood. The derivations are worth knowing because these ideas are heavily used in Artificial Neural Networks.

Our goal is to calculate the derivative of the log likelihood with respect to each theta. To start, here is the definition for the derivative of a sigmoid function with respect to its inputs:

$$\frac{\partial}{\partial z} \sigma(z) = \sigma(z)[1 - \sigma(z)] \quad \text{to get the derivative with respect to } \theta, \text{ use the chain rule}$$

Take a moment and appreciate the beauty of the derivative of the sigmoid function. The reason that sigmoid has such a simple derivative stems from the natural exponent in the sigmoid denominator.

Since the likelihood function is a sum over all of the data, and in calculus the derivative of a sum is the sum of derivatives, we can focus on computing the derivative of one example. The gradient of theta is simply the sum of this term for each training data point.

First I am going to show you how to compute the derivative the hard way. Then we are going to look at an easier method. The derivative of gradient for one data point  $(\mathbf{x}, y)$ :

$$\begin{aligned} \frac{\partial LL(\theta)}{\partial \theta_j} &= \frac{\partial}{\partial \theta_j} y \log \sigma(\theta^T \mathbf{x}) + \frac{\partial}{\partial \theta_j} (1 - y) \log[1 - \sigma(\theta^T \mathbf{x})] && \text{derivative of sum of terms} \\ &= \left[ \frac{y}{\sigma(\theta^T \mathbf{x})} - \frac{1 - y}{1 - \sigma(\theta^T \mathbf{x})} \right] \frac{\partial}{\partial \theta_j} \sigma(\theta^T \mathbf{x}) && \text{derivative of log } f(x) \\ &= \left[ \frac{y}{\sigma(\theta^T \mathbf{x})} - \frac{1 - y}{1 - \sigma(\theta^T \mathbf{x})} \right] \sigma(\theta^T \mathbf{x}) [1 - \sigma(\theta^T \mathbf{x})] x_j && \text{chain rule + derivative of } \sigma \\ &= \left[ \frac{y - \sigma(\theta^T \mathbf{x})}{\sigma(\theta^T \mathbf{x}) [1 - \sigma(\theta^T \mathbf{x})]} \right] \sigma(\theta^T \mathbf{x}) [1 - \sigma(\theta^T \mathbf{x})] x_j && \text{algebraic manipulation} \\ &= [y - \sigma(\theta^T \mathbf{x})] x_j && \text{cancelling terms} \end{aligned}$$

## Derivatives Without Tears

That was the hard way. Logistic regression is the building block of artificial neural networks. If we want to scale up, we are going to have to get used to an easier way of calculating derivatives. For that we are going to have to welcome back our old friend the chain rule (of derivatives). By the chain rule:

$$\begin{aligned} \frac{\partial LL(\theta)}{\partial \theta_j} &= \frac{\partial LL(\theta)}{\partial p} \cdot \frac{\partial p}{\partial \theta_j} && \text{where } p = \sigma(\theta^T \mathbf{x}) \\ &= \frac{\partial LL(\theta)}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial \theta_j} && \text{where } z = \theta^T \mathbf{x} \end{aligned}$$

Chain rule is the decomposition mechanism of calculus. It allows us to calculate a complicated partial derivative ( $\frac{\partial LL(\theta)}{\partial \theta_j}$ ) by breaking it down into smaller pieces.

$LL(\theta) = y \log p + (1 - y) \log(1 - p)$	where $p = \sigma(\theta^T \mathbf{x})$
$\frac{\partial LL(\theta)}{\partial p} = \frac{y}{p} - \frac{1 - y}{1 - p}$	taking the derivative
<hr/>	
$p = \sigma(z)$	where $z = \theta^T \mathbf{x}$
$\frac{\partial p}{\partial z} = \sigma(z)[1 - \sigma(z)]$	taking the derivative of the sigmoid
<hr/>	
$z = \theta^T \mathbf{x}$	as previously defined
$\frac{\partial z}{\partial \theta_j} = x_j$	only $x_j$ interacts with $\theta_j$

Each of those derivatives was much easier to calculate. Now we simply multiply them together.

$\frac{\partial LL(\theta)}{\partial \theta_j} = \frac{\partial LL(\theta)}{\partial p} \cdot \frac{\partial p}{\partial z} \cdot \frac{\partial z}{\partial \theta_j}$	
$= \left[ \frac{y}{p} - \frac{1 - y}{1 - p} \right] \cdot \sigma(z)[1 - \sigma(z)] \cdot x_j$	substituting in for each term
$= \left[ \frac{y}{p} - \frac{1 - y}{1 - p} \right] \cdot p[1 - p] \cdot x_j$	since $p = \sigma(z)$
$= [y(1 - p) - p(1 - y)] \cdot x_j$	multiplying in
$= [y - p]x_j$	expanding
$= [y - \sigma(\theta^T \mathbf{x})]x_j$	since $p = \sigma(\theta^T \mathbf{x})$