

## Logistic Regression

There is a *lot* more that could be said about linear regression. But I'm going to leave most of that for statistics courses. There is also a *lot* more that could be said about gradient-based optimization, and I'll return to some of it later. In this note we'll add a non-linearity to our model, and introduce one of the most used machine learning models.

### Transforming the output

We saw that we could attempt to fit the probability of being in a particular class with straightforward linear regression models. However, we are likely to see outputs outside the range  $[0, 1]$  for some test inputs. We will now force our function to lie in the desired  $[0, 1]$  range by transforming a linear function with a logistic sigmoid:

$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}.$$

As with linear regression, we can replace our input features  $\mathbf{x}$  with a vector of basis function values  $\phi(\mathbf{x})$ . I won't clutter the notation with this detail. Wherever you see  $\mathbf{x}$ , know that you can replace this input vector with a version that has been transformed in any way you like.

### Loss function

As before, we wish to fit the function to match the training data closely. If the labels are zero and one,  $y \in \{0, 1\}$ , we *could* minimize the square loss

$$\sum_{n=1}^N (y^{(n)} - f(\mathbf{x}^{(n)}; \mathbf{w}))^2.$$

However, the *Logistic Regression* model uses the interpretation of the function as a probability,  $f(\mathbf{x}; \mathbf{w}) = P(y=1 | \mathbf{x}, \mathbf{w})$ , more directly. Maximum likelihood fitting of this model maximizes the probability of the data:

$$L(\mathbf{w}) = \prod_{n=1}^N P(y^{(n)} | \mathbf{x}^{(n)}, \mathbf{w}),$$

for the model with parameters  $\mathbf{w}$ . Equivalently, we minimize the negative log-probability of the training labels, which for this model can be written as:

$$\text{NLL} = -\log L(\mathbf{w}) = -\sum_{n=1}^N \log \left[ \sigma(\mathbf{w}^\top \mathbf{x}^{(n)})^{y^{(n)}} (1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(n)}))^{1-y^{(n)}} \right],$$

or

$$\text{NLL} = -\sum_{n:y^{(n)}=1} \log \sigma(\mathbf{w}^\top \mathbf{x}^{(n)}) - \sum_{n:y^{(n)}=0} \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(n)})).$$

There is a trick to write the cost function more compactly. We transform the labels to be  $z^{(n)} \in \{-1, +1\}$  where  $z^{(n)} = (2y^{(n)} - 1)$ , and noticing  $\sigma(-a) = 1 - \sigma(a)$ , we can write:

$$\text{NLL} = -\sum_{n=1}^N \log \sigma(z^{(n)} \mathbf{w}^\top \mathbf{x}^{(n)}).$$

As before, either cost function can have a regularizer added to it to discourage extreme weights.

Maximum likelihood estimation has some good statistical properties. In particular, asymptotically (for lots of data) it is the most efficient estimator. Although the loss can be extreme where confident wrong predictions are made, which could mean that outliers cause more problems than with the square loss approach.

## Gradients

The final required ingredient is the gradient vector  $\nabla_{\mathbf{w}}\text{NLL}$ . A gradient-based optimizer can then find the weights that minimize our cost.

The derivative of the logistic sigmoid is as follows<sup>1</sup>:

$$\frac{\partial\sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a)),$$

which tends to zero at the asymptotes  $a \rightarrow \pm\infty$ .

I like to derive the derivatives using the third form of the cost function NLL, because it's shorter, although I think most books use the first form. We'll get an equivalent answer. For brevity, I'll use  $\sigma_n = \sigma(\mathbf{z}^{(n)} \mathbf{w}^\top \mathbf{x}^{(n)})$ . We then apply the chain rule:

$$\begin{aligned}\nabla_{\mathbf{w}}\text{NLL} &= - \sum_{n=1}^N \nabla_{\mathbf{w}} \log \sigma_n = - \sum_{n=1}^N \frac{1}{\sigma_n} \nabla_{\mathbf{w}} \sigma_n = - \sum_{n=1}^N \frac{1}{\sigma_n} \sigma_n (1 - \sigma_n) \nabla_{\mathbf{w}} \mathbf{z}^{(n)} \mathbf{w}^\top \mathbf{x}^{(n)}, \\ &= - \sum_{n=1}^N (1 - \sigma_n) \mathbf{z}^{(n)} \mathbf{x}^{(n)}.\end{aligned}$$

Interpretation:  $\sigma_n$  is the probability assigned to the correct training label. So when the classifier is confident and correct on an example, it contributes little to the gradient. Stochastic gradient descent will improve the examples the classifier gets wrong, or is less confident about, by pushing the weights parallel to the direction of the corresponding input.

## Some things to know about gradients

Whenever we can compute a differentiable cost function, it should always be possible to compute *all* of the derivatives at once in a similar number of operations to one function evaluation. That's an amazing result from the old field of *automatic differentiation*. (Caveat it might take a lot of memory.) So if your derivatives are orders of magnitude more expensive than your cost function, you are probably doing something wrong.

Despite a long history, few people use fully automatic differentiation in machine learning. There are machine learning tools like Theano, and Tensor Flow that will do most of the work for you. But some people<sup>2</sup> are still needing to do some work by hand so that those tools can work on all the models that we might build.

Whether you are differentiating by hand, or writing a compiler to compute derivatives, you need to test your code. Derivatives are easily checked by *finite differences*:

$$\frac{\partial f(w)}{\partial w} \approx \frac{f(w+\epsilon/2) - f(w-\epsilon/2)}{\epsilon},$$

and so should always be checked. Unless the weights are extreme, I'd normally set  $\epsilon = 10^{-5}$ . You have to perturb each parameter in turn to evaluate one element of a gradient vector  $\nabla_{\mathbf{w}} f(\mathbf{w})$ . Therefore, for  $D$ -dimensional vectors of derivatives, the computational cost of finite differences scales  $D$  times worse than well-written derivative code, as well as being less accurate. Finite differences are a useful check, but not for use in production.

## Check your understanding

- We need to use iterative optimizers like stochastic gradient descent to fit logistic regression. If we swapped from negative-log-likelihood to the square loss, explain whether we would be able to fit the model using a single \ or 1stsq fit.

1. It's not hard to show, but I'd give you this result in an exam if you needed it.  
2. Your lecturer is one of them: <https://arxiv.org/abs/1602.07527>.

- The expression that I've given for the logistic regression gradient  $\nabla_{\mathbf{w}} \text{NLL}$  looks different from the ones given in many textbooks. Describe how to quickly numerically check that the different-looking expressions are equivalent, without having to do any detailed mathematics. (You could also do a check!)
- We have some data where we know the target outputs  $y$  should always be positive. If we fit standard linear regression, we could predict negative values. How might we fix this deficiency? Can we still do a standard \ or lstsq fit, or do we need to swap to an iterative optimizer?

### Further Reading

All machine learning textbooks should have a treatment of logistic regression. You could read an alternative treatment to this note in Barber 17.4.1 and 17.4.4. Or Murphy: quick introduction section 1.4.6, then Chapter 8, which has an in depth treatment beyond this note.

Tom Minka has a review of alternative batch optimizers for logistic regression  
<http://research.microsoft.com/en-us/um/people/minka/papers/logreg/>  
 (Stochastic gradient methods were less popular then, and were not considered.)

For a large-scale practical tool that uses stochastic optimization, check out Vowpal Wabbit:  
[https://github.com/JohnLangford/vowpal\\_wabbit/wiki](https://github.com/JohnLangford/vowpal_wabbit/wiki)  
 Its framework includes support for logistic regression. It has various tricks to train *fast*. It can cope with data, like large-scale text corpora, where you might not know what you want your features to be until you start streaming data.

We don't have to transform a linear function with the logistic sigmoid. We could instead create a 'probit' model, which uses a Gaussian's cumulative density function (cdf) instead of the logistic sigmoid. We can also transform the function to model other data types. For example, count data can be modeled by using the underlying linear function to set the log-rate of a Poisson distribution. These alternatives can be unified as "Generalized Linear Models" (GLMs). R has a widely used glm library.

There is a cool trick with complex numbers to evaluate derivatives to machine precision, which I'd like to share:  
<http://blogs.mathworks.com/cleve/2013/10/14/complex-step-differentiation/>  
 It is no faster than finite differences though, so shouldn't be used, except as a check.