

# Logistic Regression

## Recap

Last time we considered how to make a hypothesis that approximates a real value function. All we need to do is build a cost function and then minimize it. This time we are going to use the same idea for binary function.

## Logistic Regression

Suppose you have data points

$$(\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), (\mathbf{x}^{(3)}, y^{(3)}), \dots (\mathbf{x}^{(N)}, y^{(N)})$$

where  $\mathbf{x}^{(i)}$  are feature vectors and  $y^{(i)} \in \{+1, -1\}$  are the classes.

Now the question is the input feature  $\mathbf{x}$  what is the probability that this datapoint would be of class +1.

So, the unknown target function in this case is a probability. Specifically,

$$f(\mathbf{x}) = P[y = +1 | \mathbf{x}]$$

The probability on the right reads: the probability that the given data point is of class +1 given that the data point has input feature vector  $\mathbf{x}$ . For example, the probability that a person would have more than 50k income given that he is age 35, working in computer industry, married and have 2 kids.

Note that this also imply

$$P[y = -1 | \mathbf{x}] = 1 - f(\mathbf{x})$$

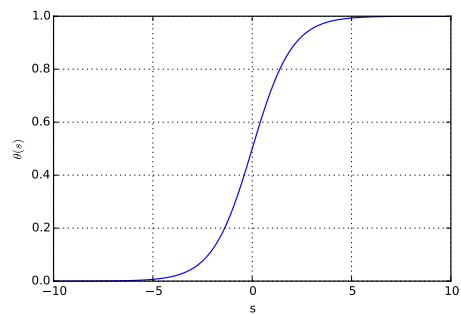
The probability is a real number between 0 and 1. We can't use linear regression for this problem since the probability ranges from +1 to -1. But the output of the linear regression is from  $-\infty$  to  $\infty$ . First attempt to fix this could be trying to take the sign of the linear combination and shift it up then scale it so it ranges from -1 to 1. The problem with this approach is that

there would be no input vector taking the value in between 0 and 1; ex: 0.8.

To solve this issue we need to introduce a soft threshold function. There many soft threshold function but the most popular one is called logistic function

$$\theta(s) = \frac{e^s}{1 + e^s} = \frac{1}{e^{-s} + 1}$$

All this function does is that it convert  $s$  from  $-\infty$  to  $\infty$  to 0 and 1 continuously. So we have a chance of making it probability. The second one is more friendly to computer round off.



Before we go on there is once property of  $\theta$  that makes this particularly useful.

$$\theta(-s) = 1 - \theta(s)$$

You can verify this easily.

Let us use this to make hypothesis for the probability of a data point  $x$  being of class +1:

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

Here I also use  $x_0 = 1$  convention to get the threshold term. To interpret this quantity as conditional probability we can do

$$\begin{aligned} P(y | \mathbf{x}) &= \begin{cases} \theta(\mathbf{w}^T \mathbf{x}) & y = +1 \\ 1 - \theta(\mathbf{w}^T \mathbf{x}) & y = -1 \end{cases} \\ &= \begin{cases} \theta(\mathbf{w}^T \mathbf{x}) & y = +1 \\ \theta(-\mathbf{w}^T \mathbf{x}) & y = -1 \end{cases} \\ &= \theta(y \mathbf{w}^T \mathbf{x}) \end{aligned}$$

54 The last line is a very neat thing to have

$$P(y|\mathbf{x}) = \theta(y\mathbf{w}^T\mathbf{x}) \quad (1)$$

55 So, now the problem is to find  $\mathbf{w}$  so that our  
 56 hypothesis best represents the conditional prob-  
 57 ability of the data. To do this, we need to in-  
 58 troduce the concept of likelihood. The likelihood  
 59 is just the probability that you will be given the  
 60 data set you were given. Specifically,

$$\mathcal{L} = \prod_{i=1}^N P(y^{(i)}|\mathbf{x}^{(i)})$$

61 If we replace the probability with the hypoth-  
 62 esis we have

$$\mathcal{L} = \prod_{i=1}^N \theta(y\mathbf{w}^T\mathbf{x})$$

63 Thus, what we need to do now is the find  $\mathbf{w}$  that  
 64 maximize the likelihood. Let me point out one  
 65 important caveat. What we are trying to do here  
 66 is find  $\mathbf{w}$  such that the data points we drawn has  
 67 the highest likelihood. This is NOT the same  
 68 thing as what is finding  $\mathbf{w}$  that the data points  
 69 we have is most likely drawn from. The second  
 70 quantity is unknown and can't really be quanti-  
 71 fied eventhough it is really what we want. But  
 72 we take the first quantity as being "good enough"  
 73 for what we are trying to do.

74 Maximizing the product is quite cumbersome.  
 75 We can change the product to sum by taking the  
 76 natural log. Moreover, since most package are  
 77 written for minimization rather than maixmiza-  
 78 tion; all we need to turn the maximization to  
 79 minimization problem is just multiplying it by a  
 80 minus sign. Thus maximizing the likelihood is  
 81 equivalent to minimizing the following quantity:

$$\text{cost}(\mathbf{w}) = -\frac{1}{N} \log \mathcal{L} \quad (2)$$

$$= -\frac{1}{N} \sum_{i=1}^N \ln [\theta(y\mathbf{w}^T\mathbf{x})] \quad (3)$$

$$= \frac{1}{N} \sum_{i=1}^N \ln \left[ \frac{1}{\theta(y\mathbf{w}^T\mathbf{x})} \right] \quad (4)$$

$$= \frac{1}{N} \sum_{i=1}^N \ln [1 + e^{-y\mathbf{w}^T\mathbf{x}}] \quad (5)$$

82 Unlike linear regression, this cost has no  
 83 closed form solution for  $\nabla \text{cost} = 0$ . So we need  
 84 to resort to numerical methods like gradient de-  
 85 scent we covered last time.

86 Once you get the  $\mathbf{w}$  that minimize the log  
 87 likelihood from gradient descent method or oth-  
 88 ers. We will have a hypothesis that approximate  
 89 the probability of it being class +1.

## 90 Hard Threshold

91 What we did in the last section was all nice that  
 92 we get an approximation for the target probabil-  
 93 ity but this doesn't tell us whether the new input  
 94 vector is of the class +1 or -1. This is easy all  
 95 you need to do is making a cut off.

$$\text{class} = \begin{cases} +1 & h(\mathbf{x}) > c \\ -1 & h(\mathbf{x}) \leq c \end{cases} \quad (6)$$

96 But, how do we pick the cut off? The answer  
 97 depends on practical application. It depends on  
 98 how much you are willing to tolerate false posi-  
 99 tive vs false negative.

## 100 ROC Curve

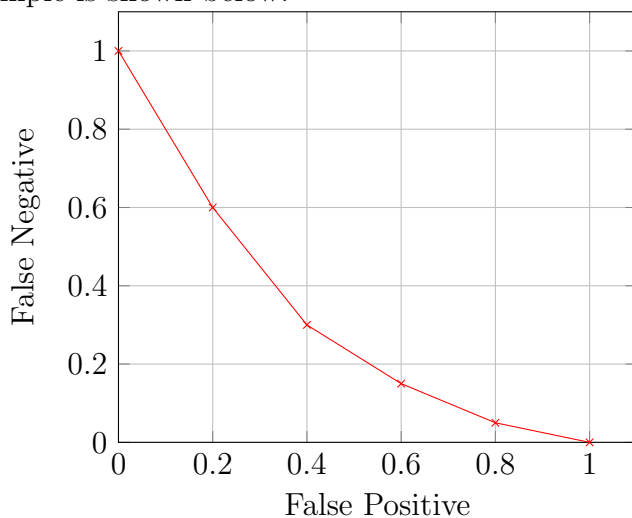
101 When you try to make a hard threshold people  
 102 usually plot false positive rate vs false negative  
 103 rate for various value of cutoff. This will let you  
 104 see the performance of your classifier. The fancy  
 105 name for this curve is called *Receiver Operating*  
 106 *Characteristic Curve* (ROC Curve). There is a  
 107 whole family of ROC curves; each with different

meaning on the axis and different names too, plus each people call it differently too.. Don't bother memorizing the name. Use common sense. The important thing is that when you see one and know what the axis means you know where a good classifier should be and where a bad classifier should be.

For example, if we care about

1. Out of 1000 true positive class how many sample get classified as negative (false negative). We will call this ratio false negative rate.
2. Out for 200 true negative class how many sample get classified as positive (false positive). We will call this ratio false positive rate.

We could plot, for each cutoff we scan over, the false positive rate versus the false negative rate. This will show us how each cut off performs. An example is shown below:



From the figure, let us start out with the lower right corner where it corresponds to cutoff value of zero. At this corner we will never get any false negative since we never classify anything as negative class. The false positive is high since we classify everything as positive. As we increase the cutoff we will classify less thing as positive and more thing as negative. This will increase false negative and decrease false positive. This will make the curve go toward the upper left. With the cut off really high, you are classifying everything as negative. You will therefore never get false positive since you never classify anything as a positive class to begin with. But you will get 100% false negative.

So if you were to ask where we want to be, it is clear that it is lower left corner. This is the place where you have no false positive and no false negative.

Please do not memorize which one increase, which one decrease and where is the perfect place though. You can just use common sense to figure this out from the first principle. There are many variations of ROC curve with different normalization each one will go toward different corner.

For example you can plot the ratio of true positive class classified as positive class versus the ratio of true negative class classified as negative class. This will give you different kind of curve which you can also use to visualize how cutoff affect the performance of the classifier. As an exercise think about which corner a good classifier would be if we plot these two ratio.

The important thing is that when you see one and you know what each axis represents you can figure out where the good classifier would be.