

## ML: Neural Networks: Learning:

## Cost Function:

We first define a few variables that we will need to know:

a)  $L$  = total number of layers in the network

b)  $s_1 = \text{number of units (not counting bias units) in layer 1}$

c)  $k = \text{number of output units / classes.}$

Recall that in neural networks, we may have many output nodes. We denote  $h_\theta(x)_k$  as being a hypothesis that results in the  $k^{\text{th}}$  output.

Our cost function for our neural network is going to be a generalization of the one we used for logistic regression.

Recall that the cost function for regularized logistic regression was:

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

For neural networks, it is going to be slightly more complicated:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ y_k^{(i)} \log((h_\theta(x^{(i)}))_k) + (1 - y_k^{(i)}) \log(1 - (h_\theta(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^m \sum_{j=1}^{s_l+1} (\theta_{j,i}^{(l)})^2$$

We have added a few nested summations to account for our multiple output nodes. In the first part of the cost function, after the square brackets, we have an additional nested summation that loops through the number of output nodes.

In the regularization part, after the square brackets, we must account for multiple theta matrices. The number of columns in our current

The matrix is equal to the number of nodes in our current layer (including the bias unit). The number of rows in our current layer matrix is equal to the number of nodes in the next layer (excluding the bias unit). As before with logistic regression, we ~~will~~ ~~not~~ ~~use~~ softmax during forward.

Note:

- The Double sum simply adds up the logistic regression costs calculated for each cell in the output layer; and
- The Triple sum simply adds up the squares of all the individual  $\Theta_{ij}$  in the weight matrix.
- The  $i$  in the Triple sum does not refer to training example;

## Backpropagation Intuition:

The cost function is:

$$J(\theta) = -\frac{1}{m} \sum_{j=1}^m \sum_{k=1}^{s_j} \left[ y_k^{(j)} \log(h_\theta(x^{(j)})_k) + (1 - y_k^{(j)}) \log(1 - h_\theta(x^{(j)})_k) \right] \\ + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_i} \sum_{j=1}^{s_{i+1}} (\theta_{j,i}^{(l)})^2$$

If we consider simple non-muticlass classification ( $k=1$ ) and disregard regularization, the cost is computed with:

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

More intuitively, we can think of that equation roughly as:

$$\text{Cost}(j) = (h_\theta(x^{(j)}) - y^{(j)})^2$$

Intuitively,  $\delta_j^{(l)}$  is the "error" for  $\theta_{j,i}^{(l)}$  (unit  $j$  in layer 1)

More formally, the delta values are obtained by the derivative of the cost function:

$$\delta_j^{(l)} = \frac{\partial \text{Cost}(j)}{\partial z_j^{(l)}}$$

Recall that our derivative is the slope of a line tangent to the cost function, so that the steeper the slope the more incorrect we are.

Note: I believe, sometimes  $i$  is used to index a training example. Sometimes  $i$  is used to index a unit in a layer. In the back propagation algorithm described here,  $i$  is used to index a training example rather than indexing the row of  $x$ .

## (2)

### Backpropagation Algorithm:

"Backpropagation" is neural-network terminology for minimizing our cost function, just like what we were doing with gradient descent in logistic and linear regression. Our goal is to compute:  
 $\min_{\theta} J(\theta)$

that is, we want to minimize our cost function  $J$  with respect to all of the parameters in  $\theta$ . In this section we'll look at the questions we use to compute the partial derivative of  $J(\theta)$ :

$$\frac{\partial}{\partial \theta^{(i)}_{j,j}} J(\theta)$$

In other words, we use the following algorithm:

### Intuition

### Backpropagation Algorithm

Given training set  $\{(x^{(1)}, y^{(1)}) \dots (x^{(m)}, y^{(m)})\}$

- Set  $\Delta_{i,j}^{(l)} := 0$  for all  $(l, i, j)$ , (hence you end up having a matrix full of zeros)

For training example  $j = 1 \dots m$ :

- Set  $a^{(1)} := x^{(j)}$

- Perform forward propagation to compute  $a^{(l)}$  for  $l = 2, 3, \dots, L$

Gradient computation:

Given one training example  $(x, y)$ :

Forward propagation:

$$a^{(1)} = x$$

$$z^{(2)} = \theta^{(1)} a^{(1)}$$

$$a^{(2)} = g(z^{(2)}) \quad (\text{odd } a_0^{(2)})$$

$$z^{(3)} = \theta^{(2)} a^{(2)}$$

$$a^{(3)} = g(z^{(3)}) \quad (\text{odd } a_0^{(3)})$$

$$z^{(4)} = \theta^{(3)} a^{(3)}$$

$$a^{(4)} = h(x) - a^{(4)}$$

3. Using  $y^{(t)}$ , compute  $\delta^{(L)} = \alpha^{(L)} - y^{(t)}$

Where  $L$  is our total number of layers and  $\alpha^{(L)}$  is the vector of outputs of the activation units for the last layer. So our "error values" for the last layer are simply the difference of our actual results in the last layer and the correct output in  $y$ . To get the delta values of the layers before the last layer, we can use an equation that steps us back from right to left:

4. Compute  $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(1)}$  using  $\delta^{(1)} = ((\Theta^{(1)})^T \delta^{(2)}) + \alpha^{(1)} * (I - \alpha^{(1)})$

The delta values of layer  $l$  are calculated by multiplying the delta values in the next layer with the  $\Theta$  matrix of layer  $l$ . We then don't have to multiply that with a function called  $\phi'$ , or  $\phi$ -prime, which is the derivative of the activation function  $\phi$  weighted with the input values given by  $z^{(1)}$ .

The  $\phi$ -prime derivative terms can also be written out as:

$$\phi'(z^{(1)}) = \alpha^{(1)} * (I - \alpha^{(1)})$$

$$5. \Delta_{i,j}^{(2)} := \Delta_{i,j}^{(1)} + \alpha_j^{(2)} \delta^{(2+1)} + \delta^{(2+1)} (\alpha^{(2)})^T$$

Or with notation,  $\Delta^{(1)} := \Delta^{(1)}$

Hence we update our new  $\Delta$  matrix:

$$\bullet D_{i,j}^{(1)} := \frac{1}{m} (\Delta_{i,j}^{(1)} + \lambda \Theta_{i,j}^{(1)}), \text{ if } j \neq 0.$$

$$\bullet D_{i,0}^{(1)} := \frac{1}{m} \Delta_{i,0}^{(1)} \text{ if } j = 0$$

The capital-delta matrix  $D$  is used as an "accumulator" to add up the values of the  $\Delta$ 's along and eventually compute our partial derivatives. Thus we get  $\frac{\partial}{\partial \Theta_{i,j}^{(1)}} J(\Theta) = D_{i,j}^{(1)}$