# Tutorial for Ex.4 Forward and Backpropagation (Spring 2014 session)

This tutorial outlines the process of accomplishing the goals for Programming Exercise 4. The purpose is to create a collection of all the useful yet scattered and obscure knowledge that otherwise would require hours of frustrating searches.This tutorial is targeted solely at vectorized implementations. If you're a looper, you're doing it the hard way, and you're on your own. I'll use the less-than-helpful greek letters and math notation from the video lectures in this tutorial, though I'll start off with a glossary so we can agree on what they are. I will also suggest some common variable names, so students can more easily get help on the Forum. It is left to the reader to convert these lines into program statements. <u>You will need to determine the correct order and transpositions for each matrix multiplication</u>.

Most of this material appears in either the video lectures, slides, course wiki, or the ex4.pdf file, though nowhere else does it all appear in one place.

Glossary: Each of these variables will have a subscript, noting which NN layer it is associated with.

$\Theta$ : A matrix of weights to compute the inner values of the neural network. When we used single-vector theta values, it was noted with the lower-case character $\theta$.

z : is the result of multiplying a data vector with a $\Theta$ matrix. A typical variable name would be "z2".

a : The "activation" output from a neural layer. This is always generated using a sigmoid function g() on a z value. A typical variable name would be "a2".

$\delta$ : lower-case delta is used for the "error" term in each layer. A typical variable name would be "d2".

$\Delta$ : upper-case delta is used to hold the sum of the product of a $\delta$ value with the previous layer's a value. In the vectorized solution, these sums are calculated automatically though the magic of matrix algebra. A typical variable name would be "Delta2".

$\Theta$ gradient : This is the thing we're looking for, the partial derivative of theta. There is one of these variables associated with each $\Delta$. These values are returned by nnCostFunction(), so the variable names must be "Theta1_grad" and "Theta2_grad".

g() is the sigmoid function. g'() is the sigmoid gradient function.

Tip: One handy method for ignoring a column of bias units is to use the notation "SomeMatrix(:,2:end)". This selects all of the rows of a matrix, and omits the entire first column.

Here we go. Nearly all of the editing in this exercise happens in nnCostFunction.m. Let's get started.

A note regarding the sizes of these data objects: See the Appendix at the bottom of the tutorial for information on the sizes of the data objects.

A note regarding bias units, regularization, and back-propagation: There are two methods for handling the bias units in the back-propagation and gradient calculations. I've described only one of them here, it's the one that I understood the best. Both methods work so choose the one that makes sense to you and avoids dimension errors. It matters not a whit whether the bias unit is dropped before or after it is calculated - both methods give the same results, though the order of operations and transpositions required may be different. Those with contrary opinions are welcome to write their own tutorial.

Forward Propagation:

We'll start by outlining the forward propagation process. Though this was already accomplished once during Exercise 3, you'll need to duplicate some of that work because computing the gradients requires some of the intermediate results from forward propagation.

Step 1 - Expand the 'y' output values into a matrix of single values (see ex4.pdf Page 5). This is most easily done using an eye() matrix of size num_labels, with vectorized indexing by 'y', as in "eye(num_labels)(y,:)". Discussions of this and other methods are available in the Course Wiki - Programming Exercises section. A typical variable name would be "y_matrix".

Step 2 - perform the forward propagation: a1 equals the X input matrix with a column of 1's added (bias units). z2 equals the product of a1 and Θ1. a2 is the result of passing z2 through g(). a2 then has a column of 1st added (bias units). z3 equals the product of a2 and Θ2. a3 is the result of passing z3 through g()

Cost Function, non-regularized

Step 3 - Compute the unregularized cost according to ex4.pdf (top of Page 5), (I had a hard time understanding this equation mainly that I had a misconception that y(i)$k$ is a vector, instead it is just simply one number) using a3, your $y$matrix, and m (the number of training examples). Cost should be a scalar value. If you get a vector of cost values, you can sum that vector to get the cost. Remember to use element-wise multiplication with the log() function. Now you can run ex4.m to check the unregularized cost is correct, then you can submit Part 1 to the grader.

Cost Regularization

Step 4 - Compute the regularized component of the cost according to ex4.pdf Page 6, using $\Theta1$ and $\Theta2$ (ignoring the columns of bias units), along with $\lambda$, and m. The easiest method to do this is to compute the regularization terms separately, then add them to the unregularized cost from Step 3. You can run ex4.m to check the regularized cost, then you can submit Part 2 to the grader.

Sigmoid Gradient and Random Initialization

Step 5 - You'll need to prepare the sigmoid gradient function g′(), as shown in ex4.pdf Page 7. You can submit Part 3 to the grader.

Step 6 - Implement the random initialization function as instructed on ex4.pdf, top of Page 8. You do not submit this function to the grader.

Backpropagation

Step 7 - Now we work from the output layer back to the hidden layer, calculating how bad the errors are. See ex4.pdf Page 9 for reference. δ3 equals the difference between a3 and the y_matrix. δ2 equals the product of δ3 and $\Theta2$ (ignoring the $\Theta2$ bias units), then multiplied element-wise by the g′() of z2 (computed back in Step 2). Note that at this point, the instructions in ex4.pdf are specific to looping implementations, so the notation there is different. Δ2 equals the product of δ3 and a2. This step calculates the product and sum of the errors. Δ1 equals the product of δ2 and a1. This step calculates the product and sum of the errors.

Gradient, non-regularized

Step 8 - Now we calculate the non-regularized theta gradients, using the sums of the errors we just computed. (see ex4.pdf bottom of Page 11). $\Theta1$ gradient equals Δ1 scaled by 1/m. $\Theta2$ gradient

equals Δ2 scaled by 1/m. The ex4.m script will also perform gradient checking for you, using a smaller test case than the full character classification example. So if you're debugging your nnCostFunction() using the "keyboard" command during this, you'll suddenly be seeing some much smaller sizes of X and the Θ values. Do not be alarmed. If the feedback provided to you by ex4.m for gradient checking seems OK, you can now submit Part 4 to the grader.

Gradient Regularization

Step 9 - For reference see ex4.pdf, top of Page 12, for the right-most terms of the equation for j>=1. Now we calculate the regularization terms for the theta gradients. The goal is that regularization of the gradient should not change the theta gradient(:,1) values (for the bias units) calculated in Step 8. There are several ways to implement this (in Steps 9a and 9b). Method 1: 9a) Calculate the regularization for indexes (:,2:end), and 9b) add them to theta gradients (:,2:end). Method 2: 9a) Calculate the regularization for the entire theta gradient, then overwrite the (:,1) value with 0 before 9b) adding to the entire matrix. Details for Steps 9a and 9b9a) Pick a method, and calculate the regularization terms as follows:$(\lambda/m)*\Theta 1$ (using either Method 1 or Method 2)...and$(\lambda/m)*\Theta 2$ (using either Method 1 or Method 2). 9b) Add these regularization terms to the appropriate Θ1 gradient and Θ2 gradient terms from Step 8 (using either Method 1 or Method 2). Avoid modifying the bias unit of the theta gradients. *Note: there is an errata in the lecture video and slides regarding some missing parenthesis for this calculation. The ex4.pdf file is correct.* The ex4.m script will provide you feedback regarding the acceptable relative difference. If all seems well, you can submit Part 5 to the grader. Now pat yourself on the back.

Appendix:

Here are the sizes for the character recognition example, using the method described in this tutorial. a1: 5000x401z2: 5000x25a2: 5000x26a3: 5000x10d3: 5000x10d2: 5000x25. Theta1, Delta1 and Theta1grad: 25x401; Theta2, Delta2 and Theta2grad: 10x26. Note that the ex4.m script uses a several test cases of different sizes, and the submit grader uses yet another different test case.

# Debugging Tip

The submit script, for all the programming assignments, does not report the line number and location of the error when it crashes. The follow method can be used to make it do so which makes debugging easier.

Open ex4/lib/submitWithConfiguration.m and replace line:

fprintf('!! Please try again later.\n');

(around 28) with:

 fprintf('Error from file:%s\nFunction:%s\nOn line:%d\n', e.stack(1,1).file,e.stack(1,1).name, e.stack(1,1).line );

That top line says '!! Please try again later' on crash, instead of that, the bottom line will give the location and line number of the error. This change can be applied to all the programming assignments.