

# Introduction to Machine Learning Program Assignment #4

Deadline : 01/10(Thu) 23:59

## 1 Neural Networks

In this homework, you will implement the feedforward propagation and backpropagation algorithm to learn the parameters for the neural network.

### 1.1 Visualizing the data

Let's get started by loading the data set. It's in MATLAB's native format, so to load it in Python we need to use a SciPy utility.

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from scipy.io import loadmat

data = loadmat('data/ex3data1.mat')
data['X'].shape, data['y'].shape

Out: ((5000L, 400L), (5000L, 1L))
```

Display it on a 2-dimensional plot (Figure 2) by calling the function `display_network.py`.

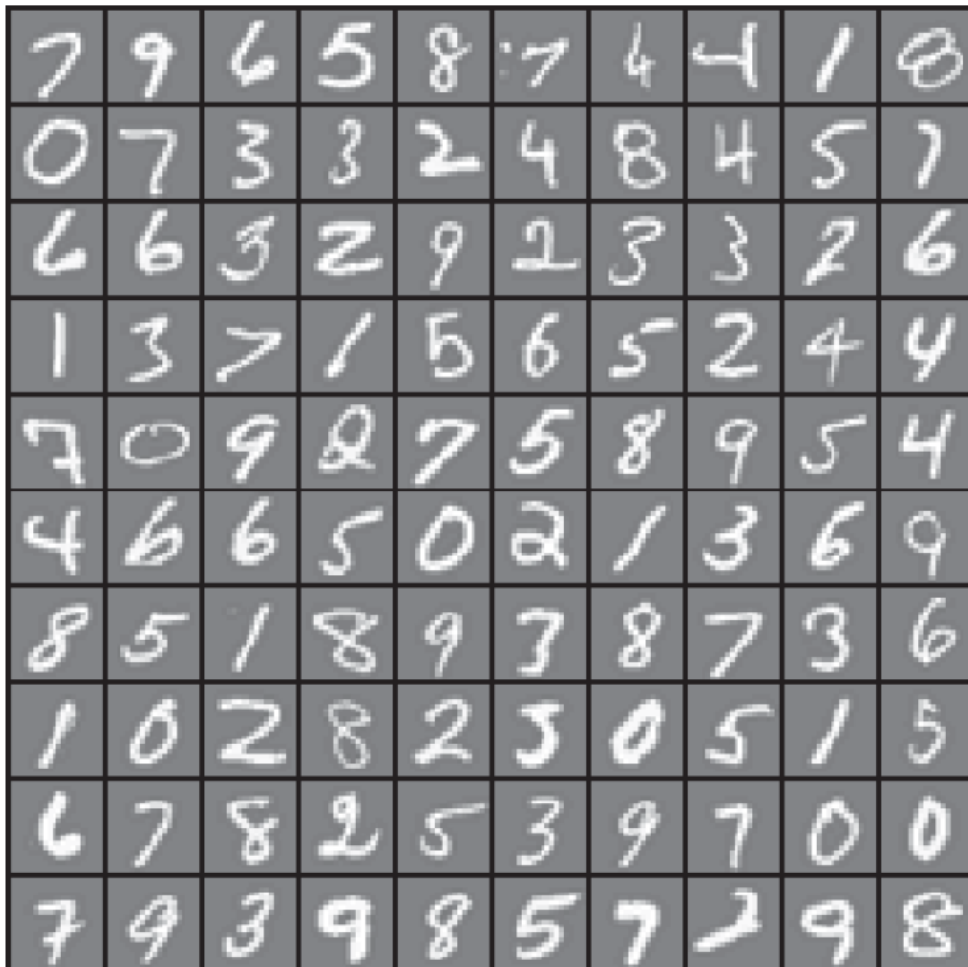


Figure 2: Examples from the dataset.

There are 5000 training examples in `ex3data1.mat`, where each training example is a 20 pixel by 20 pixel grayscale image of the digit. Each pixel is represented by a floating point number indicating the grayscale intensity at that location. The 20 by 20 grid of pixels is “unrolled” into a 400-dimensional vector. Each of these training examples becomes a single row in our data matrix  $X$ . This gives us a 5000 by 400 matrix  $X$  where every row is a training example for a handwritten digit image.

$$X = \begin{bmatrix} -(x^{(1)})^T - \\ -(x^{(2)})^T - \\ \vdots \\ -(x^{(m)})^T - \end{bmatrix}$$

The second part of the training set is a 5000-dimensional vector  $y$  that contains labels for the training set. To make things more compatible with Matlab/Matlab indexing, where there is no zero index, we have mapped the digit zero to the value ten. Therefore, a “0” digit is labeled as “10”, while the digits “1” to “9” are labeled as “1” to “9” in their natural order.

## 1.2 Model representation

Our neural network is shown in Figure 3. It has 3 layers - an input layer, a hidden layer and an output layer. Recall that our inputs are pixel values of digit images. Since the images are of size  $20 \times 20$ , this gives us 400 input layer units (not counting the extra bias unit which always outputs +1). The training data will be loaded into the variables  $X$  and  $y$  by the `ex4.m` script.

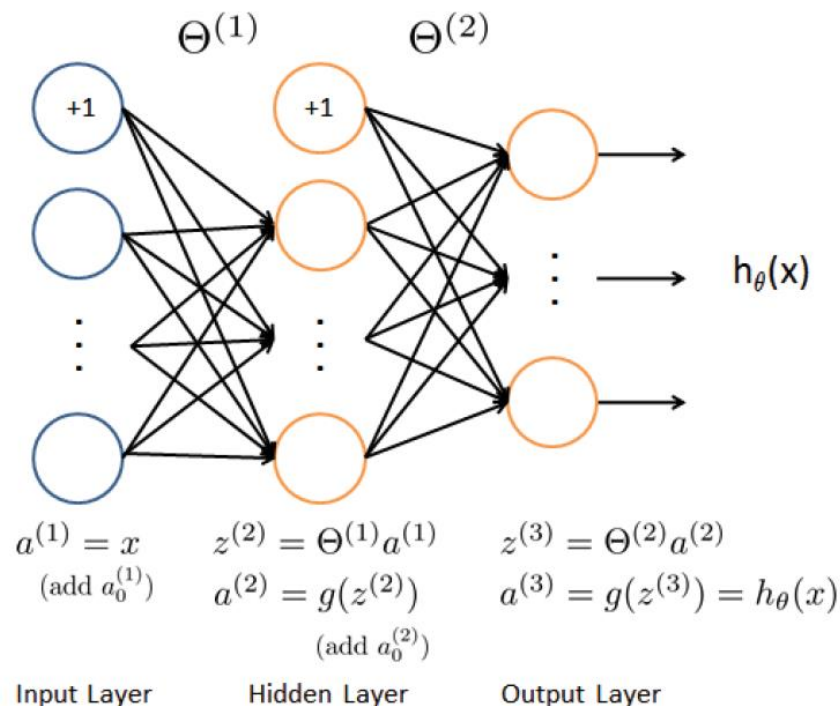


Figure 3: Neural network model.

### 1.3 Feedforward and cost function

Now you will implement the cost function and gradient for the neural network. First, complete the code in `nnCostFunction.m` to return the cost.

Recall that the cost function for the neural network (without regularization) is

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ -y^{(i)} \log((h_{\theta}(x^{(i)}))_k) - (1 - y^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right],$$

where  $h_{\theta}(x^{(i)})$  is computed as shown in the Figure 3 and  $K = 10$  is the total number of possible labels. Note that  $(h_{\theta}(x^{(i)}))_k = a_k^{(3)}$  is the activation (output value) of the  $k$ -th output unit. Also, recall that whereas the original labels (in the variable  $y$ ) were  $1, 2, \dots, 10$ , for the purpose of training a neural network, we need to recode the labels as vectors containing only values 0 or 1, so that

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \quad \dots, \quad \text{or} \quad \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

For example, if  $x^{(i)}$  is an image of the digit 5, then the corresponding  $y^{(i)}$  (that you should use with the cost function) should be a 10-dimensional vector with  $y_5 = 1$ , and the other elements equal to 0.

You should implement the feedforward computation that computes  $h_{\theta}(x^{(i)})$  for every example  $i$  and sum the cost over all examples. **Your code should also work for a dataset of any size, with any number of labels** (you can assume that there are always at least  $K \geq 3$  labels).

```
from sklearn.preprocessing import OneHotEncoder
encoder = OneHotEncoder(sparse=False)
y_onehot = encoder.fit_transform(y)
y_onehot.shape

Out: (5000L, 10L)

y[0], y_onehot[0,:]

Out:
(array([10], dtype=uint8),
 array([ 0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  0.,  1.]))
```

### 1.4 Regularized cost function

The cost function for neural networks with regularization is given by

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \sum_{k=1}^K \left[ -y^{(i)} \log((h_{\theta}(x^{(i)}))_k) - (1 - y^{(i)}) \log(1 - (h_{\theta}(x^{(i)}))_k) \right] \\ + \frac{\lambda}{2m} \left[ \sum_{j=1}^{25} \sum_{k=1}^{400} (\Theta_{j,k}^{(1)})^2 + \sum_{j=1}^{10} \sum_{k=1}^{25} (\Theta_{j,k}^{(2)})^2 \right].$$

You can assume that the neural network will only have 3 layers an input layer, a hidden layer and an output layer. However, your code should work for any number of input units, hidden units and outputs units. While we have explicitly listed the indices above for  $\Theta^{(1)}$  and  $\Theta^{(2)}$  for clarity, do note that your code should in general work with  $\Theta^{(1)}$  and  $\Theta^{(2)}$  of any size.

Note that you should not be regularizing the terms that correspond to the bias. For the matrices **Theta1** and **Theta2**, this corresponds to the first column of each matrix. You should now add regularization to your cost function.

```
def sigmoid(z):  
    return 1 / (1 + np.exp(-z))
```

**def forward\_propagate(X, theta1, theta2):**  
**#write something**

```
def cost(params, input_size, hidden_size, num_labels, X, y, learning_rate):  
    m = X.shape[0]  
    X = np.matrix(X)  
    y = np.matrix(y)  
  
    # reshape the parameter array into parameter matrices for each layer  
    theta1 = np.matrix(np.reshape(params[:hidden_size * (input_size + 1)], (hidden_size, (input_size + 1))  
    theta2 = np.matrix(np.reshape(params[hidden_size * (input_size + 1):], (num_labels, (hidden_size + 1))  
  
    # run the feed-forward pass  
    a1, z2, a2, z3, h = forward_propagate(X, theta1, theta2)  
  
    # compute the cost  
    J = 0  
    for i in range(m):  
        first_term = np.multiply(-y[i,:], np.log(h[i,:]))  
        second_term = np.multiply((1 - y[i,:]), np.log(1 - h[i,:]))  
        J += np.sum(first_term - second_term)  
  
    J = J / m  
  
    return J
```

The forward-propagate function computes the hypothesis for each training instance given the current parameters. It's output shape should match the same of our one-hot encoding for y. We can test this real quick to convince ourselves that it's working as expected (the intermediate steps are also returned as these will be useful later).

## 1.5 Random initialization

When training neural networks, it is important to randomly initialize the parameters for symmetry breaking. One effective strategy for random initialization is to randomly select values for  $\Theta^{(l)}$  uniformly in the range  $[-\epsilon_{\text{init}}, \epsilon_{\text{init}}]$ . You should use  $\epsilon_{\text{init}} = 0.12$ .<sup>1</sup> This range of values ensures that the parameters are kept small and makes the learning more efficient.

```
# initial setup
input_size = 400
hidden_size = 25
num_labels = 10
learning_rate = 1

# randomly initialize a parameter array of the size of the full network's parameters
params = (np.random.random(size=hidden_size * (input_size + 1) + num_labels * (hidden_size + 1)) - 0.5)

m = X.shape[0]
X = np.matrix(X)
y = np.matrix(y)

# unravel the parameter array into parameter matrices for each layer
theta1 = np.matrix(np.reshape(params[:hidden_size * (input_size + 1)], (hidden_size, (input_size + 1))))
theta2 = np.matrix(np.reshape(params[hidden_size * (input_size + 1):], (num_labels, (hidden_size + 1))))

theta1.shape, theta2.shape

Out: ((25L, 401L), (10L, 26L))

a1, z2, a2, z3, h = forward_propagate(X, theta1, theta2)
a1.shape, z2.shape, a2.shape, z3.shape, h.shape

Out: ((5000L, 401L), (5000L, 25L), (5000L, 26L), (5000L, 10L), (5000L, 10L))
```

The cost function, after computing the hypothesis matrix  $h$ , applies the cost equation to compute the total error between  $y$  and  $h$ .

```

def cost(params, input_size, hidden_size, num_labels, X, y, learning_rate):
    m = X.shape[0]
    X = np.matrix(X)
    y = np.matrix(y)

    # reshape the parameter array into parameter matrices for each layer
    theta1 = np.matrix(np.reshape(params[:hidden_size * (input_size + 1)], (hidden_size, (input_size + 1))
    theta2 = np.matrix(np.reshape(params[hidden_size * (input_size + 1):], (num_labels, (hidden_size + 1))

    # run the feed-forward pass
    a1, z2, a2, z3, h = forward_propagate(X, theta1, theta2)

    # compute the cost
    J = 0

    for i in range(m):
        first_term = np.multiply(-y[i,:], np.log(h[i,:]))
        second_term = np.multiply((1 - y[i,:]), np.log(1 - h[i,:]))
        J += np.sum(first_term - second_term)

    J = J / m

    # add the cost regularization term
    J += (float(learning_rate) / (2 * m)) * (np.sum(np.power(theta1[:,1:], 2)) + np.sum(np.power(theta2[:,1:], 2)))

    return J

```

```
cost(params, input_size, hidden_size, num_labels, X, y_onehot, learning_rate)
```

```
Out: 6.8281541822949299
```

## 2 Backpropagation

In this part of the homework, you will implement the backpropagation algorithm to compute the gradient for the neural network cost function. Once you have computed the gradient, you will be able to train the neural network by minimizing the cost function  $J(\Theta)$  using an advanced optimizer such as `fmincg`.

You will first implement the backpropagation algorithm to compute the gradients for the parameters for the (unregularized) neural network. After you have verified that your gradient computation for the unregularized case is correct, you will implement the gradient for the regularized neural network.

### 2.1 Sigmoid gradient

To help you get started with this part of the homework, you will first implement the sigmoid gradient function. The gradient for the sigmoid function can be computed as

$$g'(z) = \frac{d}{dz}g(z) = g(z)(1 - g(z))$$

where

$$\text{sigmoid}(z) = g(z) = \frac{1}{1 + e^{-z}}.$$

```
def sigmoid_gradient(z):
    return np.multiply(sigmoid(z), (1 - sigmoid(z)))
```

### 2.2 Backpropagation

Now, you will implement the backpropagation algorithm. Recall that the intuition behind the backpropagation algorithm is as follows. Given a training example  $(x^{(t)}, y^{(t)})$ , we will first run a “forward pass” to compute all the activations throughout the network, including the output value of the hypothesis  $h_{\Theta}(x)$ . Then, for each node  $j$  in layer  $l$ , we would like to compute an “error term”  $\delta_j^{(l)}$  that measures how much that node was “responsible” for any errors in our output.

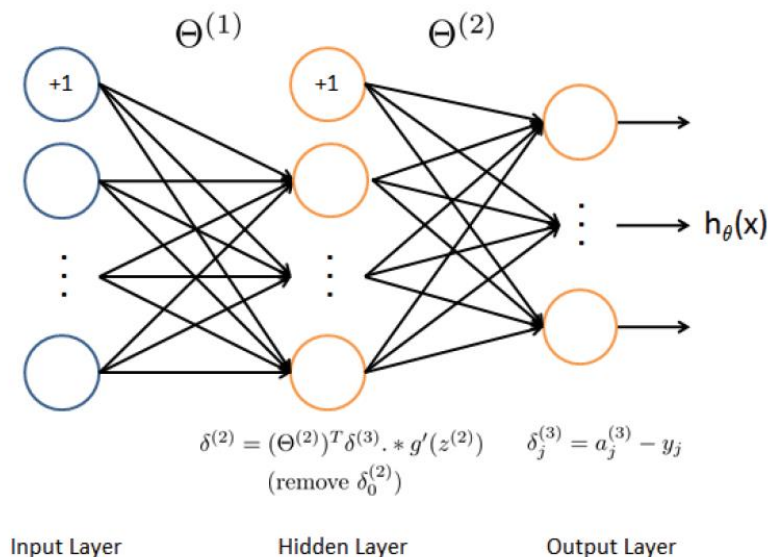


Figure 4: Backpropagation Updates.



For an output node, we can directly measure the difference between the network's activation and the true target value, and use that to define  $\delta_j^{(3)}$  (since layer 3 is the output layer). For the hidden units, you will compute  $\delta_j^{(l)}$  based on a weighted average of the error terms of the nodes in layer  $(l + 1)$ .

In detail, here is the backpropagation algorithm (also depicted in Figure 4). You should implement steps 1 to 4 in a loop that processes one example at a time. Concretely, you should implement a for-loop for  $\mathbf{t} = 1:m$  and place steps 1-4 below inside the for-loop, with the  $t$ -th iteration performing the calculation on the  $t$ -th training example  $(x^{(t)}, y^{(t)})$ . Step 5 will divide the accumulated gradients by  $m$  to obtain the gradients for the neural network cost function.

1. Set the input layer's values  $(a^{(1)})$  to the  $t$ -th training example  $x^{(t)}$ . Perform a feedforward pass (Figure 3), computing the activations  $(z^{(2)}, a^{(2)}, z^{(3)}, a^{(3)})$  for layers 2 and 3. Note that you need to add  $a + 1$  term to ensure that the vectors of activations for layers  $a^{(1)}$  and  $a^{(2)}$  also include the bias unit. In Matlab, if `a_1` is a column vector, adding one corresponds to `a_1 = [1 ; a_1]`.
2. For each output unit  $k$  in layer 3 (the output layer), set

$$\delta_k^{(3)} = (a_k^{(3)} - y_k),$$

where  $y_k \in \{0, 1\}$  indicates whether the current training example belongs to class  $k$  ( $y_k = 1$ ), or if it belongs to a different class ( $y_k = 0$ ). You may find logical arrays helpful for this task (explained in the previous programming homework).

3. For the hidden layer  $l = 2$ , set

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} .* g'(z^{(2)}).$$

4. Accumulate the gradient from this example using the following formula. Note that you should skip or remove  $\delta_0^{(2)}$ . In Matlab, removing  $\delta_0^{(2)}$  corresponds to `delta_2 = delta_2(2:end)`.

$$\Delta^{(l)} = \Delta^{(l)} + \delta^{(l+1)}(a^{(l)})^T$$

5. Obtain the (unregularized) gradient for the neural network cost function by multiplying the accumulated gradients by  $\frac{1}{m}$ :

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)}$$



## 2.4 Regularized Neural Networks

After you have successfully implemented the backpropagation algorithm, you will add regularization to the gradient. To account for regularization, it turns out that you can add this as an additional term *after* computing the gradients using backpropagation.

Specifically, after you have computed  $\Delta_{ij}^{(l)}$  using backpropagation, you should add regularization using

$$\begin{aligned}\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) &= D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} && \text{for } j = 0 \\ \frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) &= D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} + \frac{\lambda}{m} \Theta_{ij}^{(l)} && \text{for } j \geq 1\end{aligned}$$

Note that you should *not* be regularizing the first column of  $\Theta^{(l)}$  which is used for the bias term. Furthermore, in the parameters  $\Theta_{ij}^{(l)}$ ,  $i$  is indexed starting from 1, and  $j$  is indexed starting from 0. Thus,

$$\Theta^{(l)} = \begin{bmatrix} \Theta_{10}^{(l)} & \Theta_{11}^{(l)} & \cdots \\ \Theta_{20}^{(l)} & \Theta_{21}^{(l)} & \cdots \\ \vdots & & \ddots \end{bmatrix}.$$

Somewhat confusingly, indexing in Matlab starts from 1 (for both  $i$  and  $j$ ), thus `Theta1(2, 1)` actually corresponds to  $\Theta_{20}^{(l)}$  (i.e., the entry in the second row, first column of the matrix  $\Theta^{(l)}$  shown above)

**def backprop(params, input\_size, hidden\_size, num\_labels, X, y, learning\_rate):**  
**#write something**

This is a **person-based** program assignment, so **one student should only submit one source code and one report to E3**.

The code should contain the follow

1. Forward-propagate code (40%)
2. Back-propagate code (60%)

The Report should contain the follow

1. **Screenshot** of Forward-propagate code
2. **Screenshot** of Back-propagate code
3. If you finish these two functions, you need to run all codes and put the **screenshot of accuracy** to report.

Some rules

1. File name **student ID.py** + **student ID.pdf**
2. Suggest to use Python.
3. You can change all codes to C/C++/Java...
4. No cheating and plagiarizing.
5. **Deadline : 01/10(Thu) 23:59**
6. **Delay : Your score = 0**

**TA will not help you to debug and explain the way of solving the problem.**