

Programming Exercise 3

Multi-class Classification and Neural Networks

Introduction

In this exercise, you will implement one-vs-all logistic regression and neural networks to recognize handwritten digits.

All the information you need for solving this assignment is in this notebook, and all the code you will be implementing will take place within this notebook.

Before we begin with the exercises, we need to import all libraries required for this programming exercise. Throughout the course, we will be using `numpy` for all arrays and matrix operations, `matplotlib` for plotting, and `scipy` for scientific and numerical computation functions and tools.

```
In [3]: # used for manipulating directory paths
import os

# Scientific and vector computation for python
import numpy as np

# Plotting library
from matplotlib import pyplot

# Optimization module in scipy
from scipy import optimize

# will be used to load MATLAB mat datafile format
from scipy.io import loadmat

# Library written for this exercise providing additional functions for assignment submission
import utils

# tells matplotlib to embed plots within the notebook
%matplotlib inline
```

Submission and Grading

The following is a breakdown of each part of this exercise.

Section	Part	Submission function
1	Regularized Logistic Regression	<code>lrCostFunction</code>
2	One-vs-all classifier training	<code>oneVsAll</code>
3	One-vs-all classifier prediction	<code>predictOneVsAll</code>

Section	Part	Submission function
4	Neural Network Prediction Function	predict

1 Multi-class Classification

For this exercise, you will use logistic regression and neural networks to recognize handwritten digits (from 0 to 9). Automated handwritten digit recognition is widely used today - from recognizing zip codes (postal codes) on mail envelopes to recognizing amounts written on bank checks. This exercise will show you how the methods you have learned can be used for this classification task.

In the first part of the exercise, you will extend your previous implementation of logistic regression and apply it to one-vs-all classification.

1.1 Dataset

You are given a data set in `ex3data1.mat` that contains 5000 training examples of handwritten digits (This is a subset of the [MNIST](#) handwritten digit dataset). The `.mat` format means that the data has been saved in a native Octave/MATLAB matrix format, instead of a text (ASCII) format like a csv-file. We use the `.mat` format here because this is the dataset provided in the MATLAB version of this assignment. Fortunately, python provides mechanisms to load MATLAB native format using the `loadmat` function within the `scipy.io` module. This function returns a python dictionary with keys containing the variable names within the `.mat` file.

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. We start the exercise by first loading the dataset. Execute the cell below, you do not need to write any code here.

```
In [4]: # 20x20 Input Images of Digits
input_layer_size = 400

# 10 Labels, from 1 to 10 (note that we have mapped "0" to Label 10)
num_labels = 10

# training data stored in arrays X, y
```

```
data = loadmat(os.path.join('Data', 'ex3data1.mat'))
X, y = data['X'], data['y'].ravel()

# set the zero digit to 0, rather than its mapped 10 in this dataset
# This is an artifact due to the fact that this dataset was used in
# MATLAB where there is no index 0
y[y == 10] = 0

m = y.size
```

1.2 Visualizing the data

You will begin by visualizing a subset of the training set. In the following cell, the code randomly selects 100 rows from `X` and passes those rows to the `displayData` function. This function maps each row to a 20 pixel by 20 pixel grayscale image and displays the images together. We have provided the `displayData` function in the file `utils.py`. You are encouraged to examine the code to see how it works. Run the following cell to visualize the data.

```
In [5]: # Randomly select 100 data points to display
rand_indices = np.random.choice(m, 100, replace=False)
sel = X[rand_indices, :]

utils.displayData(sel)
```



1.3 Vectorizing Logistic Regression

You will be using multiple one-vs-all logistic regression models to build a multi-class classifier. Since there are 10 classes, you will need to train 10 separate logistic regression classifiers. To make this training efficient, it is important to ensure that your code is well vectorized. In this section, you will implement a vectorized version of logistic regression that does not employ any `for` loops. You can use your code in the previous exercise as a starting point for this exercise.

To test your vectorized logistic regression, we will use custom data as defined in the following cell.

```
In [6]: # test values for the parameters theta  
theta_t = np.array([-2, -1, 1, 2], dtype=float)  
  
# test values for the inputs  
X_t = np.concatenate([np.ones((5, 1)), np.arange(1, 16).reshape(5, 3, order='F')/16])  
  
# test values for the labels  
y_t = np.array([1, 0, 1, 0, 1])
```

```
# test value for the regularization parameter
lambda_t = 3
```

1.3.1 Vectorizing the cost function

We will begin by writing a vectorized version of the cost function. Recall that in (unregularized) logistic regression, the cost function is

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

To compute each element in the summation, we have to compute $h_\theta(x^{(i)})$ for every example i , where $h_\theta(x^{(i)}) = g(\theta^T x^{(i)})$ and $g(z) = \frac{1}{1+e^{-z}}$ is the sigmoid function. It turns out that we can compute this quickly for all our examples by using matrix multiplication. Let us define X and θ as

$$X = \begin{bmatrix} -(x^{(1)})^T - \\ -(x^{(2)})^T - \\ \vdots \\ -(x^{(m)})^T - \end{bmatrix} \quad \text{and} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

Then, by computing the matrix product $X\theta$, we have:

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

In the last equality, we used the fact that $a^T b = b^T a$ if a and b are vectors. This allows us to compute the products $\theta^T x^{(i)}$ for all our examples i in one line of code.

1.3.2 Vectorizing the gradient

Recall that the gradient of the (unregularized) logistic regression cost is a vector where the j^{th} element is defined as

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m \left((h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \right)$$

To vectorize this operation over the dataset, we start by writing out all the partial derivatives explicitly for all θ_j ,

$$\begin{aligned}
\left[\begin{array}{c} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{array} \right] &= \frac{1}{m} \left[\begin{array}{c} \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}) \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}) \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}) \\ \vdots \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)}) x_n^{(i)}) \end{array} \right] \\
&= \frac{1}{m} \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)}) x^{(i)}) \\
&= \frac{1}{m} X^T (h_\theta(x) - y)
\end{aligned}$$

where

$$h_\theta(x) - y = \begin{bmatrix} h_\theta(x^{(1)}) - y^{(1)} \\ h_\theta(x^{(2)}) - y^{(2)} \\ \vdots \\ h_\theta(x^{(m)}) - y^{(m)} \end{bmatrix}$$

Note that $x^{(i)}$ is a vector, while $h_\theta(x^{(i)}) - y^{(i)}$ is a scalar (single number). To understand the last step of the derivation, let $\beta_i = (h_\theta(x^{(m)}) - y^{(m)})$ and observe that:

$$\sum_i \beta_i x^{(i)} = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_m \end{bmatrix} = x^T \beta$$

where the values $\beta_i = (h_\theta(x^{(i)}) - y^{(i)})$.

The expression above allows us to compute all the partial derivatives without any loops. If you are comfortable with linear algebra, we encourage you to work through the matrix multiplications above to convince yourself that the vectorized version does the same computations.

Your job is to write the unregularized cost function `lrCostFunction` which returns both the cost function $J(\theta)$ and its gradient $\frac{\partial J}{\partial \theta}$. Your implementation should use the strategy we presented above to calculate $\theta^T x^{(i)}$. You should also use a vectorized approach for the rest of the cost function. A fully vectorized version of `lrCostFunction` should not contain any loops.

****Debugging Tip:**** Vectorizing code can sometimes be tricky. One common strategy for debugging is to print out the sizes of the matrices you are working with using the `'shape'` property of `'numpy'` arrays. For example, given a data matrix X of size 100×20 (100 examples, 20 features) and θ , a vector with size 20, you can observe that `'np.dot(X, theta)'` is a valid multiplication operation, while `'np.dot(theta, X)'` is not. Furthermore, if you have a non-vectorized version of your code, you can compare the output of your

vectorized code and non-vectorized code to make sure that they produce the same outputs.

```
In [7]: def lrCostFunction(theta, X, y, lambda_):
    """
    Computes the cost of using theta as the parameter for regularized
    logistic regression and the gradient of the cost w.r.t. to the parameters.

    Parameters
    -----
    theta : array_like
        Logistic regression parameters. A vector with shape (n, ). n is
        the number of features including any intercept.

    X : array_like
        The data set with shape (m x n). m is the number of examples, and
        n is the number of features (including intercept).

    y : array_like
        The data labels. A vector with shape (m, ).

    lambda_ : float
        The regularization parameter.

    Returns
    -----
    J : float
        The computed value for the regularized cost function.

    grad : array_like
        A vector of shape (n, ) which is the gradient of the cost
        function with respect to theta, at the current values of theta.

    Instructions
    -----
    Compute the cost of a particular choice of theta. You should set J to the cost.
    Compute the partial derivatives and set grad to the partial
    derivatives of the cost w.r.t. each parameter in theta

    Hint 1
    -----
    The computation of the cost function and gradients can be efficiently
    vectorized. For example, consider the computation

        sigmoid(X * theta)

    Each row of the resulting matrix will contain the value of the prediction
    for that example. You can make use of this to vectorize the cost function
    and gradient computations.

    Hint 2
    -----
    When computing the gradient of the regularized cost function, there are
    many possible vectorized solutions, but one solution looks like:

        grad = (unregularized gradient for logistic regression)
        temp = theta
        temp[0] = 0  # because we don't add anything for j = 0
        grad = grad + YOUR_CODE_HERE (using the temp variable)
```

```

Hint 3
-----
We have provided the implementation of the sigmoid function within
the file `utils.py`. At the start of the notebook, we imported this file
as a module. Thus to access the sigmoid function within that file, you can
do the following: `utils.sigmoid(z)`.

"""
#Initialize some useful values
m = y.size

# convert Labels to ints if their type is bool
if y.dtype == bool:
    y = y.astype(int)

# You need to return the following variables correctly
J = 0
grad = np.zeros(theta.shape)

# ===== YOUR CODE HERE =====

h = utils.sigmoid(X.dot(theta.T))

temp = theta
temp[0] = 0

J = (1 / m) * np.sum(-y.dot(np.log(h)) - (1 - y).dot(np.log(1 - h))) + (lambda_

grad = (1 / m) * (h - y).dot(X)
grad = grad + (lambda_ / m) * temp

# =====#
return J, grad

```

1.3.3 Vectorizing regularized logistic regression

After you have implemented vectorization for logistic regression, you will now add regularization to the cost function. Recall that for regularized logistic regression, the cost function is defined as

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

Note that you should not be regularizing θ_0 which is used for the bias term.

Correspondingly, the partial derivative of regularized logistic regression cost for θ_j is defined as

$$\begin{aligned} \frac{\partial J(\theta)}{\partial \theta_0} &= \frac{1}{m} \sum_{i=1}^m \left(h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x_j^{(i)} && \text{for } j = 0 \\ \frac{\partial J(\theta)}{\partial \theta_j} &= \left(\frac{1}{m} \sum_{i=1}^m \left(h_\theta\left(x^{(i)}\right) - y^{(i)} \right) x_j^{(i)} \right) + \frac{\lambda}{m} \theta_j && \text{for } j \geq 1 \end{aligned}$$

Now modify your code in `lrCostFunction` in the [previous cell](#) to account for regularization. Once again, you should not put any loops into your code.

****python/numpy Tip:**** When implementing the vectorization for regularized logistic regression, you might often want to only sum and update certain elements of θ . In `numpy`, you can index into the matrices to access and update only certain elements. For example, $A[:, 3:5] = B[:, 1:3]$ will replace the columns with index 3 to 5 of A with the columns with index 1 to 3 from B . To select columns (or rows) until the end of the matrix, you can leave the right hand side of the colon blank. For example, $A[:, 2:]$ will only return elements from the 3rd to last columns of A . If you leave the left hand side of the colon blank, you will select elements from the beginning of the matrix. For example, $A[:, :2]$ selects the first two columns, and is equivalent to $A[:, 0:2]$. In addition, you can use negative indices to index arrays from the end. Thus, $A[:, :-1]$ selects all columns of A except the last column, and $A[:, -5:]$ selects the 5th column from the end to the last column. Thus, you could use this together with the sum and power (***) operations to compute the sum of only the elements you are interested in (e.g., `np.sum(z[1:]**2)`). In the starter code, `lrCostFunction`, we have also provided hints on yet another possible method computing the regularized gradient.

Once you finished your implementation, you can call the function `lrCostFunction` to test your solution using the following cell:

```
In [8]: J, grad = lrCostFunction(theta_t, X_t, y_t, lambda_t)

print('Cost          : {:.6f}'.format(J))
print('Expected cost: 2.534819')
print('-----')
print('Gradients: ')
print(' [{:.6f}, {:.6f}, {:.6f}, {:.6f}]'.format(*grad))
print('Expected gradients: ')
print(' [0.146561, -0.548558, 0.724722, 1.398003]');
```

Cost : 2.534819
 Expected cost: 2.534819

 Gradients:
 [0.146561, -0.548558, 0.724722, 1.398003]
 Expected gradients:
 [0.146561, -0.548558, 0.724722, 1.398003]

1.4 One-vs-all Classification

In this part of the exercise, you will implement one-vs-all classification by training multiple regularized logistic regression classifiers, one for each of the K classes in our dataset. In the handwritten digits dataset, $K = 10$, but your code should work for any value of K .

You should now complete the code for the function `oneVsAll` below, to train one classifier for each class. In particular, your code should return all the classifier parameters in a matrix $\theta \in \mathbb{R}^{K \times (N+1)}$, where each row of θ corresponds to the learned logistic regression parameters for one class. You can do this with a "for"-loop from 0 to $K - 1$, training each classifier independently.

Note that the `y` argument to this function is a vector of labels from 0 to 9. When training the classifier for class $k \in \{0, \dots, K - 1\}$, you will want a K -dimensional vector of labels y ,

where $y_j \in 0, 1$ indicates whether the j^{th} training instance belongs to class k ($y_j = 1$), or if it belongs to a different class ($y_j = 0$). You may find logical arrays helpful for this task.

Furthermore, you will be using scipy's `optimize.minimize` for this exercise.

In [9]:

```
def oneVsAll(X, y, num_labels, lambda_):
    """
    Trains num_labels logistic regression classifiers and returns
    each of these classifiers in a matrix all_theta, where the i-th
    row of all_theta corresponds to the classifier for label i.

    Parameters
    -----
    X : array_like
        The input dataset of shape (m x n). m is the number of
        data points, and n is the number of features. Note that we
        do not assume that the intercept term (or bias) is in X, however
        we provide the code below to add the bias term to X.

    y : array_like
        The data labels. A vector of shape (m, ).

    num_labels : int
        Number of possible labels.

    lambda_ : float
        The logistic regularization parameter.

    Returns
    -----
    all_theta : array_like
        The trained parameters for logistic regression for each class.
        This is a matrix of shape (K x n+1) where K is number of classes
        (ie. `numlabels`) and n is number of features without the bias.

    Instructions
    -----
    You should complete the following code to train `num_labels`
    logistic regression classifiers with regularization parameter `lambda_`.

    Hint
    -----
    You can use y == c to obtain a vector of 1's and 0's that tell you
    whether the ground truth is true/false for this class.

    Note
    -----
    For this assignment, we recommend using `scipy.optimize.minimize(method='CG')`
    to optimize the cost function. It is okay to use a for-loop
    (`for c in range(num_labels):`) to loop over the different classes.

    Example Code
    -----
    # Set Initial theta
    initial_theta = np.zeros(n + 1)

    # Set options for minimize
    options = {'maxiter': 50}
```

```

# Run minimize to obtain the optimal theta. This function will
# return a class object where theta is in `res.x` and cost in `res.fun`
res = optimize.minimize(lrCostFunction,
                        initial_theta,
                        (X, (y == c), lambda_),
                        jac=True,
                        method='TNC',
                        options=options)
"""

# Some useful variables
m, n = X.shape

# You need to return the following variables correctly
all_theta = np.zeros((num_labels, n + 1))

# Add ones to the X data matrix
X = np.concatenate([np.ones((m, 1)), X], axis=1)

# ===== YOUR CODE HERE =====
for c in np.arange(num_labels):
    initial_theta = np.zeros(n + 1)
    options = {'maxiter': 50}
    res = optimize.minimize(lrCostFunction,
                            initial_theta,
                            (X, (y == c), lambda_),
                            jac=True,
                            method='CG',
                            options=options)

    all_theta[c] = res.x

# =====
return all_theta

```

After you have completed the code for `oneVsAll`, the following cell will use your implementation to train a multi-class classifier.

In [10]:

```
lambda_ = 0.1
all_theta = oneVsAll(X, y, num_labels, lambda_)
```

1.4.1 One-vs-all Prediction

After training your one-vs-all classifier, you can now use it to predict the digit contained in a given image. For each input, you should compute the “probability” that it belongs to each class using the trained logistic regression classifiers. Your one-vs-all prediction function will pick the class for which the corresponding logistic regression classifier outputs the highest probability and return the class label (0, 1, ..., K-1) as the prediction for the input example. You should now complete the code in the function `predictOneVsAll` to use the one-vs-all classifier for making predictions.

In [11]:

```
def predictOneVsAll(all_theta, X):
"""
Return a vector of predictions for each example in the matrix X.
Note that X contains the examples in rows. all_theta is a matrix where
the i-th row is a trained logistic regression theta vector for the
```

```

i-th class. You should set p to a vector of values from 0..K-1
(e.g., p = [0, 2, 0, 1] predicts classes 0, 2, 0, 1 for 4 examples) .

Parameters
-----
all_theta : array_like
    The trained parameters for logistic regression for each class.
    This is a matrix of shape (K x n+1) where K is number of classes
    and n is number of features without the bias.

X : array_like
    Data points to predict their labels. This is a matrix of shape
    (m x n) where m is number of data points to predict, and n is number
    of features without the bias term. Note we add the bias term for X in
    this function.

Returns
-----
p : array_like
    The predictions for each data point in X. This is a vector of shape (m, ).

Instructions
-----
Complete the following code to make predictions using your learned logistic
regression parameters (one-vs-all). You should set p to a vector of predictions
(from 0 to num_labels-1).

Hint
-----
This code can be done all vectorized using the numpy argmax function.
In particular, the argmax function returns the index of the max element,
for more information see '?np.argmax' or search online. If your examples
are in rows, then, you can use np.argmax(A, axis=1) to obtain the index
of the max for each row.

"""
m = X.shape[0];
num_labels = all_theta.shape[0]

# You need to return the following variables correctly
p = np.zeros(m)

# Add ones to the X data matrix
X = np.concatenate([np.ones((m, 1)), X], axis=1)

# ===== YOUR CODE HERE =====

p = np.argmax(utils.sigmoid(X.dot(all_theta.T)), axis = 1)
# =====
return p

```

Once you are done, call your `predictOneVsAll` function using the learned value of θ . You should see that the training set accuracy is about 95.1% (i.e., it classifies 95.1% of the examples in the training set correctly).

```
In [19]: pred = predictOneVsAll(all_theta, X)
print('Training Set Accuracy: {:.2f}%'.format(np.mean(pred == y) * 100))
```

Training Set Accuracy: 18.52%

2 Neural Networks

In the previous part of this exercise, you implemented multi-class logistic regression to recognize handwritten digits. However, logistic regression cannot form more complex hypotheses as it is only a linear classifier (You could add more features - such as polynomial features - to logistic regression, but that can be very expensive to train).

In this part of the exercise, you will implement a neural network to recognize handwritten digits using the same training set as before. The neural network will be able to represent complex models that form non-linear hypotheses. For this week, you will be using parameters from a neural network that we have already trained. Your goal is to implement the feedforward propagation algorithm to use our weights for prediction. In next week's exercise, you will write the backpropagation algorithm for learning the neural network parameters.

We start by first reloading and visualizing the dataset which contains the MNIST handwritten digits (this is the same as we did in the first part of this exercise, we reload it here to ensure the variables have not been modified).

In [12]:

```
# training data stored in arrays X, y
data = loadmat(os.path.join('Data', 'ex3data1.mat'))
X, y = data['X'], data['y'].ravel()

# set the zero digit to 0, rather than its mapped 10 in this dataset
# This is an artifact due to the fact that this dataset was used in
# MATLAB where there is no index 0
y[y == 10] = 0

# get number of examples in dataset
m = y.size

# randomly permute examples, to be used for visualizing one
# picture at a time
indices = np.random.permutation(m)

# Randomly select 100 data points to display
rand_indices = np.random.choice(m, 100, replace=False)
sel = X[rand_indices, :]

utils.displayData(sel)
```



2.1 Model representation

Our neural network is shown in the following figure.



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 20×20 , this gives us 400 input layer units (excluding the extra bias unit which always outputs +1). As before, the training data will be loaded into the variables X and y .

You have been provided with a set of network parameters $(\Theta^{(1)}, \Theta^{(2)})$ already trained by us. These are stored in `ex3weights.mat`. The following cell loads those parameters into `Theta1` and `Theta2`. The parameters have dimensions that are sized for a neural network with 25 units in the second layer and 10 output units (corresponding to the 10 digit classes).

```
In [15]: # Setup the parameters you will use for this exercise
input_layer_size = 400 # 20x20 Input Images of Digits
hidden_layer_size = 25 # 25 hidden units
num_labels = 10 # 10 Labels, from 0 to 9

# Load the .mat file, which returns a dictionary
```

```

weights = loadmat(os.path.join('Data', 'ex3weights.mat'))

# get the model weights from the dictionary
# Theta1 has size 25 x 401
# Theta2 has size 10 x 26
Theta1, Theta2 = weights['Theta1'], weights['Theta2']

# swap first and last columns of Theta2, due to legacy from MATLAB indexing,
# since the weight file ex3weights.mat was saved based on MATLAB indexing
Theta2 = np.roll(Theta2, 1, axis=0)

```

2.2 Feedforward Propagation and Prediction

Now you will implement feedforward propagation for the neural network. You will need to complete the code in the function `predict` to return the neural network's prediction. You should implement the feedforward computation that computes $h_{\theta}(x^{(i)})$ for every example i and returns the associated predictions. Similar to the one-vs-all classification strategy, the prediction from the neural network will be the label that has the largest output $(h_{\theta}(x))_k$.

****Implementation Note:**** The matrix X contains the examples in rows. When you complete the code in the function `predict`, you will need to add the column of 1's to the matrix. The matrices `Theta1` and `Theta2` contain the parameters for each unit in rows. Specifically, the first row of `Theta1` corresponds to the first hidden unit in the second layer. In `numpy`, when you compute $z^{(2)} = \theta^{(1)}a^{(1)}$, be sure that you index (and if necessary, transpose) X correctly so that you get $a^{(l)}$ as a 1-D vector.

```

In [16]: def predict(Theta1, Theta2, X):
    """
    Predict the label of an input given a trained neural network.

    Parameters
    -----
    Theta1 : array_like
        Weights for the first layer in the neural network.
        It has shape (2nd hidden layer size x input size)

    Theta2: array_like
        Weights for the second layer in the neural network.
        It has shape (output layer size x 2nd hidden layer size)

    X : array_like
        The image inputs having shape (number of examples x image dimensions).

    Return
    -----
    p : array_like
        Predictions vector containing the predicted label for each example.
        It has a length equal to the number of examples.

    Instructions
    -----
    Complete the following code to make predictions using your learned neural
    network. You should set p to a vector containing labels
    between 0 to (num_labels-1).

```

```

Hint
-----
This code can be done all vectorized using the numpy argmax function.
In particular, the argmax function returns the index of the max element,
for more information see '?np.argmax' or search online. If your examples
are in rows, then, you can use np.argmax(A, axis=1) to obtain the index
of the max for each row.

Note
-----
Remember, we have supplied the `sigmoid` function in the `utils.py` file.
You can use this function by calling `utils.sigmoid(z)`, where you can
replace `z` by the required input variable to sigmoid.

"""
# Make sure the input has two dimensions
if X.ndim == 1:
    X = X[None] # promote to 2-dimensions

# useful variables
m = X.shape[0]
num_labels = Theta2.shape[0]

# You need to return the following variables correctly
p = np.zeros(X.shape[0])

# ===== YOUR CODE HERE =====
X = np.concatenate([np.ones((m, 1)), X], axis=1)

a2 = utils.sigmoid(X.dot(Theta1.T))
a2 = np.concatenate([np.ones((a2.shape[0], 1)), a2], axis=1)

p = np.argmax(utils.sigmoid(a2.dot(Theta2.T)), axis = 1)

# =====
return p

```

Once you are done, call your predict function using the loaded set of parameters for `Theta1` and `Theta2`. You should see that the accuracy is about 97.5%.

```
In [17]: pred = predict(Theta1, Theta2, X)
print('Training Set Accuracy: {:.1f}%'.format(np.mean(pred == y) * 100))
```

Training Set Accuracy: 97.5%

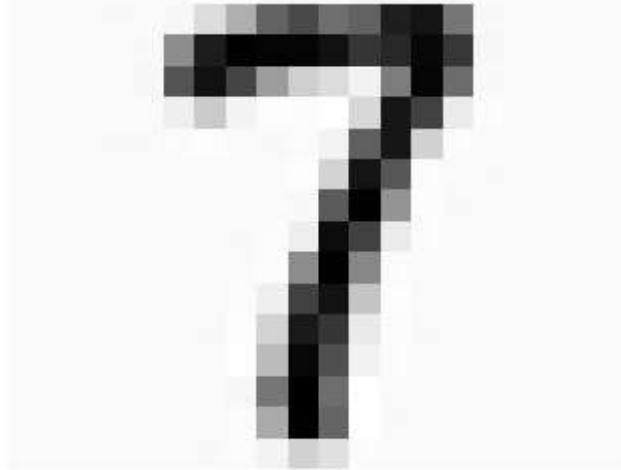
After that, we will display images from the training set one at a time, while at the same time printing out the predicted label for the displayed image.

Run the following cell to display a single image the the neural network's prediction. You can run the cell multiple time to see predictions for different images.

```
In [18]: if indices.size > 0:
    i, indices = indices[0], indices[1:]
    utils.displayData(X[i, :], figsize=(4, 4))
    pred = predict(Theta1, Theta2, X[i, :])
    print('Neural Network Prediction: {}'.format(*pred))
```

```
else:  
    print('No more images to display!')
```

Neural Network Prediction: 7



In []: