

Programming Exercise 3:

Multi-class Classification and Neural Networks

Machine Learning

Introduction

In this exercise, you will implement one-vs-all logistic regression and neural networks to recognize hand-written digits. Before starting on this programming exercise, we strongly recommend watching the video lectures and completing the review questions for the associated topics.

To get started with the exercise, you will need to download the starter code and unzip its contents to the directory where you wish to complete the exercise.

Files included in this exercise

ex3.py - Python script that steps you through part 1

ex3_nn.py - Python script that steps you through part 2

ex3data1.mat - Training set of hand-written digits

ex3weights.mat - Initial weights for the neural network exercise

displayData.py - Function to help visualize the dataset

[★] lrCostFunction.py - Logistic regression cost function

[★] oneVsAll.py - Train a one-vs-all multi-class classifier

[★] predictOneVsAll.py - Predict using a one-vs-all multi-class classifier

[★] predict.py - Neural network prediction function

★ indicates files you will need to complete

Throughout the exercise, you will be using the scripts `ex3.py` and `ex3_nn.py`. These scripts set up the dataset for the problems and make calls to functions that you will write. You do not need to modify these scripts. You are only required to modify functions in other files, by following the instructions in this assignment.

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've 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. 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. These matrices can be read directly into your program by using the loadmat command. After loading, dictionary with variable names as keys and loaded matrices as values will appear in your program's memory.

```
data = scipy.io.loadmat('ex3data1.mat')  
  
#training data stored in arrays X,y  
  
X = data['X']  
  
y = data['y']
```

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.

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

The second part of the training set is a 5000-dimensional vector y that contains labels for the training set. Pay special attention that 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 Visualizing the data

You will begin by visualizing a subset of the training set. In Part 1 of `ex3.py`, 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, and you are encouraged to examine the code to see how it works. After you run this step, you should see an image like Figure 1.

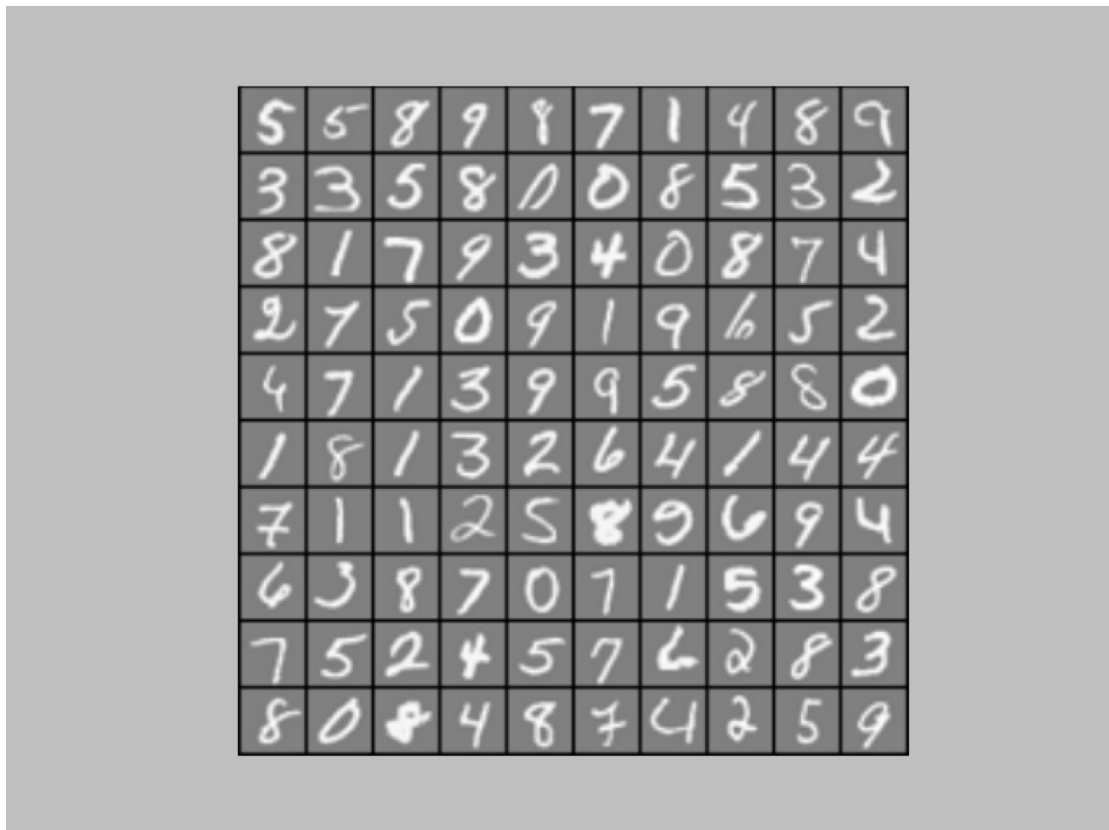


Figure 1: Examples from the dataset

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 last exercise as a starting point for this exercise.

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

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

To compute each element in the summation, we have to

compute $h_{\theta}(x^{(i)}) = g(\theta^T x^{(i)})$ for every example i , where

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

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

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.

Your job is to write the unregularized cost function in the file `lrCostFunction.py`. 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.py` should not contain any loops.

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

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

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

$$\begin{aligned}
\begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \frac{\partial J}{\partial \theta_2} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix} &= \frac{1}{m} \begin{bmatrix} \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})x_0^{(i)}) \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})x_1^{(i)}) \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})x_2^{(i)}) \\ \vdots \\ \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})x_n^{(i)}) \end{bmatrix} \\
&= \frac{1}{m} \sum_{i=1}^m ((h_\theta(x^{(i)}) - y^{(i)})x^{(i)}) \\
&= \frac{1}{m} X^T (h_\theta(x) - y). \tag{1}
\end{aligned}$$

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 $\beta_i = (h_\theta(x^{(i)}) - y^{(i)})$ 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. You should now

implement Equation 1 to compute the correct vectorized gradient. Once you are done, complete the function `lrCostFunction.py` by implementing the gradient.

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 `size` function. For example, given a data matrix X of size 100×20 (100 examples, 20 features) and θ , a vector with dimensions 20×1 , you can observe that $X\theta$ is a valid multiplication operation, while θX 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.

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

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

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

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

Now modify your code in `lrCostFunction.py` to account for regularization. Once again, you should not put any loops into your code.

Implementation Note: Since we want to use `fmin_cg` as the optimization method to calculate parameters, which takes gradient as one of the input parameters, write `lrCostFunction` and `lrGradient` to return cost and gradient separately.

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 (Figure 1). In the handwritten digits dataset, $K = 10$, but your code should work for any value of K .

You should now complete the code in `oneVsAll.py` 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 1 to K , training each classifier independently.

Note that the y argument to this function is a vector of labels from 1 to 10, where we have mapped the digit “0” to the label 10 (to avoid confusions with indexing).

When training the classifier for class $k \in \{1, \dots, K\}$, you will

want a m - 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 `fmin_cg` for this exercise (instead of `fmin`). `fmin_cg` works similarly to `fmin`, but is more efficient for dealing with a large number of parameters.

After you have correctly completed the code for `oneVsAll.py`, the script `ex3.py` will continue to use your `oneVsAll` function to train a multi-class classifier.

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 (1, 2,..., or K) as the prediction for the input example.

You should now complete the code in `predictOneVsAll.py` to use the one-vs-all classifier to make predictions.

Once you are done, `ex3.py` will call your `predictOneVsAll` function using the learned value of Θ . You should see that the training set accuracy is about 94.4% (i.e., it classifies 94.4% of the examples in the training set correctly).

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.

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.

The provided script, `ex3_nn.py`, will help you step through this exercise.

2.1 Model representation

Our neural network is shown in Figure 2. 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` and will be loaded by `ex3_nn.py` 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).

```
# Load the weights into variables Theta1 and Theta2

data = scipy.io.loadmat('ex3weights.mat')

Theta1 = data['Theta1']
```

```

Theta2 = data['Theta2']

# Theta1 has size 25*401

# Theta2 has size 10*26

```

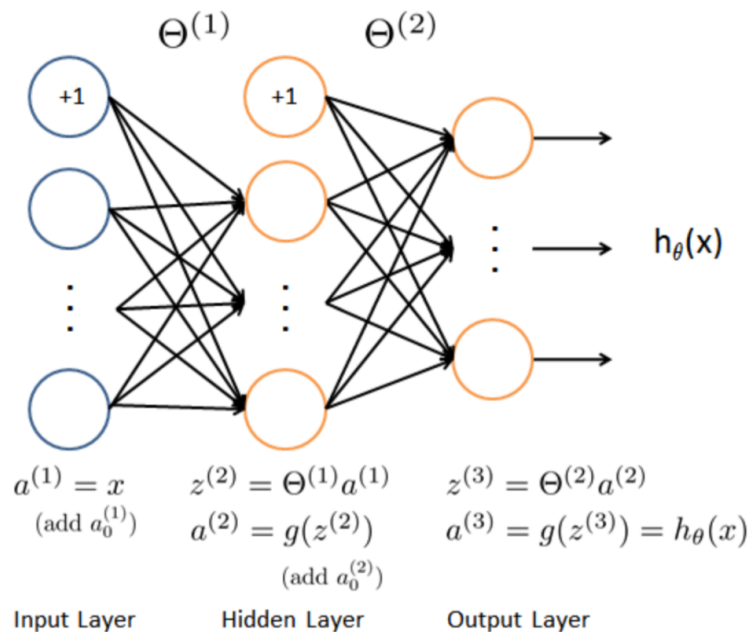


Figure 2: Neural network model.

2.2 Feedforward Propagation and Prediction

Now you will implement feedforward propagation for the neural network. You will need to complete the code in `predict.py` 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 `predict.py`, you will need to add the column of 1's to the matrix. The matrices Θ_1 and Θ_2 contain the parameters for each unit in rows. Specifically, the first row of Θ_1 corresponds to the first hidden unit in the second layer.

Once you are done, `ex3_nn.py` will call your `predict` function using the loaded set of parameters for Θ_1 and Θ_2 . You should see that the accuracy is about 97.5%. After that, an interactive sequence will launch displaying images from the training set one at a time, while the console prints out the predicted label for the displayed image. To stop the image sequence, press `Ctrl-C`.