

Programming Exercise 4: Neural Networks Learning

Introduction

In this exercise, you will implement the backpropagation algorithm for neural networks and apply it to the task of hand-written digit recognition.

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 [1]: # 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 assignments
import utils

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

Outline

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

Section	Part	Submission function
1	Feedforward and Cost Function	<code>nnCostFunction</code>
2	Regularized Cost Function	<code>nnCostFunction</code>
3	Sigmoid Gradient	<code>sigmoidGradient</code>
4	Neural Net Gradient Function (Backpropagation)	<code>nnCostFunction</code>
5	Regularized Gradient	<code>nnCostFunction</code>

Neural Networks

In the previous exercise, you implemented feedforward propagation for neural networks and used it to predict handwritten digits with the weights we provided. In this exercise, you will implement the backpropagation algorithm to learn the parameters for the neural network.

We start the exercise by first loading the dataset.

```
In [2]: # training data stored in arrays X, y
data = loadmat(os.path.join('Data', 'ex4data1.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

# Number of training examples
m = y.size
```

1.1 Visualizing the data

You will begin by visualizing a subset of the training set, using the function `displayData`, which is the same function we used in Exercise 3. It is provided in the `utils.py` file for this assignment as well. The dataset is also the same one you used in the previous exercise.

There are 5000 training examples in `ex4data1.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. The following cell randomly selects 100 images from the dataset and plots them.

```
In [3]: # Randomly select 100 data points to display
rand_indices = np.random.choice(m, 100, replace=False)
sel = X[rand_indices, :]

utils.displayData(sel)
```



1.2 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 (not counting the extra bias unit which always outputs +1). The training data was loaded into the variables `X` and `y` above.

You have been provided with a set of network parameters $(\Theta^{(1)}, \Theta^{(2)})$ already trained by us. These are stored in `ex4weights.mat` and will be loaded in the next cell of this notebook 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 [4]: # 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 weights into variables Theta1 and Theta2
weights = loadmat(os.path.join('Data', 'ex4weights.mat'))

# 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)

# Unroll parameters
nn_params = np.concatenate([Theta1.ravel(), Theta2.ravel()])

```

1.3 Feedforward and cost function

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

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

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

where $h_\theta(x^{(i)})$ is computed as shown in the neural network figure above, 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 0, 1, ..., 9, for the purpose of training a neural network, we need to encode 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).

****Implementation Note:**** The matrix X contains the examples in rows (i.e., $X[i,:]$ is the i^{th} training example $x^{(i)}$, expressed as a $n \times 1$ vector.) When you complete the code in `nnCostFunction`, you will need to add the column of 1's to the X matrix. The parameters for each unit in the neural network is represented in Θ_1 and Θ_2 as one row.

Specifically, the first row of Theta1 corresponds to the first hidden unit in the second layer. You can use a for-loop over the examples to compute the cost.

In [5]:

```
def nnCostFunction(nn_params,
                  input_layer_size,
                  hidden_layer_size,
                  num_labels,
                  X, y, lambda_=0.0):
    """
    Implements the neural network cost function and gradient for a two layer neural
    network which performs classification.

    Parameters
    -----
    nn_params : array_like
        The parameters for the neural network which are "unrolled" into
        a vector. This needs to be converted back into the weight matrices Theta1
        and Theta2.

    input_layer_size : int
        Number of features for the input layer.

    hidden_layer_size : int
        Number of hidden units in the second layer.

    num_labels : int
        Total number of labels, or equivalently number of units in output layer.

    X : array_like
        Input dataset. A matrix of shape (m x input_layer_size).

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

    lambda_ : float, optional
        Regularization parameter.

    Returns
    -----
    J : float
        The computed value for the cost function at the current weight values.

    grad : array_like
        An "unrolled" vector of the partial derivatives of the concatenation of
        neural network weights Theta1 and Theta2.

    Instructions
    -----
    You should complete the code by working through the following parts.

    - Part 1: Feedforward the neural network and return the cost in the
              variable J. After implementing Part 1, you can verify that your
              cost function computation is correct by verifying the cost
              computed in the following cell.

    - Part 2: Implement the backpropagation algorithm to compute the gradients
              Theta1_grad and Theta2_grad. You should return the partial derivative
              the cost function with respect to Theta1 and Theta2 in Theta1_grad and
              Theta2_grad, respectively. After implementing Part 2, you can check
```

that your implementation is correct by running checkNNGradients provided in the utils.py module.

Note: The vector y passed into the function is a vector of labels containing values from $0..K-1$. You need to map this vector into a binary vector of 1's and 0's to be used with the neural network cost function.

Hint: We recommend implementing backpropagation using a for-loop over the training examples if you are implementing it for the first time.

- Part 3: Implement regularization with the cost function and gradients.

Hint: You can implement this around the code for backpropagation. That is, you can compute the gradients for the regularization separately and then add them to Theta1_grad and Theta2_grad from Part 2.

Note

We have provided an implementation for the sigmoid function in the file `utils.py` accompanying this assignment.

"""

Reshape nn_params back into the parameters Theta1 and Theta2, the weight matrix for our 2 layer neural network

```
Theta1 = np.reshape(nn_params[:hidden_layer_size * (input_layer_size + 1)],  
                    (hidden_layer_size, (input_layer_size + 1)))
```

```
Theta2 = np.reshape(nn_params[(hidden_layer_size * (input_layer_size + 1)):],  
                    (num_labels, (hidden_layer_size + 1)))
```

Setup some useful variables

```
m = y.size
```

You need to return the following variables correctly

```
J = 0
```

```
Theta1_grad = np.zeros(Theta1.shape)
```

```
Theta2_grad = np.zeros(Theta2.shape)
```

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

```
sig = utils.sigmoid
```

```
yBin = np.zeros((m, num_labels))  
yBin[np.arange(m, step=1), y] = 1
```

```
ones_col = np.ones((m, 1))
```

```
X = np.concatenate([ones_col, X], axis=1)
```

```
Z2 = X@Theta1.T
```

```
A2 = np.concatenate([ones_col, sig(Z2)], axis=1)
```

```
Z3 = A2@Theta2.T
```

```
A3 = sig(Z3)
```

```
J = np.sum(-yBin*np.log(A3)-(1-yBin)*np.log(1-A3)) / m +\n        (lambda_/(2*m)) * (np.sum(Theta1[:, 1:]**2) + np.sum(Theta2[:, 1:]**2))
```

try:

```
delta3 = A3 - yBin
```

```
delta2 = delta3@Theta2 * np.concatenate([ones_col, sigmoidGradient(Z2)], axis=1)
```

```

Theta2_grad = delta3.T@A2 / m
Theta1_grad = delta2[:, 1:].T@X / m

Theta1_grad[:, 1:] += (lambda_/m) * Theta1[:, 1:]
Theta2_grad[:, 1:] += (lambda_/m) * Theta2[:, 1:]

except:
    pass

# =====
# Unroll gradients
# grad = np.concatenate([Theta1_grad.ravel(order='F'), Theta2_grad.ravel(order='F')])
grad = np.concatenate([Theta1_grad.ravel(), Theta2_grad.ravel()])

return J, grad

```

Use the following links to go back to the different parts of this exercise that require to modify the function `nnCostFunction`.

Back to:

- Feedforward and cost function
- Regularized cost
- Neural Network Gradient (Backpropagation)
- Regularized Gradient

Once you are done, call your `nnCostFunction` using the loaded set of parameters for `Theta1` and `Theta2`. You should see that the cost is about 0.287629.

```
In [6]: lambda_ = 0
J, _ = nnCostFunction(nn_params, input_layer_size, hidden_layer_size,
                      num_labels, X, y, lambda_)
print('Cost at parameters (loaded from ex4weights): %.6f ' % J)
print('The cost should be about : 0.287629.')
```

Cost at parameters (loaded from ex4weights): 0.287629
The cost should be about : 0.287629.

1.4 Regularized cost function

The cost function for neural networks with regularization is given by:

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

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. Notice that you can first compute the unregularized cost function J using your existing `nnCostFunction` and then later add the cost for the regularization terms.

[Click here to go back to `nnCostFunction` for editing.](#)

Once you are done, the next cell will call your `nnCostFunction` using the loaded set of parameters for `Theta1` and `Theta2`, and $\lambda = 1$. You should see that the cost is about 0.383770.

```
In [7]: # Weight regularization parameter (we set this to 1 here).
lambda_ = 1
J, _ = nnCostFunction(nn_params, input_layer_size, hidden_layer_size,
                      num_labels, X, y, lambda_)

print('Cost at parameters (loaded from ex4weights): %.6f' % J)
print('This value should be about : 0.383770.')
```

```
Cost at parameters (loaded from ex4weights): 0.383770
This value should be about : 0.383770.
```

2 Backpropagation

In this part of the exercise, you will implement the backpropagation algorithm to compute the gradient for the neural network cost function. You will need to update the function `nnCostFunction` so that it returns an appropriate value for `grad`. 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 `scipy`'s `optimize.minimize`. 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 exercise, 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

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

Now complete the implementation of `sigmoidGradient` in the next cell.

```
In [8]: def sigmoidGradient(z):
    """
    Computes the gradient of the sigmoid function evaluated at z.
    This should work regardless if z is a matrix or a vector.
    """
```

```
In particular, if z is a vector or matrix, you should return  
the gradient for each element.
```

```
Parameters
```

```
-----
```

```
z : array_like
```

```
    A vector or matrix as input to the sigmoid function.
```

```
Returns
```

```
-----
```

```
g : array_like
```

```
    Gradient of the sigmoid function. Has the same shape as z.
```

```
Instructions
```

```
-----
```

```
Compute the gradient of the sigmoid function evaluated at  
each value of z (z can be a matrix, vector or scalar).
```

```
Note
```

```
----
```

```
We have provided an implementation of the sigmoid function  
in `utils.py` file accompanying this assignment.
```

```
"""
```

```
g = np.zeros(z.shape)
```

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

```
sig = utils.sigmoid
```

```
g = sig(z)*(1-sig(z))
```

```
# =====
```

```
return g
```

When you are done, the following cell call `sigmoidGradient` on a given vector `z`. Try testing a few values by calling `sigmoidGradient(z)`. For large values (both positive and negative) of `z`, the gradient should be close to 0. When $z = 0$, the gradient should be exactly 0.25. Your code should also work with vectors and matrices. For a matrix, your function should perform the sigmoid gradient function on every element.

```
In [9]: z = np.array([-1, -0.5, 0, 0.5, 1])  
g = sigmoidGradient(z)  
print('Sigmoid gradient evaluated at [-1 -0.5 0 0.5 1]:\n ')  
print(g)
```

Sigmoid gradient evaluated at [-1 -0.5 0 0.5 1]:

```
[0.19661193 0.23500371 0.25 0.23500371 0.19661193]
```

2.2 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 $[-\epsilon_{init}, \epsilon_{init}]$. You should use $\epsilon_{init} = 0.12$. This range of values ensures that the parameters are kept small and makes the learning more efficient.

One effective strategy for choosing ϵ_{init} is to base it on the number of units in the network. A good choice of ϵ_{init} is $\epsilon_{init} = \frac{\sqrt{6}}{\sqrt{L_{in} + L_{out}}}$ where $L_{in} = s_l$ and $L_{out} = s_{l+1}$ are the number of units in the layers adjacent to Θ^l .

Your job is to complete the function `randInitializeWeights` to initialize the weights for Θ . Modify the function by filling in the following code:

```
# Randomly initialize the weights to small values
W = np.random.rand(L_out, 1 + L_in) * 2 * epsilon_init - epsilon_init
```

Note that we give the function an argument for ϵ with default value `epsilon_init = 0.12`.

```
In [10]: def randInitializeWeights(L_in, L_out, epsilon_init=0.12):
    """
    Randomly initialize the weights of a layer in a neural network.

    Parameters
    -----
    L_in : int
        Number of incoming connections.

    L_out : int
        Number of outgoing connections.

    epsilon_init : float, optional
        Range of values which the weight can take from a uniform
        distribution.

    Returns
    -----
    W : array_like
        The weight initialized to random values. Note that W should
        be set to a matrix of size(L_out, 1 + L_in) as
        the first column of W handles the "bias" terms.

    Instructions
    -----
    Initialize W randomly so that we break the symmetry while training
    the neural network. Note that the first column of W corresponds
    to the parameters for the bias unit.
    """

    # You need to return the following variables correctly
    W = np.zeros((L_out, 1 + L_in))

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

    W = np.random.rand(L_out, 1 + L_in) * (2 * epsilon_init) - epsilon_init

    # =====
    return W
```

Execute the following cell to initialize the weights for the 2 layers in the neural network using the `randInitializeWeights` function.

```
In [11]: print('Initializing Neural Network Parameters ...')
```

```

initial_Theta1 = randInitializeWeights(input_layer_size, hidden_layer_size)
initial_Theta2 = randInitializeWeights(hidden_layer_size, num_labels)

# Unroll parameters
initial_nn_params = np.concatenate([initial_Theta1.ravel(), initial_Theta2.ravel()])

```

Initializing Neural Network Parameters ...

2.4 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.

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 the figure above). 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 t in range(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, 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 `numpy`, if a 1 is a column matrix, adding one corresponds to `a_1 = np.concatenate([np.ones((m, 1)), a_1], axis=1)`.

2. For each output unit k in layer 3 (the output layer), set

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

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 exercise).

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

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

Note that the symbol `*` performs element wise multiplication in `numpy`.

1. Accumulate the gradient from this example using the following formula. Note that you should skip or remove $\delta_0^{(2)}$. In `numpy`, removing $\delta_0^{(2)}$ corresponds to `delta_2 = delta_2[1:]`.

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

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

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

****Python/Numpy tip**:** You should implement the backpropagation algorithm only after you have successfully completed the feedforward and cost functions. While implementing the backpropagation algorithm, it is often useful to use the `shape` function to print out the shapes of the variables you are working with if you run into dimension mismatch errors.

[Click here to go back and update the function `nnCostFunction` with the backpropagation algorithm.](#)

After you have implemented the backpropagation algorithm, we will proceed to run gradient checking on your implementation. The gradient check will allow you to increase your confidence that your code is computing the gradients correctly.

2.4 Gradient checking

In your neural network, you are minimizing the cost function $J(\Theta)$. To perform gradient checking on your parameters, you can imagine “unrolling” the parameters $\Theta^{(1)}, \Theta^{(2)}$ into a long vector θ . By doing so, you can think of the cost function being $J(\theta)$ instead and use the following gradient checking procedure.

Suppose you have a function $f_i(\theta)$ that purportedly computes $\frac{\partial}{\partial \theta_i} J(\theta)$; you'd like to check if f_i is outputting correct derivative values.

$$\text{Let } \theta^{(i+)} = \theta + \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \epsilon \\ \vdots \\ 0 \end{bmatrix} \quad \text{and} \quad \theta^{(i-)} = \theta - \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \epsilon \\ \vdots \\ 0 \end{bmatrix}$$

So, $\theta^{(i+)}$ is the same as θ , except its i^{th} element has been incremented by ϵ . Similarly, $\theta^{(i-)}$ is the corresponding vector with the i^{th} element decreased by ϵ . You can now numerically verify $f_i(\theta)$'s correctness by checking, for each i , that:

$$f_i(\theta) \approx \frac{J(\theta^{(i+)}) - J(\theta^{(i-)})}{2\epsilon}$$

The degree to which these two values should approximate each other will depend on the details of J . But assuming $\epsilon = 10^{-4}$, you'll usually find that the left- and right-hand sides of the above will agree to at least 4 significant digits (and often many more).

We have implemented the function to compute the numerical gradient for you in `computeNumericalGradient` (within the file `utils.py`). While you are not required to modify the file, we highly encourage you to take a look at the code to understand how it works.

In the next cell we will run the provided function `checkNNGradients` which will create a small neural network and dataset that will be used for checking your gradients. If your backpropagation implementation is correct, you should see a relative difference that is less than 1e-9.

****Practical Tip:**** When performing gradient checking, it is much more efficient to use a small neural network with a relatively small number of input units and hidden units, thus having a relatively small number of parameters. Each dimension of θ requires two evaluations of the cost function and this can be expensive. In the function `checkNNGradients`, our code creates a small random model and dataset which is used with `computeNumericalGradient` for gradient checking. Furthermore, after you are confident that your gradient computations are correct, you should turn off gradient checking before running your learning algorithm.

Practical Tip: Gradient checking works for any function where you are computing the cost and the gradient. Concretely, you can use the same `computeNumericalGradient` function to check if your gradient implementations for the other exercises are correct too (e.g., logistic regression's cost function).

```
In [12]: utils.checkNNGradients(nnCostFunction)
```

```

[[ -9.27825235e-03 -9.27825236e-03]
 [ -3.04978487e-06 -3.04978914e-06]
 [ -1.75060084e-04 -1.75060082e-04]
 [ -9.62660640e-05 -9.62660620e-05]
 [ 8.89911959e-03 8.89911960e-03]
 [ 1.42869427e-05 1.42869443e-05]
 [ 2.33146358e-04 2.33146357e-04]
 [ 1.17982666e-04 1.17982666e-04]
 [-8.36010761e-03 -8.36010762e-03]
 [-2.59383093e-05 -2.59383100e-05]
 [-2.87468731e-04 -2.87468729e-04]
 [-1.37149712e-04 -1.37149706e-04]
 [ 7.62813550e-03 7.62813551e-03]
 [ 3.69883235e-05 3.69883234e-05]
 [ 3.35320351e-04 3.35320347e-04]
 [ 1.53247082e-04 1.53247082e-04]
 [-6.74798370e-03 -6.74798370e-03]
 [-4.68759742e-05 -4.68759769e-05]
 [-3.76215583e-04 -3.76215587e-04]
 [-1.66560292e-04 -1.66560294e-04]
 [ 3.14544970e-01 3.14544970e-01]
 [ 1.64090819e-01 1.64090819e-01]
 [ 1.64567932e-01 1.64567932e-01]
 [ 1.58339334e-01 1.58339334e-01]
 [ 1.51127527e-01 1.51127527e-01]
 [ 1.49568335e-01 1.49568335e-01]
 [ 1.11056588e-01 1.11056588e-01]
 [ 5.75736494e-02 5.75736493e-02]
 [ 5.77867378e-02 5.77867378e-02]
 [ 5.59235296e-02 5.59235296e-02]
 [ 5.36967009e-02 5.36967009e-02]
 [ 5.31542052e-02 5.31542052e-02]
 [ 9.74006970e-02 9.74006970e-02]
 [ 5.04575855e-02 5.04575855e-02]
 [ 5.07530173e-02 5.07530173e-02]
 [ 4.91620841e-02 4.91620841e-02]
 [ 4.71456249e-02 4.71456249e-02]
 [ 4.65597186e-02 4.65597186e-02]]

```

The above two columns you get should be very similar.
 (Left-Your Numerical Gradient, Right-Analytical Gradient)

If your backpropagation implementation is correct, then
 the relative difference will be small (less than 1e-9).

Relative Difference: 2.49247e-11

2.5 Regularized Neural Network

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

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{for } j = 0 \quad (1)$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)} = \frac{1}{m} \Delta_{ij}^{(l)} + \frac{\lambda}{m} \Theta_{ij}^{(l)} \quad \text{for } j \geq 1 \quad (2)$$

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,

$$\Theta^{(l)} = \begin{bmatrix} \Theta_{1,0}^{(i)} & \Theta_{1,1}^{(l)} & \dots \\ \Theta_{2,0}^{(i)} & \Theta_{2,1}^{(l)} & \dots \\ \vdots & & \ddots \end{bmatrix}$$

Now modify your code that computes grad in `nnCostFunction` to account for regularization.

After you are done, the following cell runs gradient checking on your implementation. If your code is correct, you should expect to see a relative difference that is less than 1e-9.

```
In [13]: # Check gradients by running checkNNGradients
lambda_ = 3
utils.checkNNGradients(nnCostFunction, lambda_)

# Also output the costFunction debugging values
debug_J, _ = nnCostFunction(nn_params, input_layer_size,
                           hidden_layer_size, num_labels, X, y, lambda_)

print('\n\nCost at (fixed) debugging parameters (w/ lambda = %f): %f' % (lambda_,
print(' (for lambda = 3, this value should be about 0.576051)')
```

```

[[ -9.27825235e-03 -9.27825236e-03]
 [ -1.67679797e-02 -1.67679797e-02]
 [ -6.01744725e-02 -6.01744725e-02]
 [ -1.73704651e-02 -1.73704651e-02]
 [  8.89911959e-03  8.89911960e-03]
 [  3.94334829e-02  3.94334829e-02]
 [ -3.19612287e-02 -3.19612287e-02]
 [ -5.75658668e-02 -5.75658668e-02]
 [ -8.36010761e-03 -8.36010762e-03]
 [  5.93355565e-02  5.93355565e-02]
 [  2.49225535e-02  2.49225535e-02]
 [ -4.51963845e-02 -4.51963845e-02]
 [  7.62813550e-03  7.62813551e-03]
 [  2.47640974e-02  2.47640974e-02]
 [  5.97717617e-02  5.97717617e-02]
 [  9.14587966e-03  9.14587966e-03]
 [ -6.74798370e-03 -6.74798370e-03]
 [ -3.26881426e-02 -3.26881426e-02]
 [  3.86410548e-02  3.86410548e-02]
 [  5.46101548e-02  5.46101547e-02]
 [  3.14544970e-01  3.14544970e-01]
 [  1.18682669e-01  1.18682669e-01]
 [  2.03987128e-01  2.03987128e-01]
 [  1.25698067e-01  1.25698067e-01]
 [  1.76337550e-01  1.76337550e-01]
 [  1.32294136e-01  1.32294136e-01]
 [  1.11056588e-01  1.11056588e-01]
 [  3.81928711e-05  3.81928696e-05]
 [  1.17148233e-01  1.17148233e-01]
 [ -4.07588279e-03 -4.07588279e-03]
 [  1.13133142e-01  1.13133142e-01]
 [ -4.52964427e-03 -4.52964427e-03]
 [  9.74006970e-02  9.74006970e-02]
 [  3.36926556e-02  3.36926556e-02]
 [  7.54801264e-02  7.54801264e-02]
 [  1.69677090e-02  1.69677090e-02]
 [  8.61628953e-02  8.61628953e-02]
 [  1.50048382e-03  1.50048382e-03]]

```

The above two columns you get should be very similar.
 (Left-Your Numerical Gradient, Right-Analytical Gradient)

If your backpropagation implementation is correct, then
 the relative difference will be small (less than 1e-9).

Relative Difference: 2.38614e-11

Cost at (fixed) debugging parameters (w/ lambda = 3.000000): 0.576051
 (for lambda = 3, this value should be about 0.576051)

2.6 Learning parameters using `scipy.optimize.minimize`

After you have successfully implemented the neural network cost function and gradient computation, the next step we will use `scipy`'s minimization to learn a good set of parameters.

```
In [15]: # After you have completed the assignment, change the maxiter to a Larger
# value to see how more training helps.
options= {'maxiter': 100}
```

```

# You should also try different values of Lambda
lambda_ = 1

# Create "short hand" for the cost function to be minimized
costFunction = lambda p: nnCostFunction(p, input_layer_size,
                                         hidden_layer_size,
                                         num_labels, X, y, lambda_)

# Now, costFunction is a function that takes in only one argument
# (the neural network parameters)
res = optimize.minimize(costFunction,
                        initial_nn_params,
                        jac=True,
                        method='TNC',
                        options=options)

# get the solution of the optimization
nn_params = res.x

# Obtain Theta1 and Theta2 back from nn_params
Theta1 = np.reshape(nn_params[:hidden_layer_size * (input_layer_size + 1)],
                    (hidden_layer_size, (input_layer_size + 1)))

Theta2 = np.reshape(nn_params[(hidden_layer_size * (input_layer_size + 1)):], 
                    (num_labels, (hidden_layer_size + 1)))

```

After the training completes, we will proceed to report the training accuracy of your classifier by computing the percentage of examples it got correct. If your implementation is correct, you should see a reported training accuracy of about 95.3% (this may vary by about 1% due to the random initialization). It is possible to get higher training accuracies by training the neural network for more iterations. We encourage you to try training the neural network for more iterations (e.g., set `maxiter` to 400) and also vary the regularization parameter λ . With the right learning settings, it is possible to get the neural network to perfectly fit the training set.

```
In [16]: pred = utils.predict(Theta1, Theta2, X)
print('Training Set Accuracy: %f' % (np.mean(pred == y) * 100))
```

Training Set Accuracy: 94.20000

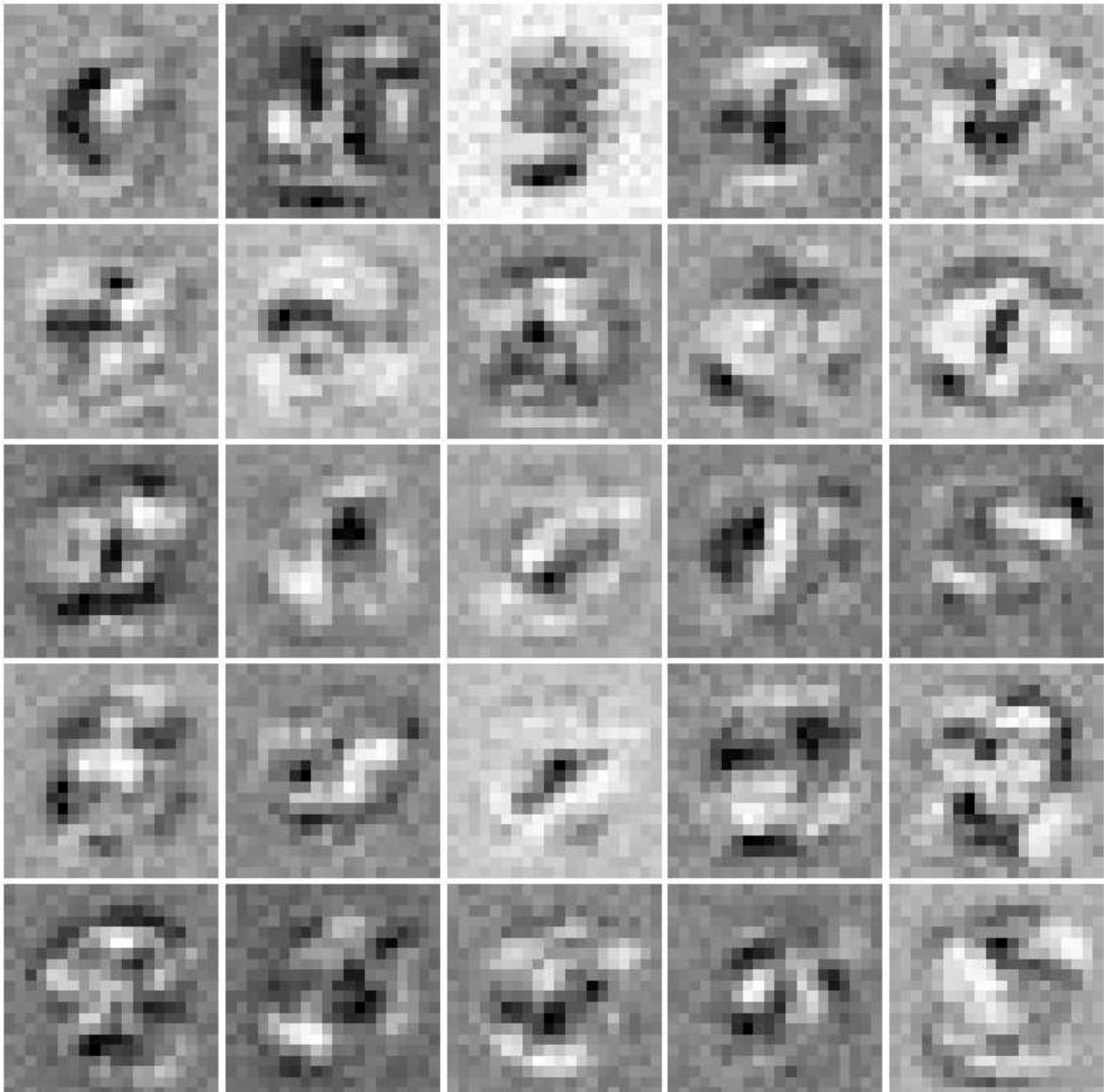
3 Visualizing the Hidden Layer

One way to understand what your neural network is learning is to visualize what the representations captured by the hidden units. Informally, given a particular hidden unit, one way to visualize what it computes is to find an input x that will cause it to activate (that is, to have an activation value ($a_i^{(l)}$) close to 1). For the neural network you trained, notice that the i^{th} row of $\Theta^{(1)}$ is a 401-dimensional vector that represents the parameter for the i^{th} hidden unit. If we discard the bias term, we get a 400 dimensional vector that represents the weights from each input pixel to the hidden unit.

Thus, one way to visualize the “representation” captured by the hidden unit is to reshape this 400 dimensional vector into a 20×20 image and display it (It turns out that this is equivalent to finding the input that gives the highest activation for the hidden unit, given a “norm” constraint on the input (i.e., $\|x\|_2 \leq 1$)).

The next cell does this by using the `displayData` function and it will show you an image with 25 units, each corresponding to one hidden unit in the network. In your trained network, you should find that the hidden units corresponds roughly to detectors that look for strokes and other patterns in the input.

```
In [17]: utils.displayData(Theta1[:, 1:])
```



3.1 Optional exercise

In this part of the exercise, you will get to try out different learning settings for the neural network to see how the performance of the neural network varies with the regularization parameter λ and number of training steps (the `maxiter` option when using `scipy.optimize.minimize`). Neural networks are very powerful models that can form highly complex decision boundaries. Without regularization, it is possible for a neural network to “overfit” a training set so that it obtains close to 100% accuracy on the training

set but does not as well on new examples that it has not seen before. You can set the regularization λ to a smaller value and the `maxiter` parameter to a higher number of iterations to see this for yourself.