

# CSE474/574 Introduction to Machine Learning

## Handwritten Digits Classification

### 1 Introduction

In this assignment, your task is to implement a Multilayer Perceptron Neural Network and evaluate its performance in classifying handwritten digits. After completing this assignment, you are able to understand:

- How Neural Network works and use Feed Forward, Back Propagation to implement Neural Network?
- How to setup a Machine Learning experiment on real data?
- How *regularization* plays a role in the *bias-variance* tradeoff?

To get started with the exercise, you will need to download the supporting files and unzip its contents to the directory you want to complete this assignment.

**Warning:** In this project, you will have to handle many computing intensive tasks such as training a neural network. Our suggestion is to use the CSE server Metallica (this server is dedicated to intensive computing tasks) to run your computation. In addition, training such a big dataset will take a very long time, maybe many hours or even days to complete. Therefore, we suggest that you should start doing this project as soon as possible so that the computer will have time to do heavy computational jobs.

#### 1.1 File included in this exercise

- *mnist\_all.mat*: original dataset from MNIST. In this file, there are 10 matrices for testing set and 10 matrices for training set, which corresponding to 10 digits. You will have to split the training data into training and validation data.
- *nnScript.py*: Python script for this programming project. Contains function definitions -
  - *preprocess()*: performs some preprocess tasks, and output the preprocessed train, validation and test data with their corresponding labels. *You need to make changes to this function.*
  - *sigmoid()*: compute sigmoid function. The input can be a scalar value, a vector or a matrix. *You need to make changes to this function.*
  - *nnObjFunction()*: compute the error function of Neural Network. *You need to make changes to this function.*
  - *nnPredict()*: predicts the label of data given the parameters of Neural Network. *You need to make changes to this function.*
  - *initializeWeights()*: return the random weights for Neural Network given the number of unit in the input layer and output layer.

## 1.2 Datasets

The MNIST dataset [1] consists of a training set of 60000 examples and test set of 10000 examples. All digits have been size-normalized and centered in a fixed image of  $28 \times 28$  size. In original dataset, each pixel in the image is represented by an integer between 0 and 255, where 0 is black, 255 is white and anything between represents different shade of gray.

In many research papers, the official training set of 60000 examples is divided into an actual training set of 50000 examples and validation set of 10000 examples. So we suggest that you should follow this convention.

## 2 Your tasks

- Implement **Neural Network** (forward pass and back propagation)
- Incorporate regularization on the weights ( $\lambda$ )
- Use validation set to tune hyper-parameters for Neural Network (number of units in the hidden layer and  $\lambda$ ).
- Write a report to explain the experimental results.

## 3 Some practical tips in implementation

### 3.1 Preprocessing

The preprocessing needs to do three things:

1. Stack all training matrices into one  $60000 \times 784$  matrix. Do the same for test matrices.
2. Create a 60000 length vector with true labels (digits) for each training example. Same for test data.
3. Normalize the training matrix and test matrix so that the values are between 0 and 1.
4. Randomly split the  $60000 \times 784$  normalized matrix into two matrices: training matrix ( $50000 \times 784$ ) and validation matrix ( $10000 \times 784$ ). Make sure you split the true labels vector into two parts as well.
5. Feature selection - see next subsection.

### 3.2 Feature selection

In the dataset, one can observe that there are many features which values are exactly the same for all data points in the training set. With those features, the classification models cannot gain any more information about the difference (or variation) between data points. Therefore, we can ignore those features in the pre-processing step.

Later on in this course, you will learn more sophisticated models to reduce the dimension of dataset (but not for this assignment).

### 3.3 Neural Network

#### 3.3.1 Neural Network Representation

Neural network can be graphically represented as in Figure 1.

As observed in the Figure 1, there are totally 3 layers in the neural network:

- The first layer comprises of  $(d + 1)$  units, each represents a feature of image (there is one extra unit representing the bias).

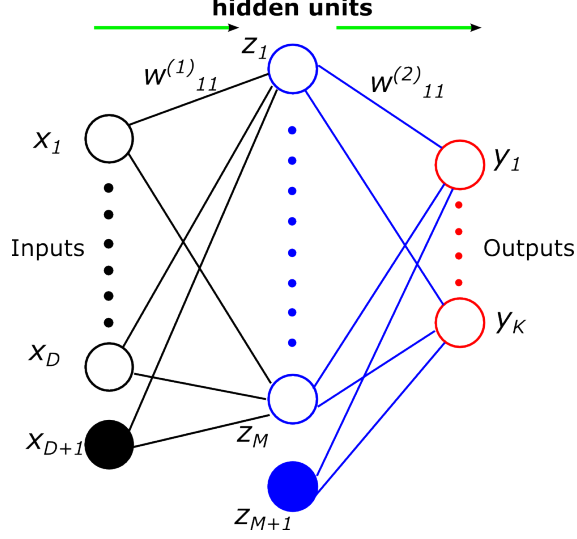


Figure 1: Neural network

- The second layer in neural network is called the hidden units. In this document, we denote  $m + 1$  as the number of hidden units in hidden layer. There is an additional bias node at the hidden layer as well. Hidden units can be considered as the learned features extracted from the original data set. Since number of hidden units will represent the dimension of learned features in neural network, it's our choice to choose an appropriate number of hidden units. Too many hidden units may lead to the slow training phase while too few hidden units may cause the under-fitting problem.
- The third layer is also called the output layer. The value of  $l^{th}$  unit in the output layer represents the probability of a certain hand-written image belongs to digit  $l$ . Since we have 10 possible digits, there are 10 units in the output layer. In this document, we denote  $k$  as the number of output units in output layer.

The parameters in Neural Network model are the weights associated with the hidden layer units and the output layers units. In our standard Neural Network with 3 layers (input, hidden, output), in order to represent the model parameters, we use 2 matrices:

- $W^{(1)} \in \mathbb{R}^{m \times (d+1)}$  is the weight matrix of connections from input layer to hidden layer. Each row in this matrix corresponds to the weight vector at each hidden layer unit.
- $W^{(2)} \in \mathbb{R}^{k \times (m+1)}$  is the weight matrix of connections from hidden layer to output layer. Each row in this matrix corresponds to the weight vector at each output layer unit.

We also further assume that there are  $n$  training samples when performing learning task of Neural Network. In the next section, we will explain how to perform learning in Neural Network.

### 3.3.2 Feedforward Propagation

In Feedforward Propagation, given parameters of Neural Network and a feature vector  $\mathbf{x}$ , we want to compute the probability that this feature vector belongs to a particular digit.

Suppose that we have totally  $m$  hidden units. Let  $a_j$  for  $1 \leq j \leq m$  be the linear combination of input data and let  $z_j$  be the output from the hidden unit  $j$  after applying an activation function (in this exercise, we use sigmoid as an activation function). For each hidden unit  $j$  ( $j = 1, 2, \dots, m$ ), we can compute its value as follow:

$$a_j = \sum_{i=1}^{D+1} w_{ji}^{(1)} x_i \quad (1)$$

$$z_j = \sigma(a_j) = \frac{1}{1 + \exp(-a_j)} \quad (2)$$

where  $w_{ji}^{(1)} = W^{(1)}[j][i]$  is the weight of connection from unit  $i$  in input layer to unit  $j$  in hidden layer. Note that we do not compute the output for the bias hidden node ( $m+1$ );  $z_{m+1}$  is directly set to 1.

The third layer in neural network is called the output layer where the learned features in hidden units are linearly combined and a sigmoid function is applied to produce the output. Since in this assignment, we want to classify a hand-written digit image to its corresponding class, we can use the one-vs-all binary classification in which each output unit  $l$  ( $l = 1, 2, \dots, 10$ ) in neural network represents the probability of an image belongs to a particular digit. For this reason, the total number of output unit is  $k = 10$ . Concretely, for each output unit  $l$  ( $l = 1, 2, \dots, 10$ ), we can compute its value as follow:

$$b_l = \sum_{j=1}^{m+1} w_{lj}^{(2)} z_j \quad (3)$$

$$o_l = \sigma(b_l) = \frac{1}{1 + \exp(-b_l)} \quad (4)$$

Now we have finished the **Feedforward pass**.

### 3.3.3 Error function and Backpropagation

The error function in this case is the negative log-likelihood error function which can be written as follow:

$$J(W^{(1)}, W^{(2)}) = -\frac{1}{n} \sum_{i=1}^n \sum_{l=1}^k (y_{il} \ln o_{il} + (1 - y_{il}) \ln(1 - o_{il})) \quad (5)$$

where  $y_{il}$  indicates the  $l^{th}$  target value in 1-of-K coding scheme of input data  $i$  and  $o_{il}$  is the output at  $l^{th}$  output node for the  $i^{th}$  data example (See (4)).

Because of the form of error function in equation (5), we can separate its error function in terms of error for each input data  $\mathbf{x}_i$ :

$$J(W^{(1)}, W^{(2)}) = \frac{1}{n} \sum_{i=1}^n J_i(W^{(1)}, W^{(2)}) \quad (6)$$

where

$$J_i(W^{(1)}, W^{(2)}) = -\sum_{l=1}^k (y_{il} \ln o_{il} + (1 - y_{il}) \ln(1 - o_{il})) \quad (7)$$

One way to learn the model parameters in neural networks is to initialize the weights to some random numbers and compute the output value (feed-forward), then compute the error in prediction, transmits this error backward and update the weights accordingly (error backpropagation).

The feed-forward step can be computed directly using formula (1), (2), (3) and (4).

On the other hand, the error backpropagation step requires computing the derivative of error function with respect to the weight.

Consider the derivative of error function with respect to the weight from the hidden unit  $j$  to output unit  $l$  where  $j = 1, 2, \dots, m+1$  and  $l = 1, \dots, 10$ :

$$\frac{\partial J_i}{\partial w_{lj}^{(2)}} = \frac{\partial J_i}{\partial o_l} \frac{\partial o_l}{\partial b_l} \frac{\partial b_l}{\partial w_{lj}^{(2)}} \quad (8)$$

$$= \delta_l z_j \quad (9)$$

where

$$\delta_l = \frac{\partial J_i}{\partial o_l} \frac{\partial o_l}{\partial b_l} = -\left(\frac{y_l}{o_l} - \frac{1 - y_l}{1 - o_l}\right)(1 - o_l)o_l = o_l - y_l$$

Note that we are dropping the subscript  $i$  for simplicity. The error function (log loss) that we are using in (5) is different from the squared loss error function that we have discussed in class. Note that the choice of the error function has “simplified” the expressions for the error!

On the other hand, the derivative of error function with respect to the weight from the input unit  $l$  to output unit  $j$  where  $p = 1, 2, \dots, d + 1$  and  $j = 1, \dots, m$  can be computed as follow:

$$\frac{\partial J_i}{\partial w_{jp}^{(1)}} = \sum_{l=1}^k \frac{\partial J_i}{\partial o_l} \frac{\partial o_l}{\partial b_l} \frac{\partial b_l}{\partial z_j} \frac{\partial z_j}{\partial a_j} \frac{\partial a_j}{\partial w_{jp}^{(1)}} \quad (10)$$

$$= \sum_{l=1}^k \delta_l w_{lj}^{(2)} (1 - z_j) z_j x_p \quad (11)$$

$$= (1 - z_j) z_j \left( \sum_{l=1}^k \delta_l w_{lj}^{(2)} \right) x_p \quad (12)$$

Note that we do not compute the gradient for the weights at the bias hidden node.

After finish computing the derivative of error function with respect to weight of each connection in neural network, we now can write the formula for the gradient of error function:

$$\nabla J(W^{(1)}, W^{(2)}) = \frac{1}{n} \sum_{i=1}^n \nabla J_i(W^{(1)}, W^{(2)}) \quad (13)$$

We again can use the gradient descent to update each weight (denoted in general as  $w$ ) with the following rule:

$$w^{new} = w^{old} - \gamma \nabla J(w^{old}) \quad (14)$$

### 3.3.4 Regularization in Neural Network

In order to avoid overfitting problem (the learning model is best fit with the training data but give poor generalization when test with validation data), we can add a regularization term into our error function to control the magnitude of parameters in Neural Network. Therefore, our objective function can be rewritten as follow:

$$\tilde{J}(W^{(1)}, W^{(2)}) = J(W^{(1)}, W^{(2)}) + \frac{\lambda}{2n} \left( \sum_{j=1}^m \sum_{p=1}^{d+1} (w_{jp}^{(1)})^2 + \sum_{l=1}^k \sum_{j=1}^{m+1} (w_{lj}^{(2)})^2 \right) \quad (15)$$

where  $\lambda$  is the regularization coefficient.

With this new objective function, the partial derivative of new objective function with respect to weight from hidden layer to output layer can be calculated as follow:

$$\frac{\partial \tilde{J}}{\partial w_{lj}^{(2)}} = \frac{1}{n} \left( \sum_{i=1}^n \frac{\partial J_i}{\partial w_{lj}^{(2)}} + \lambda w_{lj}^{(2)} \right) \quad (16)$$

Similarly, the partial derivative of new objective function with respect to weight from input layer to hidden layer can be calculated as follow:

$$\frac{\partial \tilde{J}}{\partial w_{jp}^{(1)}} = \frac{1}{n} \left( \sum_{i=1}^n \frac{\partial J_i}{\partial w_{jp}^{(1)}} + \lambda w_{jp}^{(1)} \right) \quad (17)$$

With this new formulas for computing objective function (15) and its partial derivative with respect to weights (16) (17), we can again use gradient descent to find the minimum of objective function.

### 3.3.5 Python implementation of Neural Network

In the supporting files, we have provided the base code for you to complete. In particular, you have to complete the following function in Python:

- *sigmoid*: compute sigmoid function. The input can be a scalar value, a vector or a matrix.
- *nnObjFunction*: compute the objective function of Neural Network *with regularization* and the gradient of objective function.
- *nnPredict*: predicts the label of data given the parameters of Neural Network.

Details of how to implement the required functions is explained in Python code.

**Optimization:** In general, the learning phase of Neural Network consists of 2 tasks. First task is to compute the value and gradient of error function given Neural Network parameters. Second task is to optimize the error function given the value and gradient of that error function. As explained earlier, we can use gradient descent to perform the optimization problem. In this assignment, you have to use the Python scipy function: **scipy.optimize.minimize** (using the option *method='CG'* for conjugate gradient descent), which performs the conjugate gradient descent algorithm to perform optimization task. In principle, conjugate gradient descent is similar to gradient descent but it chooses a more sophisticated learning rate  $\gamma$  in each iteration so that it will converge faster than gradient descent. Details of how to use *minimize* are provided here: <http://docs.scipy.org/doc/scipy-0.14.0/reference/generated/scipy.optimize.minimize.html>.

We use regularization in Neural Network to avoid overfitting problem (more about this will be discussed in class). You are encouraged to change different value of  $\lambda$  to see its effect in prediction accuracy in validation set. Your report should include diagrams to explain the relation between  $\lambda$  and performance of Neural Network. Moreover, by plotting the value of  $\lambda$  with respect to the accuracy of Neural Network, you should explain in your report how to choose an appropriate hyper-parameter  $\lambda$  to avoid both underfitting and overfitting problem. You can vary  $\lambda$  from 0 (no regularization) to 1 in increments of 0.1 or 0.2.

You are also encouraged to try different number hidden units to see its effect to the performance of Neural Network. Since training Neural Network is very slow, especially when the number hidden units in Neural Network is large. You should try with small hidden units and gradually increase the size and see how it effects the training time. Your report should include some diagrams to explain relation between number of hidden units and training time. Recommended values: 4, 8, 12, 16, 20.

## 4 Submission

You are required to submit a single file called *proj1.zip* using UBLearn.

File *proj1.zip* must contain 2 folders: *report* and *code*.

- Folder *report* contains your report file (in pdf format). Please indicate the team members and your course number on the top of the report.
- Folder *code* must contains the following updated files: *nnScript.py* and *params.pickle*<sup>1</sup>. File *params.pickle* contains the learned parameters of Neural Network. Concretely, file *params.pickle* must contain the following variables: optimal *n\_hidden* (number of units in hidden layer), *w1* (matrix of weight  $W^{(1)}$  as mentioned in section 3.2.1), *w2* (matrix of weight  $W^{(2)}$  as mentioned in section 3.2.1), optimal  $\lambda$  (regularization coefficient  $\lambda$  as mentioned in section 3.2.4).<sup>2</sup>

<sup>1</sup>Check this to learn how to pickle objects in Python: <https://wiki.python.org/moin/UsingPickle>

<sup>2</sup>If you want to write more supporting functions to complete the required functions, you should include these supporting functions and a README file which explains your supporting functions.

**Using UBLearns Submission:** In the groups page of the UBLearns website you will see groups called “Programming Assignment 1 Group x”. Please choose any available group number for your group and join the group. All project group members must join the same group. Please do not join any other group on UBLearns that you are not part of. You should submit one solution per group through the groups page.

**Project report:** The hard-copy of report will be collected in class at due date. Your report should include the following:

- Explanation of how to choose the hyper-parameters for Neural Network (number of hidden units, regularization term  $\lambda$ ).

## 5 Grading scheme

- Successfully implement Neural Network: 60 points (*preprocess()* [10 points], *sigmoid()* [10 points], *nnObjFunction()* [30 points], *nnPredict()* [10 points]).
- Project report: 30 points
  - Explanation with supporting figures of how to choose the hyper-parameter for Neural Network: 30 points
- Accuracy of classification method on the test data: 10 points

## References

- [1] LeCun, Yann; Corinna Cortes, Christopher J.C. Burges. “MNIST handwritten digit database”.
- [2] Bishop, Christopher M. “Pattern recognition and machine learning (information science and statistics).” (2007).