

**NANYANG  
TECHNOLOGICAL  
UNIVERSITY**

---

**SINGAPORE**

**CZ042 NEURAL NETWORKS  
ASSIGNMENT 1**

**ZHANG WANLU U1420638G  
ZHU YIMIN**

**SCHOOL OF COMPUTER ENGINEERING**

# 1 Part 1

## 1.1 Introduction

The first part of this project is to use neural network to deal with classification problem. The aim is to predict class labels in the test dataset after training the neural networks on the training data. The dataset contains multispectral values of pixels in a 3x3 neighbourhoods in a satellite images and class labels of the centre pixels in each neighbourhood. Since different parameters will affect the overall performance, it's necessary to try out several of them and select the one given the best result. In this section, we use multilayer feedforward networks to solve the problem.

## 1.2 Notes

For Part 1, the main tools we use are Theano, numpy and matplotlib. All the training and performance measurements are done on a notebook with 4.0 GHz Intel Core i7-6700k.

## 1.3 Data Preparation

The training data is read from file *sat\_train.txt* while the test data is read from file *sat\_test.txt*. Then we perform scaling on train and test data to make them distributed within a certain range. The min-max value method is used here. The code for data preparation is here.

```
1 # scale data
2 def scale(X, X_min, X_max):
3     return (X - X_min)/(X_max-X_min)
4
5 #read train data
6 train_input = np.loadtxt('sat_train.txt', delimiter=' ')
7 trainX, train_Y = train_input[:, :36], train_input[:, -1].astype(int)
8 trainX_min, trainX_max = np.min(trainX, axis=0), np.max(trainX, axis=0)
9 trainX = scale(trainX, trainX_min, trainX_max)
10
11 train_Y[train_Y == 7] = 6
12 trainY = np.zeros((train_Y.shape[0], 6))
13 trainY[np.arange(train_Y.shape[0]), train_Y-1] = 1
14
15 #read test data
16 test_input = np.loadtxt('sat_test.txt', delimiter=' ')
17 testX, test_Y = test_input[:, :36], test_input[:, -1].astype(int)
18
19 testX_min, testX_max = np.min(testX, axis=0), np.max(testX, axis=0)
20 testX = scale(testX, testX_min, testX_max)
21
22 test_Y[test_Y == 7] = 6
23 testY = np.zeros((test_Y.shape[0], 6))
24 testY[np.arange(test_Y.shape[0]), test_Y-1] = 1
```

## 1.4 Question 1

Question 1 is to design and implement a 3-layer feedforward neural network with a hidden layer of 10 neurons. The learning curves and testing results are in the result section.

### 1.4.1 Method

The model contains 3 layers, input layer, hidden layer and output layer. The activation function for the hidden layer is logistic activation. The output layer uses softmax function. The weight and bias are random initialized. We set the default batch size as 32. We use Stochastic gradient descent (shortened to SGD) with mini batch to train the network. The cost function is defined as

the mean of the crossentropy of predicted and actual value then plus regularization term. Once the predefined Theano function *train* is called, the cost will be computed then the weights and biases will be updated by the SGD function. The learning rate  $\alpha = 0.01$ . The decay parameter  $\beta = 10^{-6}$ .

```

1  # theano expressions
2  X = T.matrix() #features
3  Y = T.matrix() #output
4
5  #weights and biases from input to hidden layer
6  w1, b1 = init_weights(36, hnno), init_bias(hnno)
7  #weights and biases from hidden to output layer
8  w2, b2 = init_weights(hnno, 6, logistic=False), init_bias(6)
9
10 h1 = T.nnet.sigmoid(T.dot(X, w1) + b1)
11 py = T.nnet.softmax(T.dot(h1, w2) + b2)
12
13 y_x = T.argmax(py, axis=1)
14
15 cost = T.mean(T.nnet.categorical_crossentropy(py, Y)) + decay*(T.sum(T.sqr(w1))+T.
16     sum(T.sqr(w2)))
17 params = [w1, b1, w2, b2]
18 updates = sgd(cost, params, learning_rate)
19
20 # compile
21 train = theano.function(inputs=[X, Y],
22     outputs=cost,
23     updates=updates,
24     allow_input_downcast=True)
25 predict = theano.function(inputs=[X],
26     outputs=y_x,
27     allow_input_downcast=True)

```

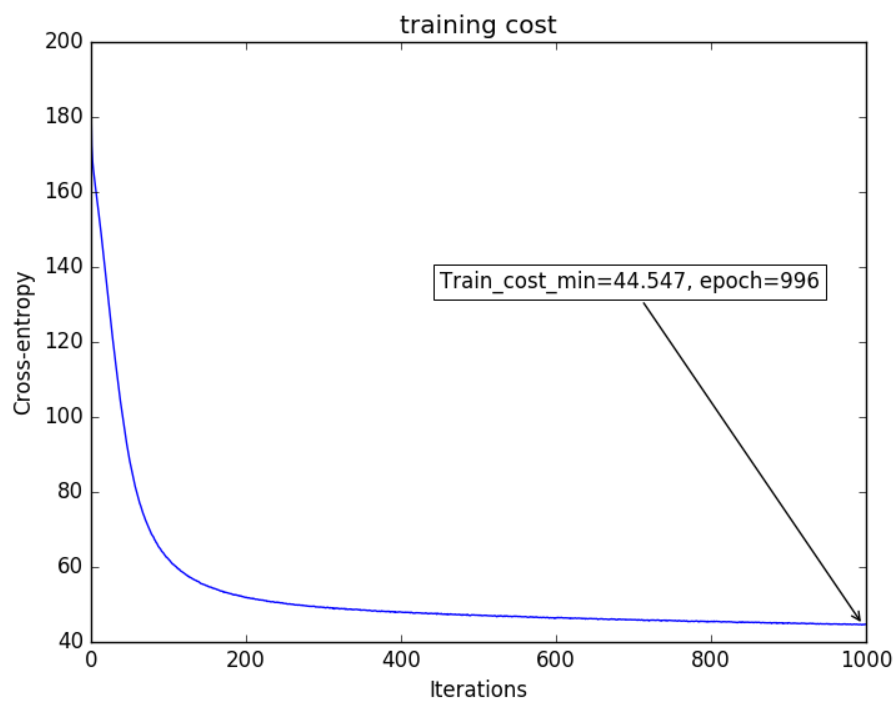
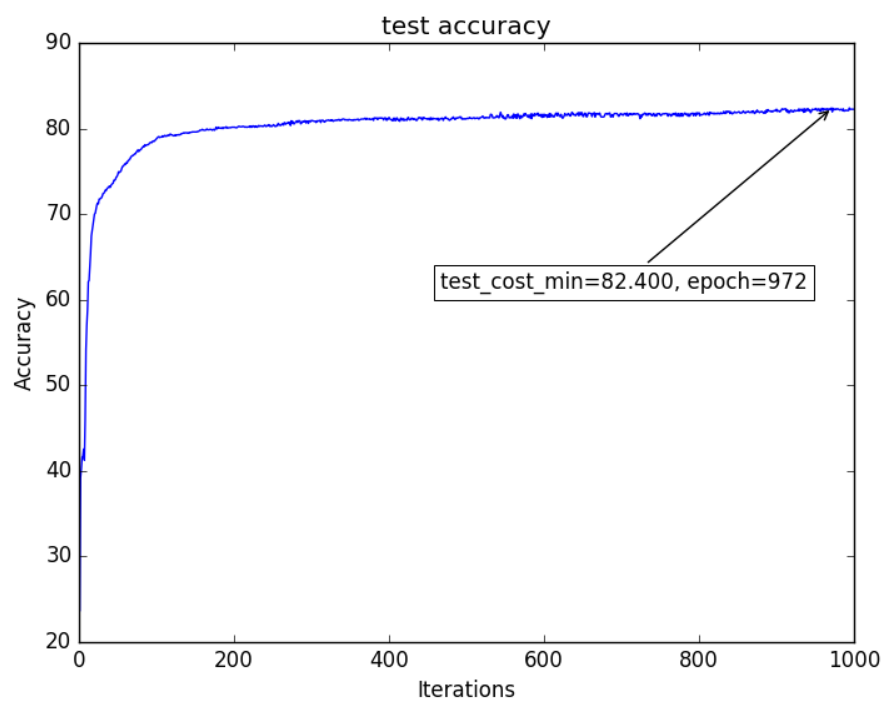
At each epoch, the data will be shuffled. An average cost for this epoch is computed and stored in the list *train\_cost*.

```

1  # train and test
2  n = len(trainX)
3  test_accuracy = []
4  train_cost = []
5  for i in range(epochs):
6      if i % 1000 == 0:
7          print(i)
8
9      trainX, trainY = shuffle_data(trainX, trainY)
10     cost = 0.0
11     for start, end in zip(range(0, n, batch_size), range(batch_size, n, batch_size)):
12         cost += train(trainX[start:end], trainY[start:end])
13     train_cost = np.append(train_cost, cost/(n // batch_size))
14
15     test_accuracy = np.append(test_accuracy, np.mean(np.argmax(testY, axis=1) ==
16         predict(testX)))
17 print('%1f accuracy at %d iterations'%(np.max(test_accuracy)*100, np.argmax(
18     test_accuracy)+1))

```

### 1.4.2 Result



## 1.5 Question 2

Question 2 is to find the optimal batch size for mini-batch gradient descent.

### 1.5.1 Method

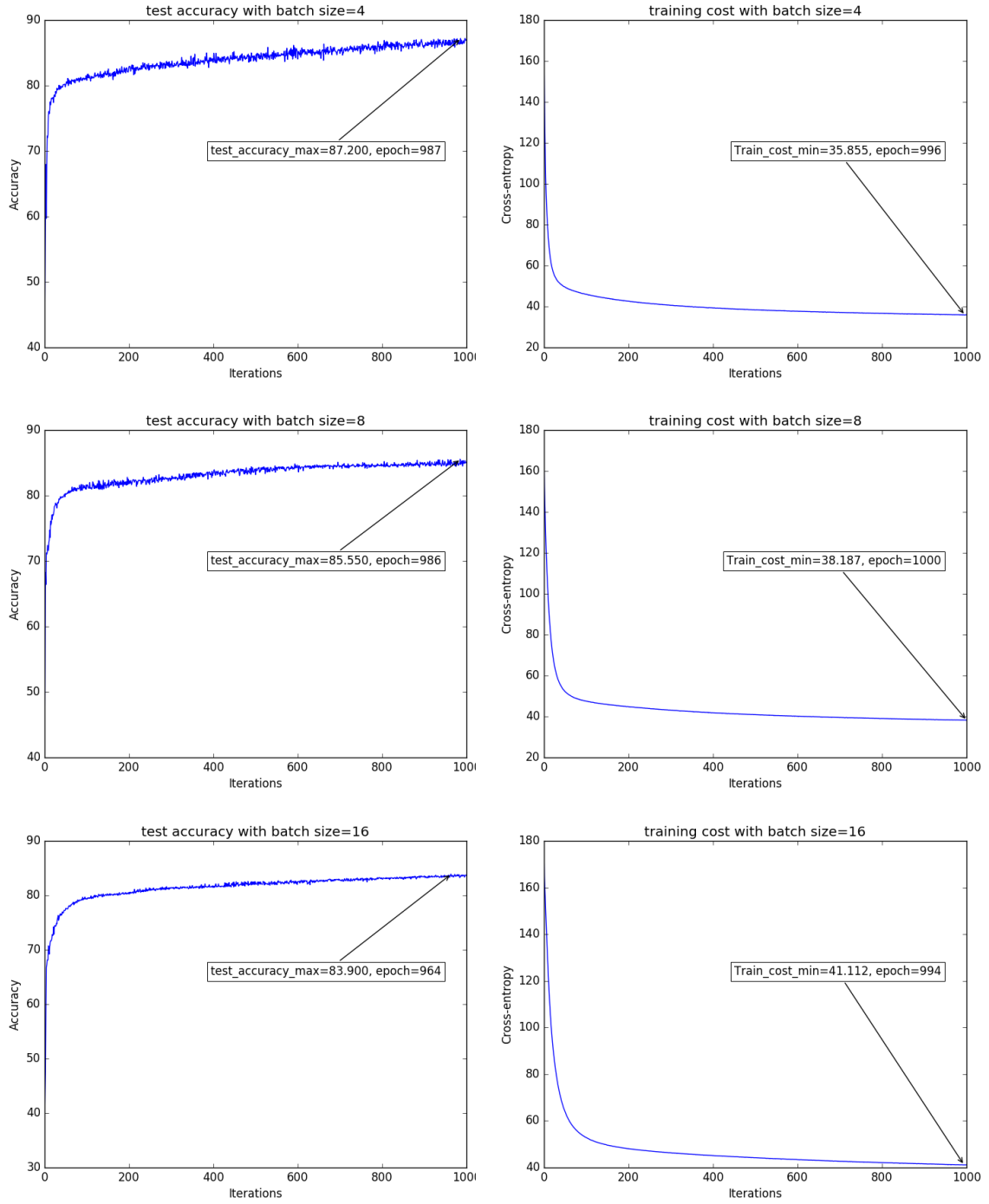
We tried different batch size from the list `batch_sizeA` {4,8,16,32,64}. A for loop is used to iterate through the list and execute the code. In each iteration, a new model is generated and initialized for training. The list `times` is used to record the time taken to update parameters of the network for each batch size.

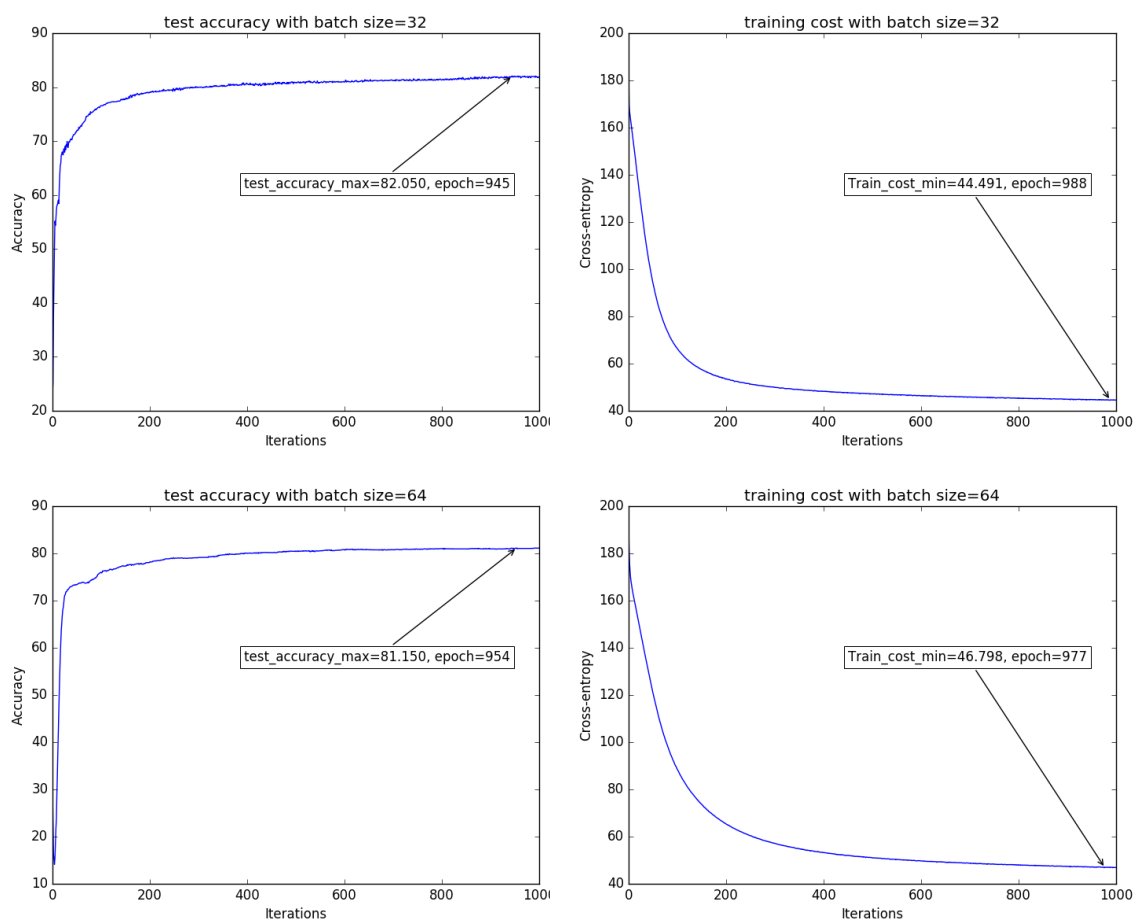
```
1 def tnt(batch_size, hnno, decay, trainX, trainY):
2     initialize(hnno, decay)
3     global test_accuracy, train_cost, epochs, times, Accuracy_plot
4     t1 = time.time()
5     n = len(trainX)
6     for i in range(epochs):
7         trainX, trainY = shuffle_data(trainX, trainY)
8         cost = 0.0
9         for start, end in zip(range(0, n, batch_size), range(batch_size, n, batch_size)):
10             cost += train(trainX[start:end], trainY[start:end])
11             train_cost = np.append(train_cost, cost/(n // batch_size))
12             test_accuracy = np.append(test_accuracy, np.mean(np.argmax(testY, axis=1) ==
13                 predict(testX)))
14             times.append(time.time()-t1)
15             print('Time: %.2fs'%times[-1])
16             accuracy = (np.max(test_accuracy)*100, np.argmax(test_accuracy)+1)
17             print('%.2f accuracy at %d iterations'%accuracy)
18
19 times = []
20 batch_sizeA = [4,8,16,32,64]
21
22 for batch_size in batch_sizeA:
23     hnno = 10
24     decay = 1e-6
25     test_accuracy = []
26     train_cost = []
27     print('Batch size = %d running'%batch_size)
28     tnt(batch_size, hnno, decay, trainX, trainY)
29     plots(batch_size)
30 timeplot(batch_sizeA)
```

The function `tnt` is also used in the program for following questions.

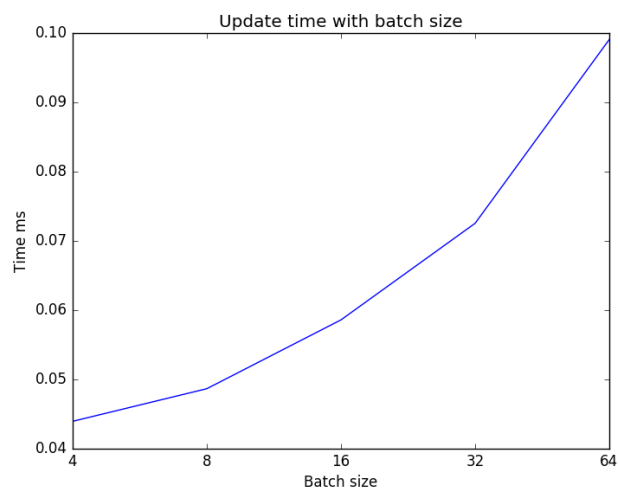
### 1.5.2 Result

The figures for training errors and test accuracies against the number of epochs with different batch size.





As shown above, we can see that the performance is best when **batch size = 4**. Batch size determines how many examples are used to update the weight. The lower it is, the noisier the training signal is going to be, the higher it is, the longer it will take to compute the gradient for each step. Furthermore, too big batch size may result in local optimum. The following figure shows the time taken to update parameters.



## 1.6 Question 3

Question 3 is to find the optimal number of hidden neurons. According to last part, the batch size is set as 4.

### 1.6.1 Method

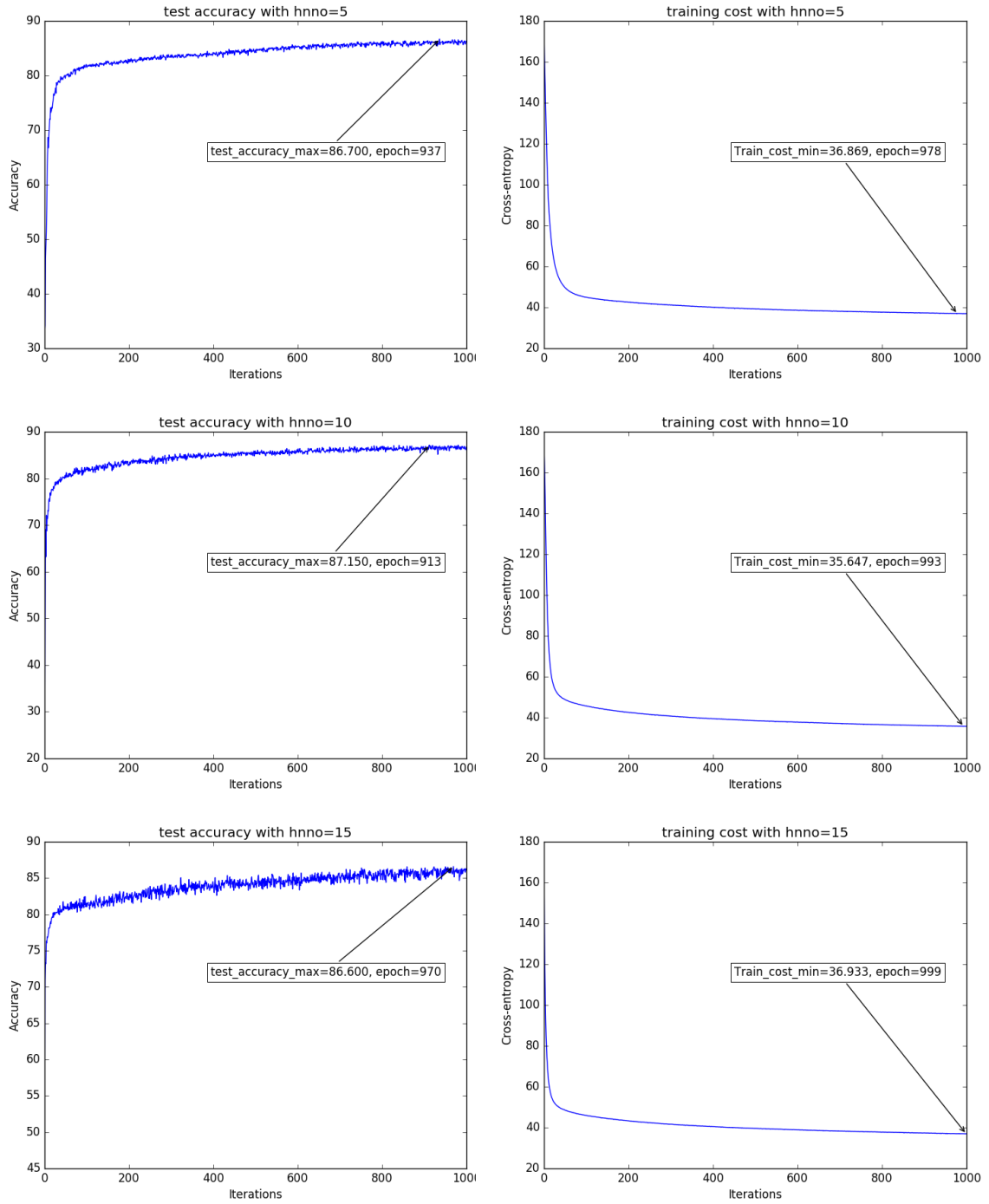
Method is similar as the batch size part. The number of hidden neurons are from the list hnnoA.

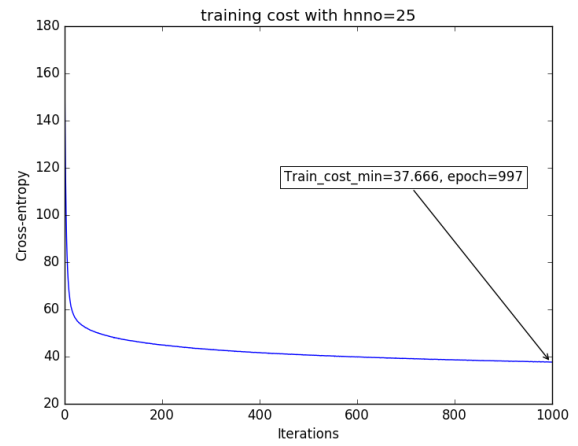
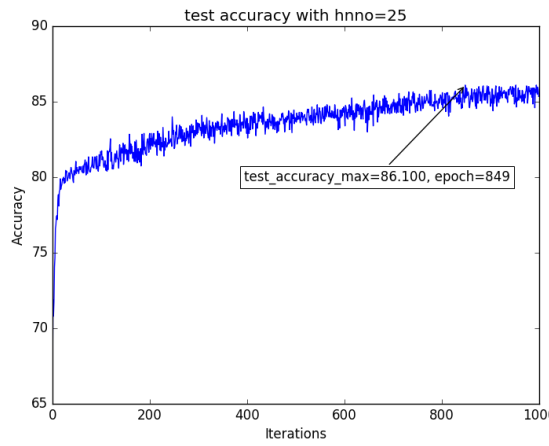
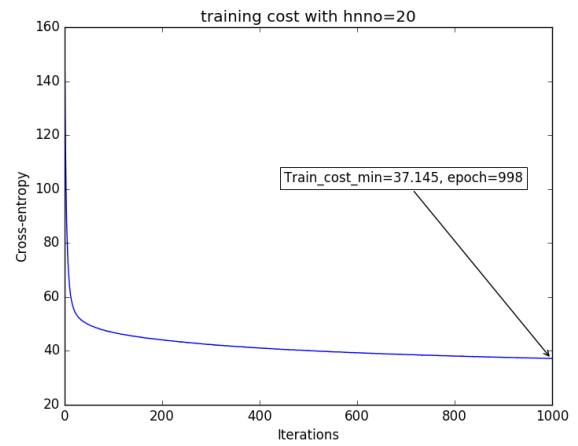
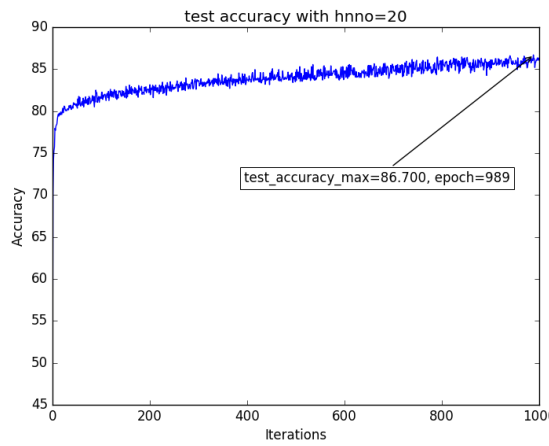
```
1 hnnoA = [5,10,15,20,25]
2 times = []
3
4 for hnno in hnnoA:
5     batch_size = 4
6     decay = 1e-6
7     test_accuracy = []
8     train_cost = []
9     print('No. of hidden neurons = %d running'%hnno)
10    tnt(batch_size, hnno, decay, trainX, trainY)
11    plots(hnno)
12 timeplot(hnnoA)
```



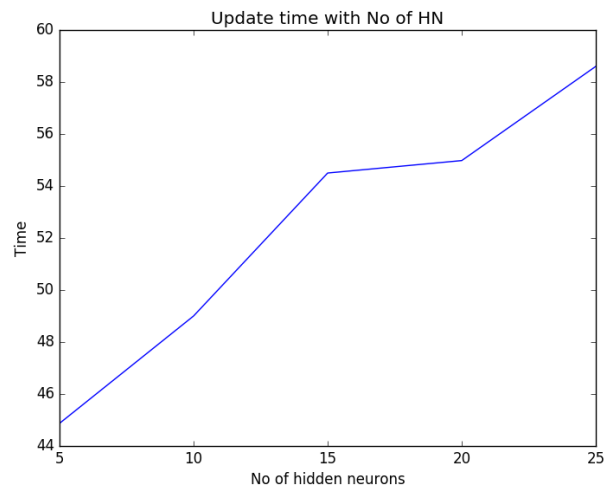
### 1.6.2 Result

The following figures show the training errors and test accuracies against the number of epochs with different number of hidden-layer neurons.





From the figures above, we can see that when  $n = 10$  the performance is best. The number of neurons may lead to underfitting while too many neurons may result in overfitting. The following figure shows the time to update parameters.



## 1.7 Question 4

Question 4 is to find the optimal decay parameter. According to last 2 parts, the batch size is set as 4 and number of neurons is 10.

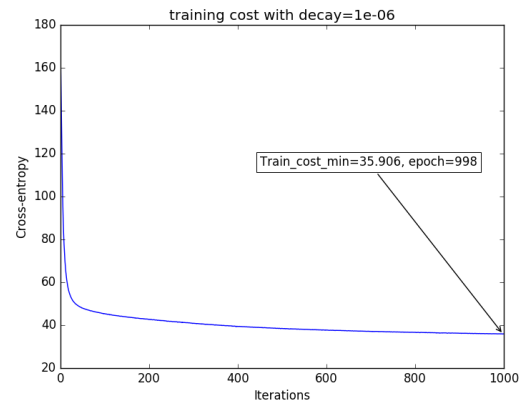
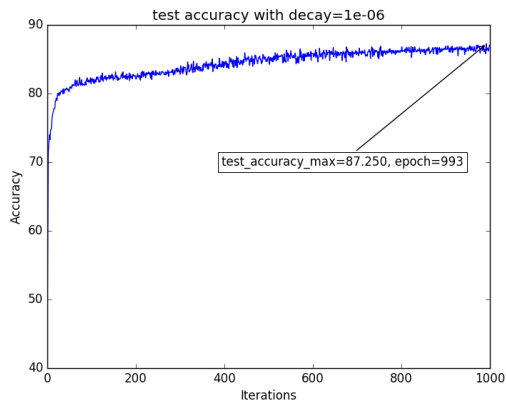
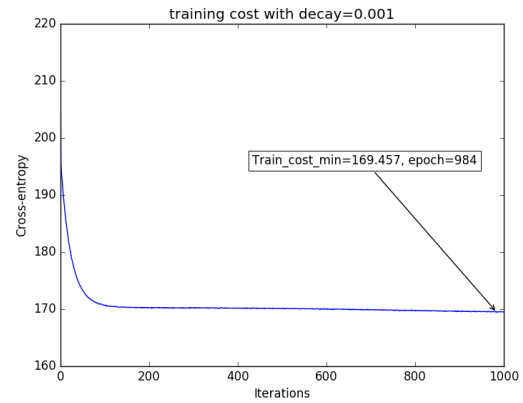
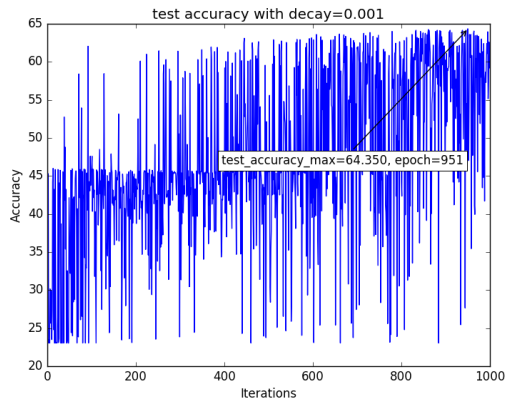
### 1.7.1 Method

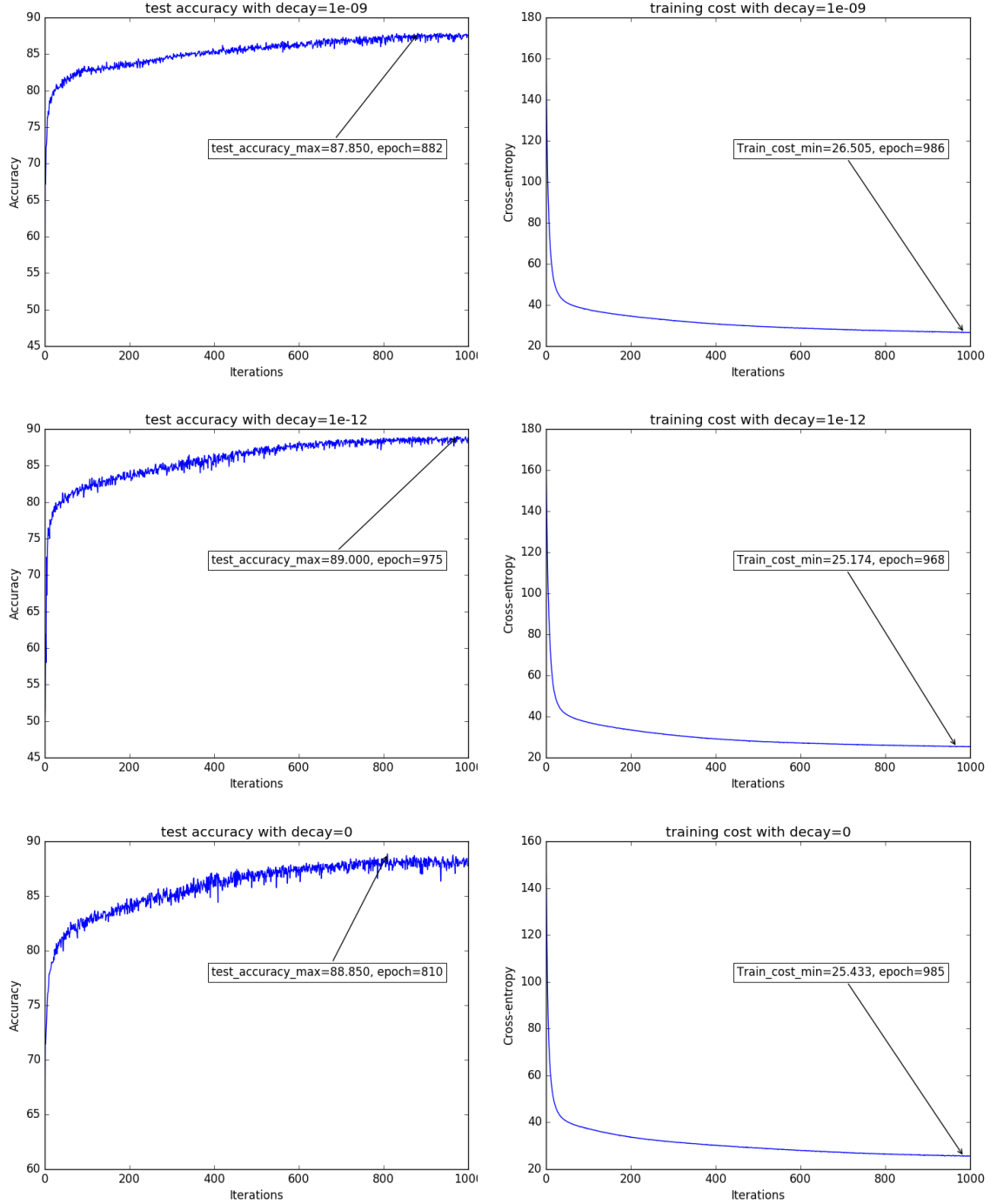
Method is similar as the last part. The decay parameters are from the list decayA.

```
1 decayA = [1e-3,1e-6,1e-9,1e-12,0]
2 times = []
3 Accuracy_plot = []
4
5 for decay in decayA:
6     batch_size = 4
7     hnno = 10
8     test_accuracy = []
9     train_cost = []
10    print('Decay parameter = '+str(decay)+' running')
11    tnt(batch_size, hnno, decay, trainX, trainY)
12    plots(decay)
13 decayplot(decayA)
```

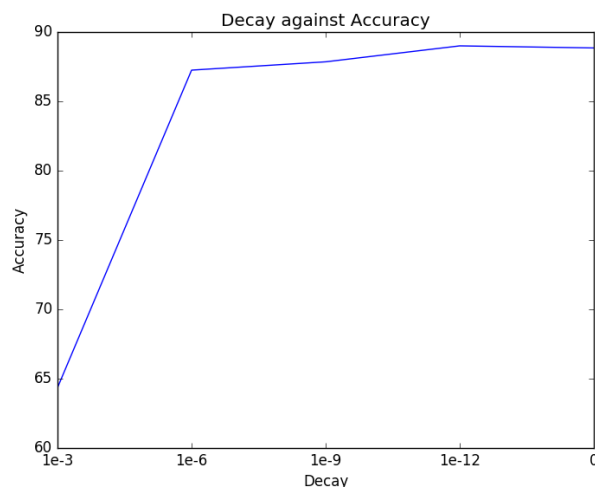
### 1.7.2 Result

The following figures show the training errors and test accuracies against the number of epochs with different number of hidden-layer neurons.





Weight decay is a regularization term that penalizes big weights. From the above figures, we can see that the result is extremely bad when decay = 0.001 is used. This might be due to underfitting since the regularization term is relatively large for the network. When choose decay =  $10^{-12}$ , the the training performance is the best. The following curve shows the test accuracy for different decay values. As the figure shown, decay =  $10^{-12}$  also provides the best test result.



## 1.8 Question 5

Question 5 is to design a 4-layer network with two hidden-layers and compare the performance of it with 3-layers network.

### 1.8.1 Method

The model contains 4 layers, input layer, 2 hidden layers and output layer. The activation function for two hidden layers is logistic activation. The output layer uses softmax function. The weights and bias are random initialized. We set the default batch size as 32. We use Stochastic gradient descent (shortened to SGD) with mini batch to train the network. The cost function is defined as the mean of the crossentropy of prediction and actual value then plus regularization term. Once the predefined Theano function train is called, the cost will be computed then the weight and bias will be updated by the SGD function. The learning rate  $\alpha = 0.01$ . The decay parameter  $\beta = 10^{-6}$ .

```

1  # theano expressions
2  X = T.matrix() #features
3  Y = T.matrix() #output
4
5  #weights and biases from input to 1st-hidden layer
6  w1, b1 = init_weights(36, 10), init_bias(10)
7  #weights and biases from 1st-hidden to 2nd-hidden layer
8  w2, b2 = init_weights(10, 10), init_bias(10)
9  #weights and biases from 2nd-hidden to output layer
10 w3, b3 = init_weights(10, 6, logistic=False), init_bias(6)
11
12 h1 = T.nnet.sigmoid(T.dot(X, w1) + b1)
13 h2 = T.nnet.sigmoid(T.dot(h1, w2) + b2)
14 py = T.nnet.softmax(T.dot(h2, w3) + b3)
15
16 y_x = T.argmax(py, axis=1)
17
18 cost = T.mean(T.nnet.categorical_crossentropy(py, Y)) + decay*(T.sum(T.sqr(w1))+T.
19     sum(T.sqr(w2))+T.sum(T.sqr(w3))))
20 params = [w1, b1, w2, b2, w3, b3]
21 updates = sgd(cost, params, learning_rate)
22
23 # compile
24 train = theano.function(inputs=[X, Y],
25     outputs=cost,
26     updates=updates,
```

```

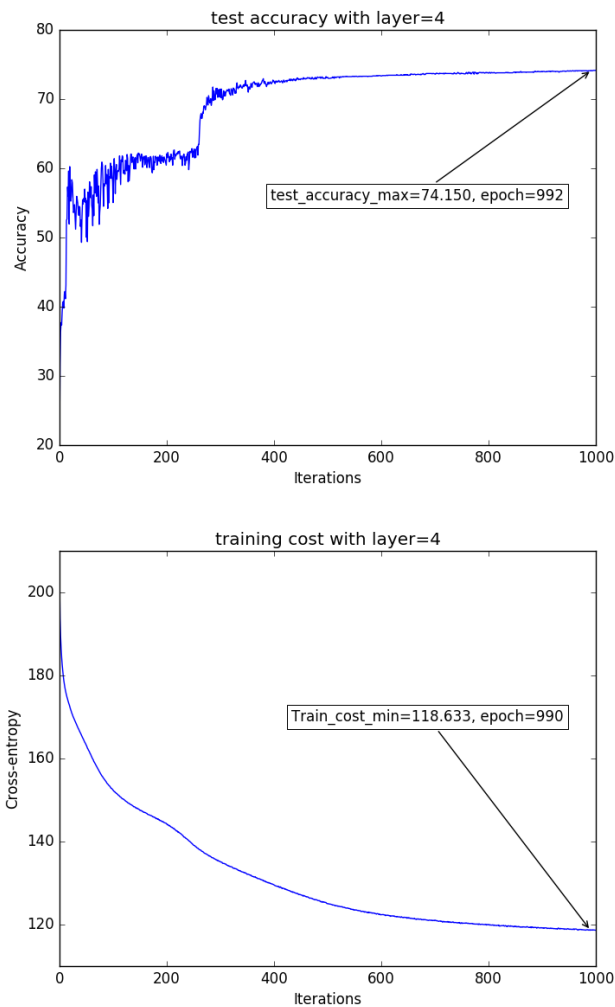
26         allow_input_downcast=True)
27     predict = theano.function(inputs=[X],
28                               outputs=y_x,
29                               allow_input_downcast=True)

```

At each epoch, the data will be shuffled. An average cost for this epoch is computed and stored in the list `train_cost`. This part of code is same as 1.4.1.

### 1.8.2 Result

The following figures show the accuracy and train cost of 4-layers network.



### 1.8.3 Discussion

Obviously, 4-layers network cannot perform as well as 3-layers network. The 3-layers network gets about 5 percentage point higher than 4-layers. This might be because we didn't optimize the parameters used for 4-layers. Overfitting since more neurons are involved can also be a reason. According to the figures, it seems that 4-layers network fell into some local optimum.

## 2 Part 2

### 2.1 Introduction

The second part of the project is to use Neuron Networks to deal with regression problem. The goal is to predict the house price in California by training the network using the provided California Housing dataset. Since different parameters will affect the overall performance, it's necessary to try out several of them and select the one given the best result. In this section, we use multilayer feedforward networks to solve the problem.

### 2.2 Notes

For Part 2, the main tools we use are Theano, numpy, scikit-learn and matplotlib. All the training and performance measurements are done on a Macbook Pro with 2.6 GHz Intel Core i5.

### 2.3 Data Preparation

The data is read from file *cal.housing.data* and separated into two parts with training and testing ratio of 0.7:0.3. Then perform scaling on training and testing data to let the transformed data be distributed within a certain range. Note that we use the min-max value computed from train data for the scaling on test data. This is because we pretend that the test data is unseen data which is only used for evaluating the performance of the network.

```
1 def scale(X, X_min, X_max):
2     return (X - X_min)/(X_max - X_min)
3
4 def read_data():
5     #read and divide data into test and train sets
6     cal_housing = np.loadtxt('cal_housing.data', delimiter=',')
7     X_data, Y_data = cal_housing[:, :8], cal_housing[:, -1]
8     Y_data = (np.asmatrix(Y_data)).transpose()
9
10    X_data, Y_data = shuffle_data(X_data, Y_data)
11    #separate train and test data
12    m = 3*X_data.shape[0] // 10
13    testX, testY = X_data[:m], Y_data[:m]
14    trainX, trainY = X_data[m:], Y_data[m:]
15
16    # scale data
17    trainX_max, trainX_min = np.max(trainX, axis=0), np.min(trainX, axis=0)
18
19    trainX = scale(trainX, trainX_min, trainX_max)
20    testX = scale(testX, trainX_min, trainX_max)
21
22    return trainX, trainY, testX, testY
```

### 2.4 Question 1

The first question is to develop a basic 3-layer feedforward neuron network with the hidden layer having 30 neurons. The plots for training error and prediction error are shown in the result section.

#### 2.4.1 Method

We set the parameters for the model as follow,

batch\_size = 32; learning\_rates = 1e-4; epochs = 1000; no\_hidden\_neurons = 30.

The model contains 3 layers, input layer, hidden layer and output layer, and the activation functions for the hidden layer and output layer are sigmoid and linear respectively with random initialized weights and bias. We use mini-batch gradient descent to train the network with the loss function

defined as the absolute mean square error between the generated predictions( $\hat{d}$ ), and the actual value( $y$ ). Once the predefined theano function *train* is being called, the corresponding cost would be computed, meanwhile, the weights and bias are updated by the gradient values.

```

1 # initialize weights and biases for hidden layer(s) and output layer
2 w_o = theano.shared(np.random.randn(no_hidden1)*.01, floatX )
3 b_o = theano.shared(np.random.randn()*0.01, floatX)
4 w_h1 = theano.shared(np.random.randn(no_features , no_hidden1)*.01, floatX )
5 b_h1 = theano.shared(np.random.randn(no_hidden1)*0.01, floatX)
6
7 #Define mathematical expression:
8 h1_out = T.nnet.sigmoid(T.dot(x, w_h1) + b_h1)
9 y = T.dot(h1_out, w_o) + b_o
10
11 cost = T.abs_(T.mean(T.sqr(d - y)))
12 accuracy = T.mean(d == y)
13
14 #define gradients
15 dw_o, db_o, dw_h, db_h = T.grad(cost, [w_o, b_o, w_h1, b_h1])
16
17 train = theano.function(
18     inputs = [x, d],
19     outputs = cost,
20     updates = [[w_o, w_o - alpha*dw_o],
21                [b_o, b_o - alpha*db_o],
22                [w_h1, w_h1 - alpha*dw_h],
23                [b_h1, b_h1 - alpha*db_h]],
24     allow_input_downcast=True
25 )

```

At each epoch, shuffle the data and train the model using mini batches. An average cost for this epoch is computed and stored in the *train\_cost* list. The training procedure is shown as follow,

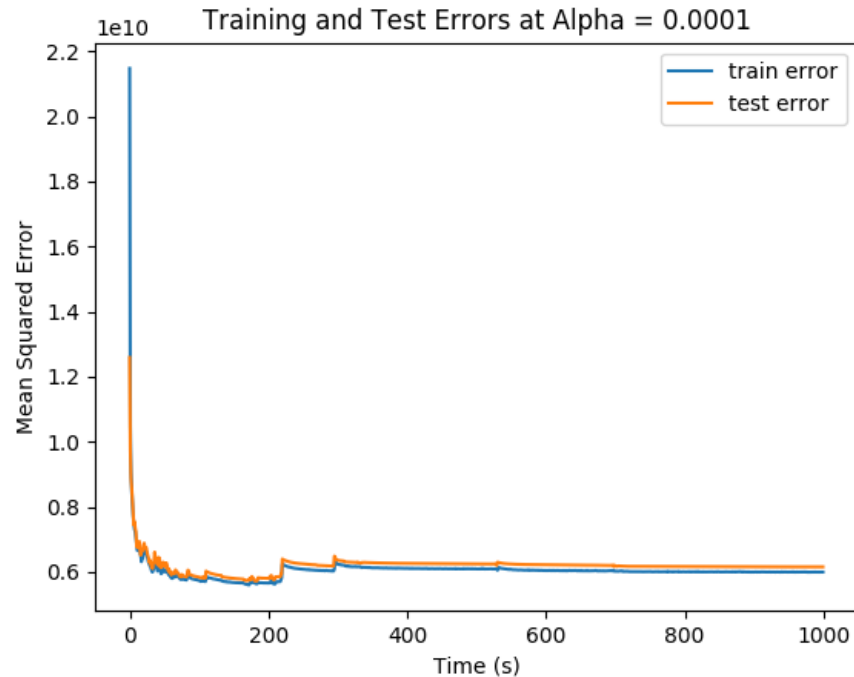
```

1 for iter in range(epochs):
2     if iter % 100 == 0:
3         print(iter)
4
5     trainX, trainY = shuffle_data(trainX, trainY)
6     n = len(trainX)
7     cost = 0
8     for start in range(0, n, batch_size):
9         end = start + batch_size
10        cost += train(trainX[start:end], np.transpose(trainY[start:end]))
11    train_cost[iter] = cost/(n // batch_size)
12    pred, test_cost[iter], test_accuracy[iter] = test(testX, np.transpose(testY))

```



### 2.4.2 Result



As we can see, the error curve is a bit noisy at beginning. This probably because when use mini-batch, the frequent updates can result in a noisy gradient signal, which may cause the model parameters and in turn the model error to jump around (having a higher variance over training epochs). One positive point for the noisy update process is that it can allow the model to avoid local optimum. The model eventually converges at a later point.

## 2.5 Question 2

By choosing different learning rates to do gradient update, the performance varies. The second question is to find the optimum learning rate for the network.

### 2.5.1 Method

Five-fold cross validation is implemented using Kfold method from scikit-learn. The basic idea behind is that the training data will be divided into five equally sized parts. Each time one set of the data will be picked for validation purpose, the rest is used for training. Repeat the procedure for 5 times, with each of the 5 parts is used exactly once as the validation data. In the end, an average result will be generated from the previous five as a single estimate. Note that in each iteration, a new model is generated for training. We tried different learning rate from the list  $\{10^3, 0.5 \times 10^3, 10^4, 0.5 \times 10^4, 10^5\}$ . Use a for loop to iterate through the list and execute the code.

```
1 for l in learning_rates:
2     alpha.set_value(l)
3     for train_index, val_index in kf.split(trainX):
4         val_set_X = trainX[val_index]
5         val_set_Y = trainY[val_index]
6         train_set_X = trainX[train_index]
7         train_set_Y = trainY[train_index]
```

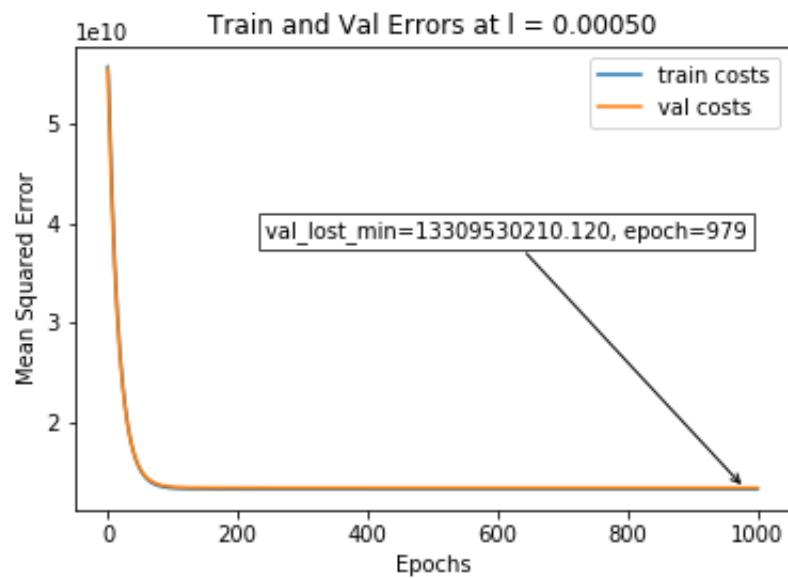
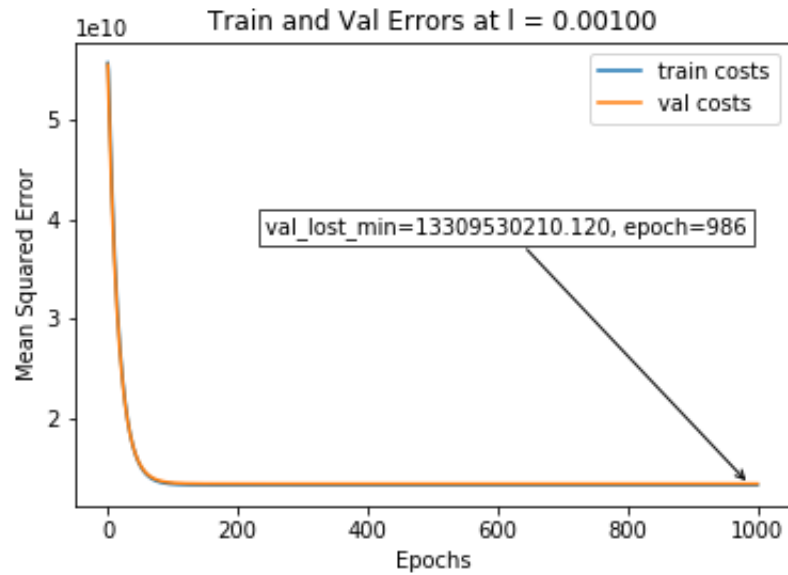
```

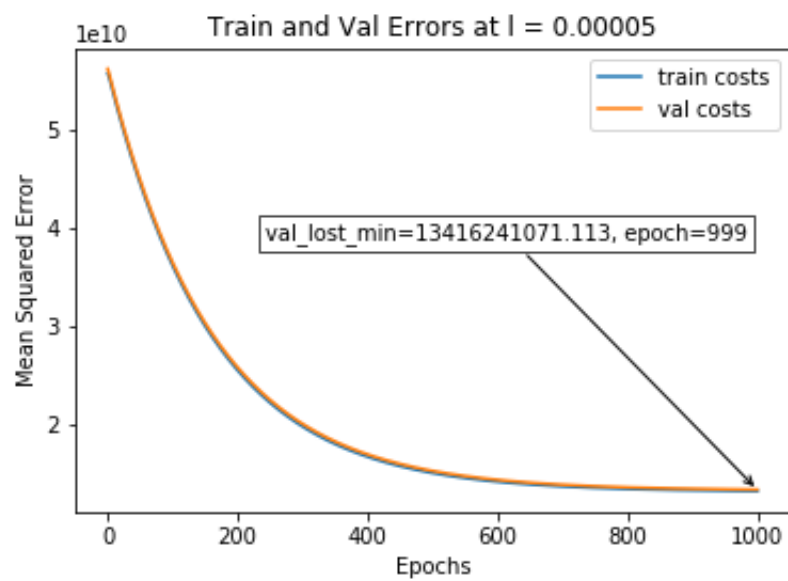
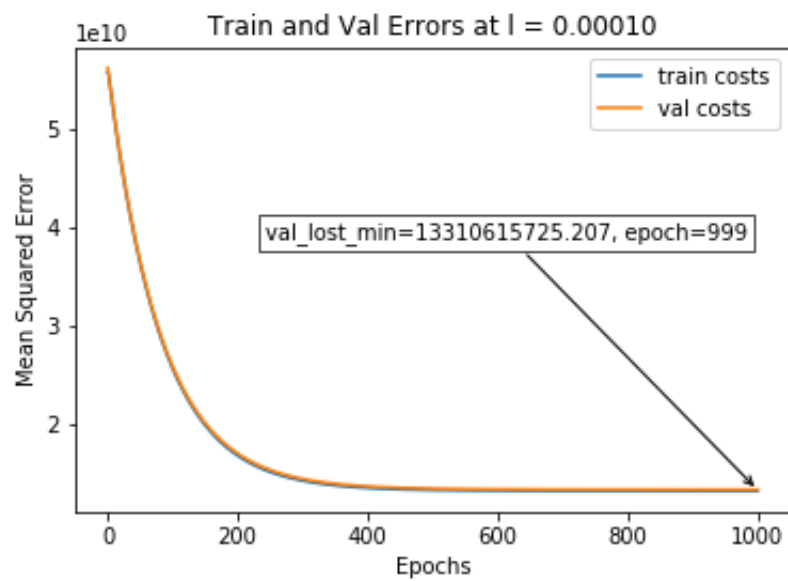
8     new_model()
9     for iter in range(epochs):
10        # ... followed by the same code

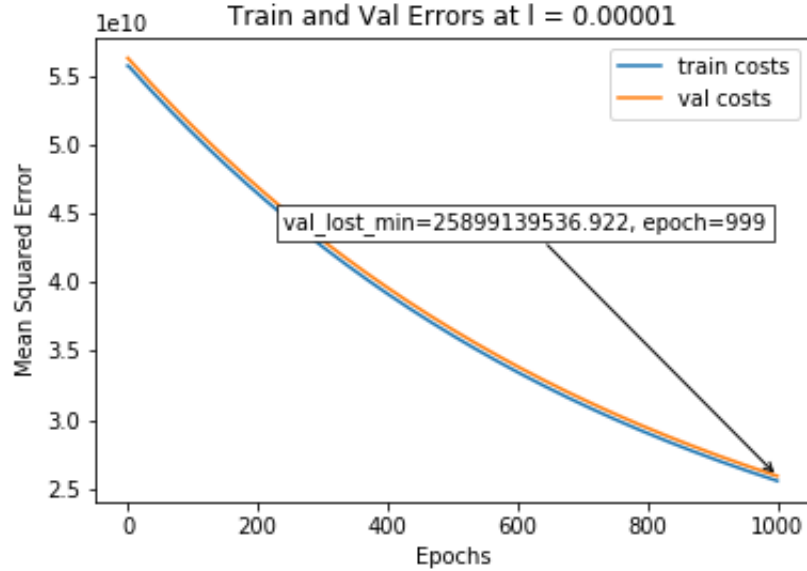
```

## 2.5.2 Result

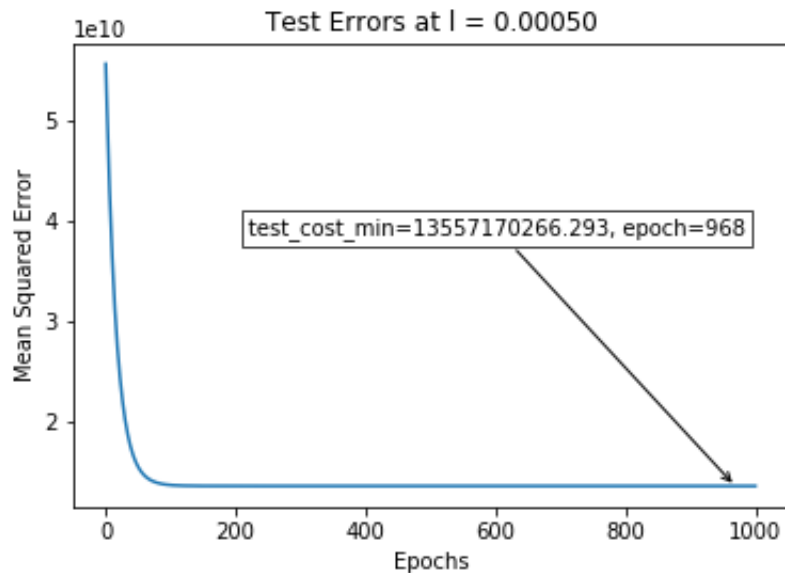
The following figures show the training and validation errors for models under different learning rates.







As shown above, we can see that  $l = 5e-4$  gives the best performance since it takes much less time to converge. Having very large learning rate may not decrease the gradient on every iteration and thus it may not converge well, or it might take some time to converge. While although small learning rate will lead to convergence in the end, it's always much slower. For example, when  $l=1e-5$ , it's obviously not converged at epoch 1000. The test errors curve with optimum learning rate is shown as follow,



## 2.6 Question 3

The number of neurons also determines the final performance, thus it's important to select a proper amount of neurons for the hidden layer.

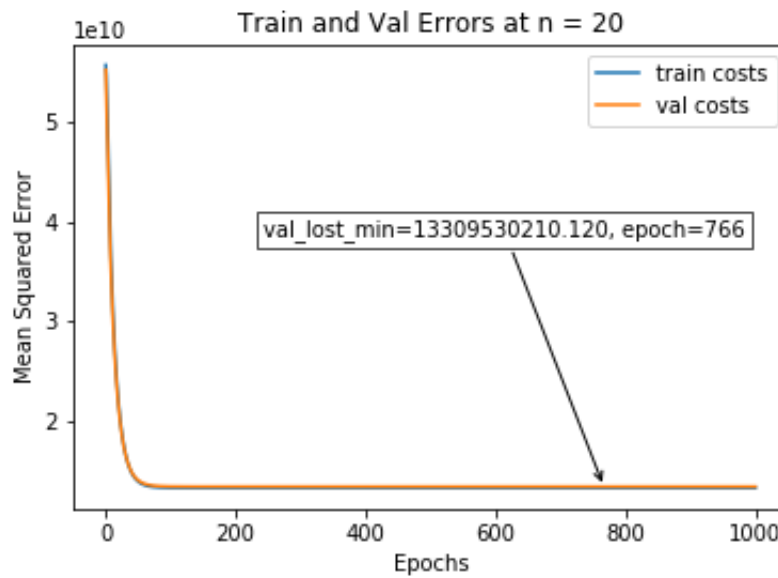
### 2.6.1 Method

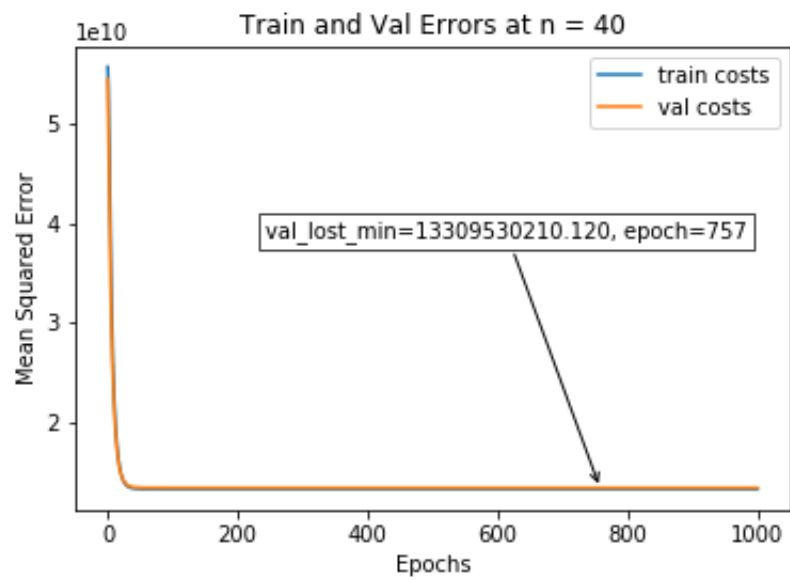
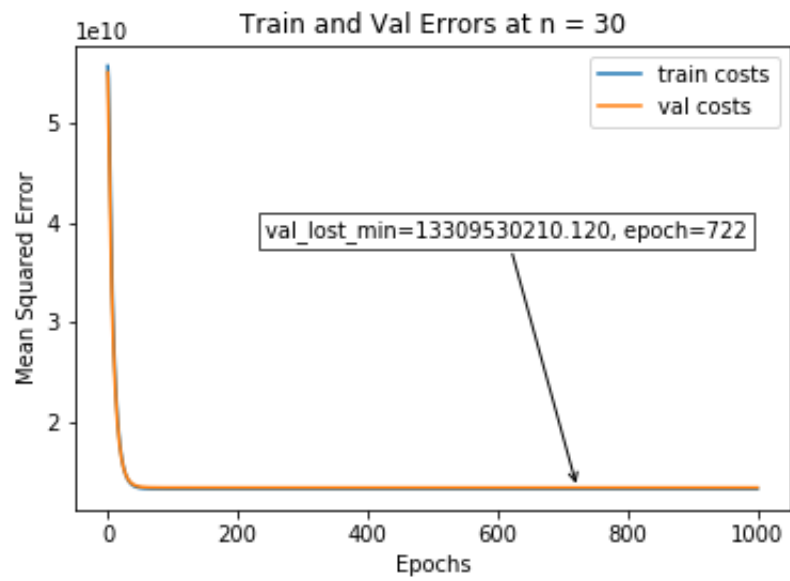
We set learning rate to be  $5e-4$ , which is the optimal one derived from Question 2, and use a for loop to evaluate the model with different number of hidden neurons. Same as Q2, we perform the 5 fold cross validation on the train data for choosing the best model.

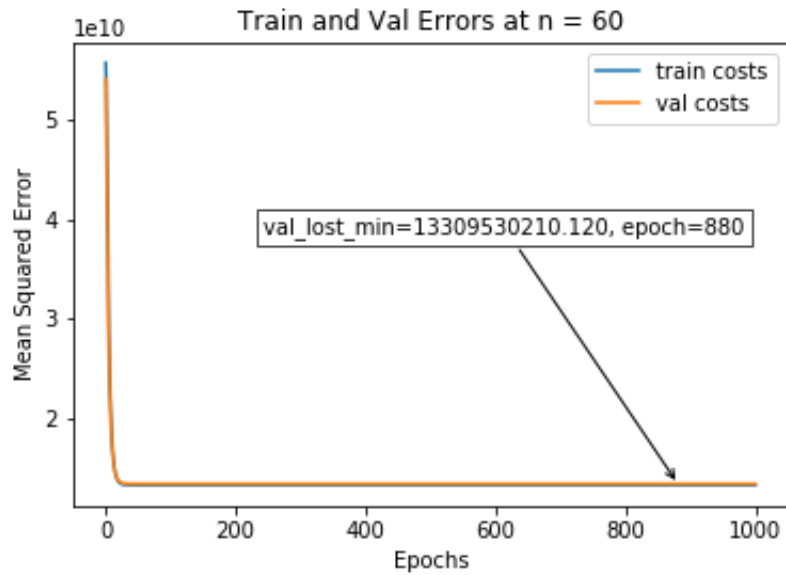
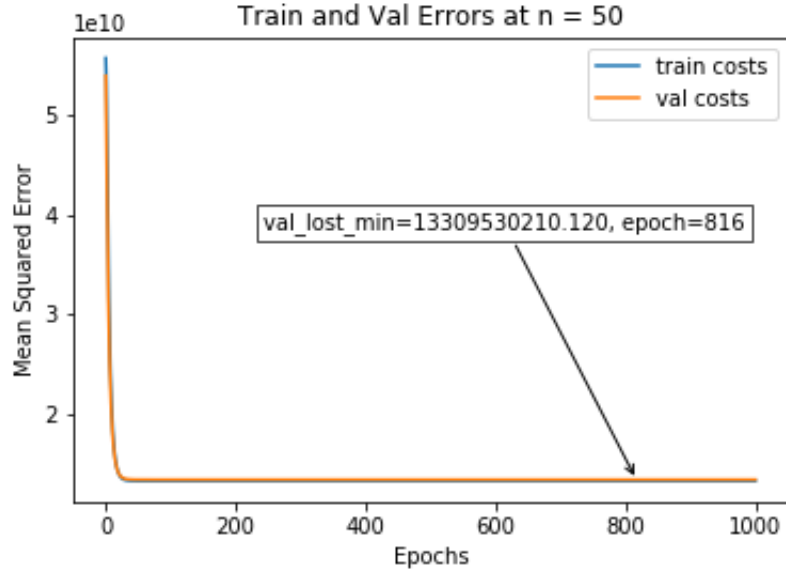
```
1 for no_hidden1 from no_hidden1_list:
2     for train_index, val_index in kf.split(trainX):
3         val_set_X = trainX[val_index]
4         val_set_Y = trainY[val_index]
5         train_set_X = trainX[train_index]
6         train_set_Y = trainY[train_index]
7         new_model()
8         for iter in range(epochs):
9             # ... followed by the same code
```

### 2.6.2 Result

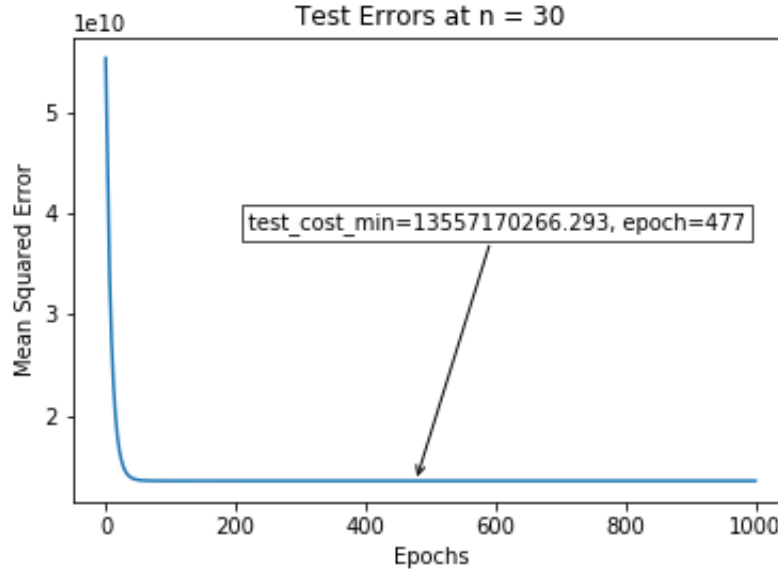
The following figures show the training and validation error for the models under different hidden-layer neurons.







As we can see, the validation costs are the same for all the models. But each converges at a different time. The overall trend is that along with the number of neurons in a hidden layer getting larger, the time taken for convergence is longer. But it seems that  $n=30$  converges faster than  $n=20$ . This is probably because that a bit more neurons would help to learn more relevant features from data, and thus could be earlier to reach optimum. Based on the result, we pick the one with  $n=30$ . The figure for test errors with `no_hidden = 30` is shown as follow,



## 2.7 Question 4

The main purpose for this part is to see whether a deeper network will bring a better result or not. Implement a four-layer neural network and five-layer neural network using the optimal parameters.

### 2.7.1 Method

On top of the previous implementation, we further add in another hidden layer with 20 neurons. The modifications are done based on the network structure as well as the connections between each layer. Similar modification is done for the five-layer neural network. For learning rate as well as number of neurons in the first hidden layer, we use  $l = 0.0001$ , `no_hidden1 = 30` which are the optimal ones chosen before. Note that in the section, we removed the k-fold cross validation because we are purely looking at the impact of depth of the network, rather than selecting the best model.

```

1 # define gradients
2 dw_o, db_o, dw_h1, db_h1, dw_h2, db_h2 = T.grad(cost, [w_o, b_o, w_h1, b_h1, w_h2,
3               b_h2])
4
5 # define mathematical expression:
6 h1_out = T.nnet.sigmoid(T.dot(x, w_h1) + b_h1)
7 h2_out = T.nnet.sigmoid(T.dot(h1_out, w_h2) + b_h2)
8 y = T.dot(h2_out, w_o) + b_o
9
10 # modify the train function
11 train = theano.function(
12     inputs = [x, d],
13     outputs = cost,
14     updates = [[w_o, w_o - alpha*dw_o],
15                [b_o, b_o - alpha*db_o],
16                [w_h1, w_h1 - alpha*dw_h1],
17                [b_h1, b_h1 - alpha*db_h1],
18                [w_h2, w_h2 - alpha*dw_h2],
19                [b_h2, b_h2 - alpha*db_h2]],
20     allow_input_downcast=True)
21
22 # generate new weight and bias

```



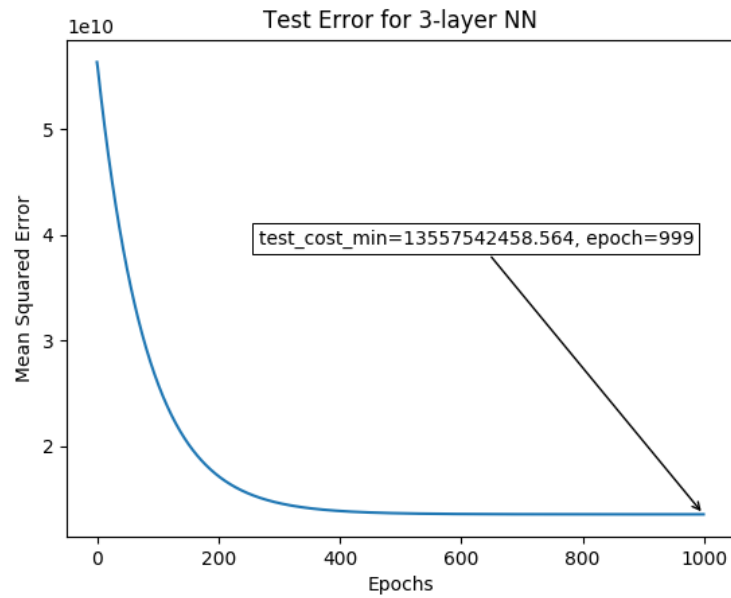
```

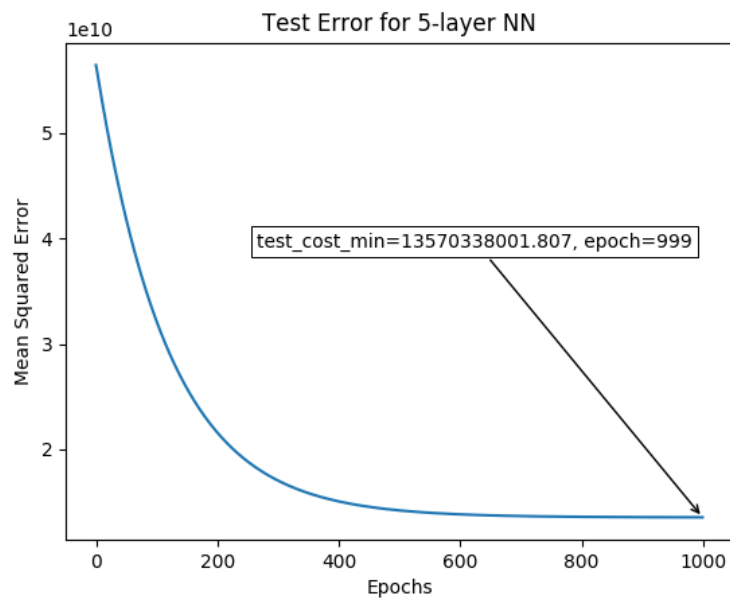
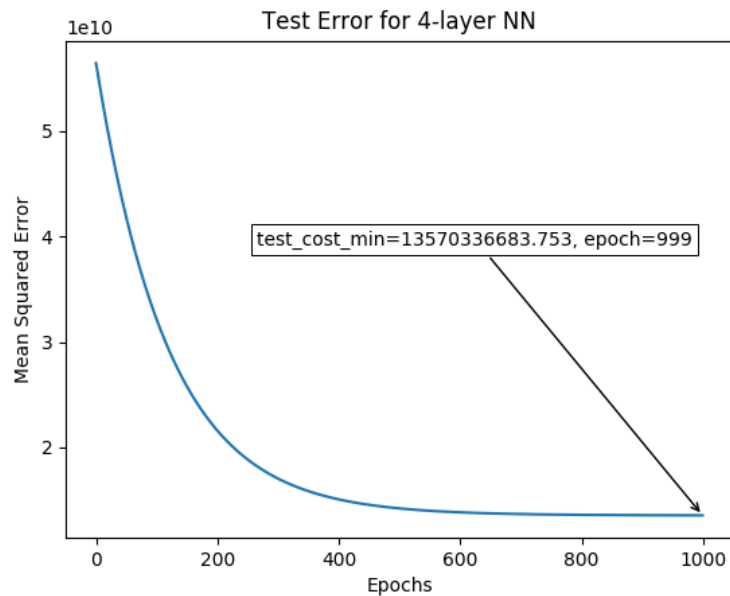
22 def new_model(no_hidden1, no_hidden2, no_features):
23     w_o.set_value(np.random.randn(no_hidden2)*.01)
24     b_o.set_value(np.random.randn()*0.01)
25     w_h1.set_value(np.random.randn(no_features, no_hidden1)*.01)
26     b_h1.set_value(np.random.randn(no_hidden1)*.01)
27     w_h2.set_value(np.random.randn(no_hidden1, no_hidden2)*.01)
28     b_h2.set_value(np.random.randn(no_hidden2)*.01)

```

### 2.7.2 Result

The test errors curves for the networks with different depth are shown as follow,





### 2.7.3 Discussion

Compare with the three-layer network, it seems that the overall performance drop a bit when the network getting deeper since it takes longer time to converge. The result may be better if we choose another set of parameters for deeper networks. But if with a certain predefined parameters, deeper network may not always outperform others, especially with increasingly computational effort.