

Automatical Learning and Data Mining

HW4 Report

Yifan Zhao

April 1, 2018

1 Description

The data set is the same as the second task described in the homework text. We downloaded about 1300 stocks' per minute price sequences during a 20-day period from google finance websites. Then select those who have more than 5000 cases (Some stocks may miss some cases for delisted). Then we have 757 stocks left. We chose two of them as the target stocks. And decide the different classes according to the up and down of these two stocks:

- 1: up up;
- 2: up down;
- 3: down up;
- 3: down down.

Then the number of descriptors is 755. The number of cases is 5000. The number of classes is 4. All descriptors are real and continuous. Expected outputs are 4-dimension 0-1 vectors valued 1 if and only if the index is equal to the correct class index, otherwise, equal to 0. The training set has 4000 cases, the test set has 1000 cases. Before separate the training set and the test set, we shuffle the whole data set in order to remove the influence of chronical trend.

2 Architecture of Auto-encoder

2.1

We construct the auto-encoder based on the MLP paradigm. The auto-encoder has three layers:

$$INPUT \rightarrow HIDDEN \rightarrow OUTPUT$$

where INPUT has the same size as OUTPUT, i.e., the number of descriptors.

2.2 PCA of Original Data

Before computing the correlation matrix of the data, we do following manipulations for each descriptor: Subtract its median(or mean) and divide it with its standard deviation. Subtracting median instead of mean could make the data unbiased because of the definition of median, so we choose median here. After performing the PCA, the 3 cut-off numbers of eigenvalues are as follows :

- R1=50%: 5
- R2=75%: 30
- R3=90%: 91.

So the sizes of auto-encoder's hidden layer are $h_1=5$, $h_2=30$, $h_3=91$.

The graphs of eigenvalues and their cumulative ratios are as follows:

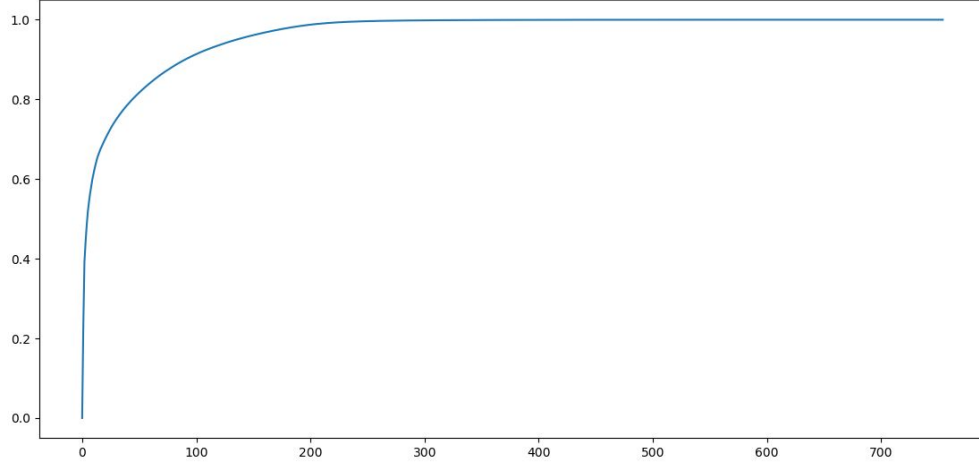


Figure 1: Ratio

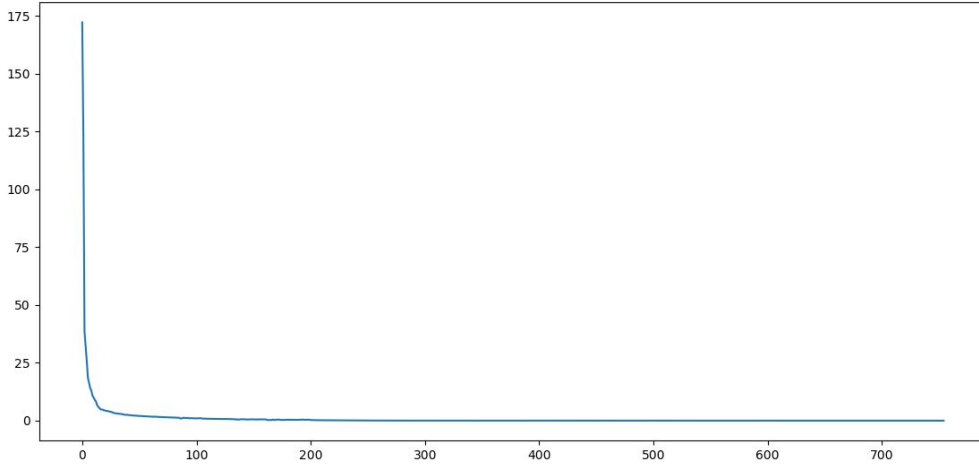


Figure 2: Eigenvalues

2.3 Pretreatment of Original Data

We use the sigmoid response function. Since we have already subtracted the median and divided with standard deviation of the data. And the main non-linear part of sigmoid function is in $(-6,6)$, we do not subtract the minimum of the data, nor divide with the range(maximum minus minimum). The nonlinearity of sigmoid function could be better developed.

3 Automatic Learning for 3 AEs

3.1 Software Options

We use tensorflow as the software for automatic learning. There are two usual tools for minimizing the MSE on the training set with validation on the test set, ADAM optimizer(function in tensorflow: `tf.train.AdamOptimizer`) and Gradient Descent optimizer(function in tensorflow: `tf.train.GradientDescentOptimizer`). The former is a specialized gradient descent optimizer for classification tasks, and the latter is the generic gradient descent optimizer. So we use Gradient Descent optimizer in the auto-encoders.

With tensorflow, we could generate various distributed including but not limited to normal distributed tensors for

initialization of the wights and assign parameters like means and standard deviations to specify the initialization. One of the appropriate specifications we found is that mean equals to 0.1 and standard deviation equals to 0.01. For batch learning, we use the function `np.random.choice` in Python package `numpy` to obtain a batch. This function returns a random sample from the population with a given size. The current tested appropriate batch size is around 100.

For the learning rate, we assign a learning rate variable valued $\frac{\epsilon}{n}$ to control the step size of gradient, where n the current number of global learning steps or epochs(the latter is better according to several tests).

The stopping criterion of learning is when the norm of gradient is smaller than a given level: e.g., 0.03. For intermediary outputs to monitor learning quality, we compute and print the learning loss per 10 batches. The loss is defined by the mean squared error between reconstructed data and the original input data out of the whole batch.

3.2 RMSE, ΔW , Gradient

The computation of RMSE, ΔW and gradient is similar to ordinary MLPs. RMSE is just the square root of MSE, it can be computed by the function `tf.sqrt`, then normalize it with its dimension. For the $||\Delta W||$, we first compute the norm of the difference of weights and biases of each layer separately. Then we add the squared norms together and normalize it with its dimension. In this way we could avoid to manipulate the matrices of weights and biases, which are pretty complicated. The computation of gradient could be easier after we have the weights' difference, just divide it with the learning rate and normalize it with its dimension.

Then we could plot the three curves:

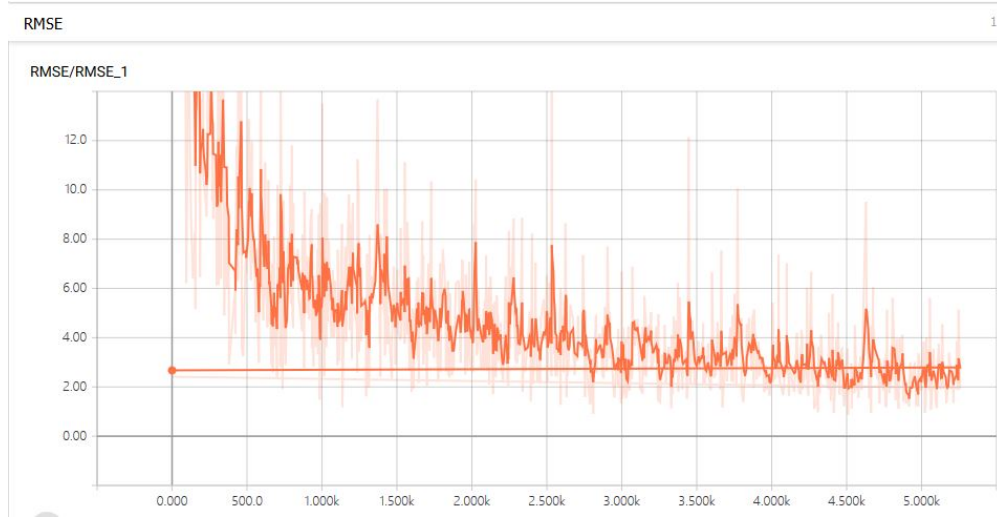


Figure 3: RMSE

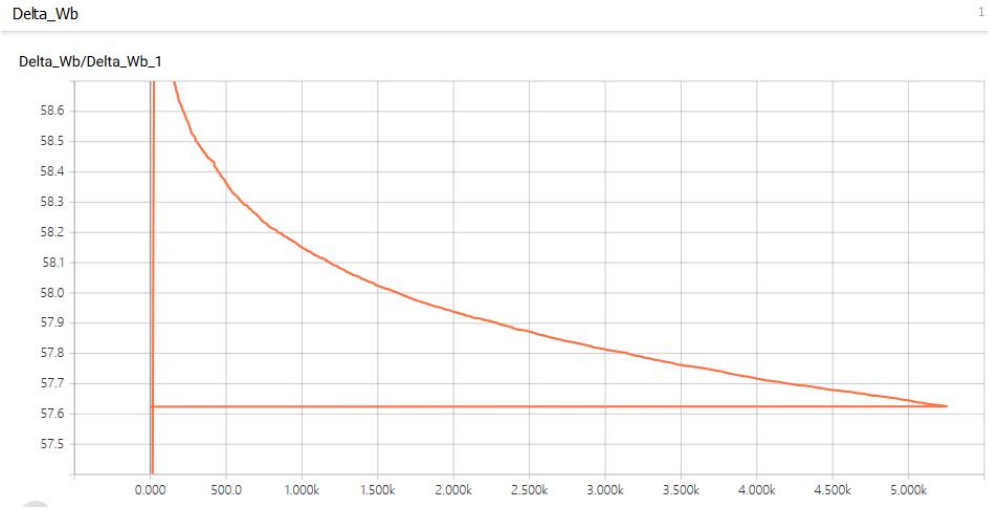


Figure 4: $\|\Delta W\|$

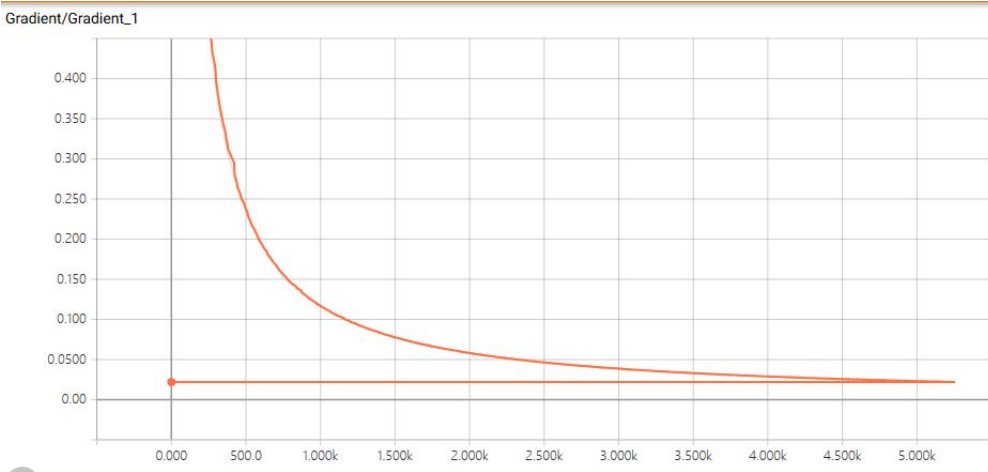


Figure 5: Gradient

As we could see, all three curves decline in an exponential shape as the global steps grow. The fluctuation of RMSE is the greatest among these three curves. The declination of gradient is relatively insignificant comparing with the other two. The learning rate influences the gradient a lot. Therefore, we increased the number of epochs in order to guarantee the auto-encoder converge sufficiently.

3.3 Global Performance

Here is the graph of two GRMSEs on the same plot.

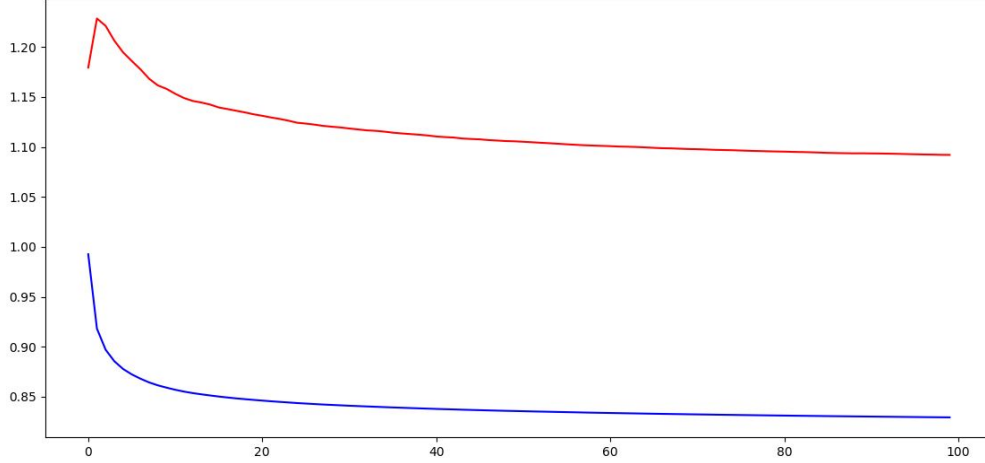


Figure 6: GRMSE

The blue one represents the training set, the red one represents the test set. As we could see from the graph, the performances on both set are improving as the number of epochs grows. The GRMSE on the test set is slightly greater than that on the training set, yet still comparable.

3.4 Performances of Three AEs

Since we have 3 different auto-encoders with hidden layer dimensions 5, 30 and 91, we perform several runs of each auto-encoder. The best result is from the 91-dimensional one with batch loss(MSE) converges to 0.06. Obviously, the first two dimensions are too small for reconstructing data with 755 dimensions. Additionally, 91 is a relatively successful compression of data compared with 755. So we select the third auto-encoder as the best.

4 Activity Analysis of Hidden Layer

4.1 PCA of Hidden Layer

After the training, we save the terminal values of weights and biases. Then we compute the terminal values of hidden layer on the training set. Then implement PCA for these hidden layer vectors. The number of eigenvalues of 90% energy for the first auto-encoder is 3, the second is 11 and the third is 26. Here is the scatter graph of the cloud of the third auto-encoder. Since graphs for different auto-encoders are not significantly distinct, so we only display the one of 91 dimensions. Points in different classes have different colors.

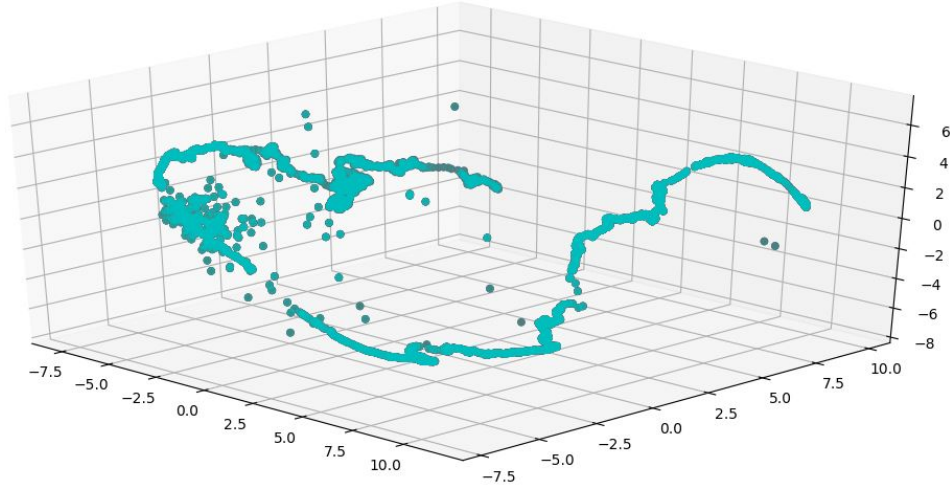


Figure 7: The Cloud of Hidden Layer Vectors

As we could see from the scatter, points of different classes have not been separated enough, for auto-encoder is only a pretreatment and simplification of the data. So the implementation of MLP is still necessary.

4.2 The Maximum and Activity of Each Node

We compute the maximum of each node with following codes. First, we initialize the list variable to store the indices of maximums of each node. Then we find the input such that one specific node attains its maximum by the function `np.argmax`, it returns the argument(index) such that the input attains its maximum. Then we find the class it belongs to with return value of `np.argmax` and the array of each input's class index. At last, we count the number of these inputs belonging to each class and return with the list `hid_max_sum`. Here is the list of the number of nodes attributed to each class:

[14, 4, 0, 73]

. As we could see, the fourth class is the most active.

The activity of each node is also computed by the following codes. First, we initialize an activity list variable to store the activities. Then for each node, we compute the mean of the node based on the whole set with the function `np.mean`.

```

1  hid_max=[]
2  activity=[]
3  class_act=[[[]],[[]],[[]],[[]]]
4  for i in range(h):
5      hid_max_index=np.argmax(hidden[:,i])
6      hid_max.append(np.argwhere(label[hid_max_index,:]))
7      activity.append(np.mean(hidden[:,i]))
8      for j in range(n_classes):
9          class_act[j].append(np.mean(hidden[label_train_index[j],i]))
10 hid_max=np.asarray(hid_max)
11 hid_max_sum=[]
12 for i in range(n_classes):
13     hid_max_sum.append(np.count_nonzero(hid_max==[[i]]))
14 print(hid_max_sum)

```

Then we reorder the activities in increasing order and plot:

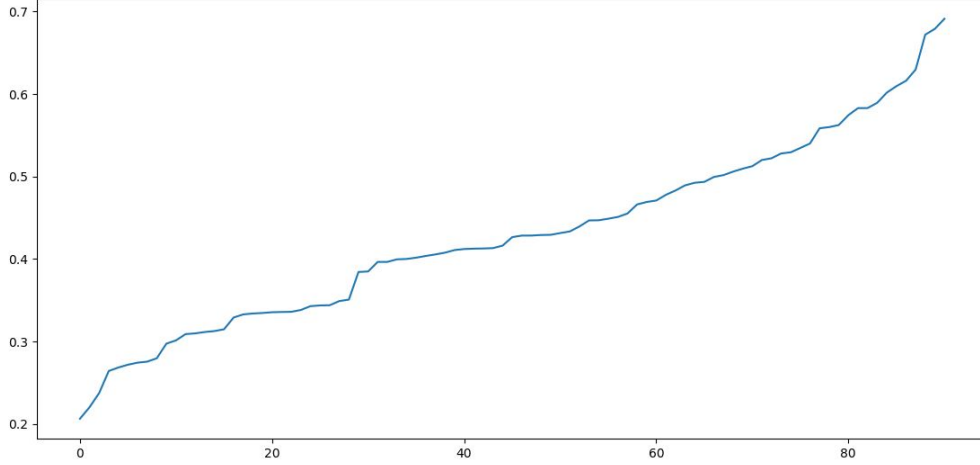


Figure 8: Activity

As we could see, the activities are all between 0.2 and 0.7, and the number of low activity nodes is small, thus no extreme situations. This is an acceptable result. To identify the nodes with low activities, we first use the function `np.argsort` to compute the index after reordering `Q` in increasing order. Then extract those nodes which are in the first 5% of the array. The 5% quantile of `Q` is 4 nodes: 13, 35, 46, 71.

If we remove these nodes, the loss on the whole set will improve from 0.82 to 0.69. Relevant Python codes are as follows. We first initialize the weights variable to store new weights and assign it with all-zero matrix temporarily. Then assign it with the trained terminal weights only in those active nodes indices. Thus we have got the weights matrix with specific rows equal to zero. Then compute the loss with this new matrix.

```

1 | act_sort=np.argsort(activity)
2 | inact_sort=(0.05*len(act_sort))
3 | act_sort=act_sort[int(inact_sort):]
4 | W_act=np.zeros([n_des,h])
5 | b_act=np.zeros([h])
6 | W_act[:,act_sort]=W[:,act_sort]
7 | b_act[act_sort]=b[act_sort]
8 | hid_act=sess.run(tf.sigmoid(np.matmul(train,W_act)+b_act))
9 | out_act=np.matmul(hid_act,W_out)+b_out
10 | loss_act = sess.run(tf.losses.mean_squared_error(labels=train, predictions=out_act))
11 | print("Loss after removing inactive nodes: {}".format(loss_act))

```

Similar analysis for each class could be achieved with following codes:

```

1 | for j in range(n_classes):
2 |     class_act[j].append(np.mean(hidden[label_train_index[j],i]))

```

We compute the activities for each class separately, then we could do the same analysis with the former codes for each class.

5 Auto-encoder Insertion in MLP Classifiers

5.1 Architecture of MLP

First, we compute the values of AE's hidden layer on the whole set with the terminal weights. Then use these vectors as new inputs. Thus the weights in the auto-encoder are frozen. According to the PCA of hidden layer, we choose the size of MLP's hidden layer `hh=26`. After several implementations of the MLP, we modify `hh=21` as the best size of hidden layer.

5.2 Confusion Matrix

We first assign a variable that signifies the class that the cases truly belonged to. Then we compare the class predicted by the output vector with this index variable. Then we compute the ratio of each kind of prediction. Finally after running the session, we could store the results to two matrices and output them.

The training confusion matrix:

```
[0.7238636363636364, 0.003977272727272727, 0.006439393939393939 0.265719696969697]
[0.000946969696969697, 0.8005681818181818, 0.08958333333333333, 0.08844696969696969]
[0.0028414472437961733, 0.37715476415987875, 0.4995264254593673, 0.04546315590073877]
[0.005871212121212121, 0.11022727272727273, 0.014772727272727272, 0.8615530303030303]
```

The test confusion matrix:

```
[0.7328125, 0.005729166666666667, 0.004166666666666667, 0.25729166666666664]
[0.00625, 0.6572916666666667, 0.17447916666666666, 0.16145833333333334]
[0.42239583333333336, 0.26666666666666666, 0.28645833333333333, 0.0234375]
[0.0067708333333333336, 0.2140625, 0.022916666666666665, 0.7473958333333334]
```

6 Appendix: Python Codes

```
1 import tensorflow as tf
2 import numpy as np
3 import pandas as pd
4 from matplotlib.mlab import PCA
5 import matplotlib.pyplot as plt
6 from ops import fc_layer, batch_ae, batch
7 from mpl_toolkits.mplot3d import Axes3D
8
9 count=0
10 df=pd.read_csv('F:/Files/DataMining/HW4/companylist.csv',
11                usecols=[0], skiprows=[0], header=None)
12 list=df.values
13 for i in range(len(list)):
14     filename = str(list[i])[2:-2] + '.csv'
15     try:
16         df = pd.read_csv('F:/Files/DataMining/HW4/data/' + filename,
17                          usecols=[2], skiprows=[0, 1, 2, 3, 4, 5, 6], header=None)
18     except FileNotFoundError:
19         pass
20     except pd.io.common.EmptyDataError:
21         pass
22     if len(df.values)>=5010:
23         col=df.values[0:5010]
24         col=np.nan_to_num(col)
25         m=np.median(col)
26         sigma = np.sqrt(np.var(col))
27         col=(col-m)/sigma
28         if count==0:
29             data=col
30             count=1
31         else:
32             data=np.append(data,col,axis=1)
33
34 label=np.empty([5000,4])
35 for j in range(5000):
36     if data[j,755]<data[j+10,755]:
37         if data[j,756]<data[j+10,756]:
```



```

38         label[j,:]=[1,0,0,0]
39     else:
40         label[j, :] = [0, 1, 0, 0]
41     else:
42         if data[j,756]<data[j+10,756]:
43             label[j,:]=[0,0,1,0]
44         else:
45             label[j, :] = [0, 0, 0, 1]
46
47 label_train=label[0:4000,:]
48 label_test=label[4000:5000,:]
49 data=data[0:5000,0:755]
50 data=np.nan_to_num(data)
51 n_des=data.shape[1]
52 n_cases=data.shape[0]
53 data_label=np.append(data,label,axis=1)
54 train=data[0:4000,:]
55 test=data[4000:5000,:]
56 label_train_index=[[],[],[],[],[]]
57 label_test_index=[[],[],[],[],[]]
58 for i in range(4):
59     label_train_index[i]= label_train[:,i]==1
60     label_test_index[i]=label_test[:,i]==1
61
62 pca=PCA(data)
63 co=np.corrcoef(np.transpose(data))
64 eig=np.linalg.eig(co)
65
66 for i in range(n_des):
67     if np.sum(eig[0][0:i])>=0.5*np.sum(eig[0]):
68         r1=i
69         break
70 for i in range(n_des):
71     if np.sum(eig[0][0:i])>=0.75*np.sum(eig[0]):
72         r2=i
73         break
74 for i in range(n_des):
75     if np.sum(eig[0][0:i])>=0.9*np.sum(eig[0]):
76         r3=i
77         break
78 rat=[]
79 for i in range(755):
80     rat.append(np.sum(pca.fracs[0:i]))
81
82
83 epochs=100
84 test_size=1000
85 batch_size=100
86 display_freq=10
87 h=r3
88 hh=20
89 n_classes=4
90
91 x = tf.placeholder(tf.float64, shape=[None, n_des], name='X')
92 y = tf.placeholder(tf.float64, shape=[None, n_des], name='Y')
93 x1=tf.placeholder(tf.float64, shape=[None, h], name='X1')
94 z=tf.placeholder(tf.float64, shape=[None, n_classes], name='Z')
95 learning_rate = tf.placeholder(tf.float64, shape=None, name='epsilon')
96 W_hid2out_form=tf.Variable(tf.zeros(shape=[h,n_des]),name='W_hid2out_form')

```

```

97 W_in2hid_form=tf.Variable(tf.zeros(shape=[n_des,h]),name='W_in2hid_form')
98 b_hid2out_form=tf.Variable(tf.zeros(shape=[n_des]),name='b_hid2out_form')
99 b_in2hid_form=tf.Variable(tf.zeros(shape=[h]),name='b_in2hid_form')
100 W_hid2out_form=tf.cast(W_hid2out_form,tf.float64)
101 W_in2hid_form=tf.cast(W_in2hid_form,tf.float64)
102 b_hid2out_form=tf.cast(b_hid2out_form,tf.float64)
103 b_in2hid_form=tf.cast(b_in2hid_form,tf.float64)
104
105 fc1 = fc_layer(x, h, 'Hidden_layer', use_relu=True)
106 output_logits = fc_layer(fc1[0], n_des, 'Output_layer', use_relu=False)
107 fc11 = fc_layer(x1, hh, 'Hidden_layer1', use_relu=True)
108 output_logits1=fc_layer(fc11[0], n_classes, 'Output_layer1', use_relu=False)
109
110 #define the loss defined as mean squared error
111 with tf.name_scope('MSE'):
112     loss = tf.losses.mean_squared_error(labels=y, predictions=output_logits[0])
113     tf.summary.scalar('MSE',loss)
114 loss1 = tf.losses.mean_squared_error(labels=z, predictions=output_logits1[0])
115 #compute the root MSE
116 with tf.name_scope('RMSE'):
117     rmse=tf.sqrt(loss,name='RMSE')
118     tf.summary.scalar('RMSE',rmse)
119 #define the optimizer
120 with tf.name_scope('Optimizer'):
121     optimizer = tf.train.GradientDescentOptimizer(learning_rate=learning_rate, name='Op').minimize
122     optimizer1 = tf.train.GradientDescentOptimizer(learning_rate=learning_rate, name='Op').minimize
123 #compute the norm of the difference between the weights of the former and the present batches
124 delta_W=tf.square(tf.norm(fc1[1]-W_in2hid_form))+tf.square(tf.norm(output_logits[1]-W_hid2out_f
125 delta_b=tf.square(tf.norm(fc1[2]-b_in2hid_form))+tf.square(tf.norm(output_logits[2]-b_hid2out_f
126 with tf.name_scope('Delta_Wb'):
127     difference=tf.sqrt(delta_W+delta_b,name='Delta_Wb')
128     tf.summary.scalar('Delta_Wb',difference)
129 #compute the gradient
130 with tf.name_scope('Gradient'):
131     gradient=tf.multiply(1/learning_rate,difference,name='Gradient')
132     tf.summary.scalar('Gradient', gradient)
133
134 #compute the number of correct predictions
135 correct_prediction = tf.equal(tf.argmax(output_logits1[0], 1), tf.argmax(z, 1), name='co_pre')
136 #compute the percentage of correct predictions in the whole batch
137 accuracy = tf.reduce_mean(tf.cast(correct_prediction, tf.float64), name='accuracy')
138
139 #the variable that signifies the class of denominator
140 indice=tf.placeholder(tf.int64,shape=[1],name='Index')
141 #initialization of the two matrices
142 entry_train=[[[]], [[]], [[]], [[]], [[]]]
143 entry_test=[[[]], [[]], [[]], [[]], [[]]]
144 #compare the output vector to the given class index and count the numbers in order to compute t
145 prediction = tf.equal(tf.argmax(output_logits1[0], 1), indice)
146 #compute the ratio of prediction and the number of cases in one class
147 ratio = tf.reduce_mean(tf.cast(prediction, tf.float64))
148
149 W=[]
150 b=[]
151 grmse_train=[]
152 grmse_test=[]
153 init_op = tf.global_variables_initializer()
154 merged=tf.summary.merge_all()
155

```

```

156 with tf.Session() as sess:
157     sess.run(init_op)
158     train_writer = tf.summary.FileWriter(logdir="./logs/", graph=sess.graph)
159     num_iter = 40
160     global_step = 0
161     for epoch in range(epochs):
162         train_index = np.arange(len(train))
163         print("Training Epoch: {}".format(epoch + 1))
164         for iteration in range(num_iter):
165             global_step += 1
166             x_batch, train_index = batch_ae(batch_size, train[train_index, :])
167             feed_dict_batch = {x: x_batch, y: x_batch, learning_rate: 2/(epoch+1)}
168             _, summary_tr = sess.run([optimizer, merged], feed_dict=feed_dict_batch)
169             train_writer.add_summary(summary_tr, global_step)
170
171             diff_batch, grad_batch = sess.run([difference, gradient], feed_dict=feed_dict_batch)
172
173             W_hid2out_present = sess.run(output_logits, feed_dict=feed_dict_batch)[1]
174             W_in2hid_present = sess.run(fc1, feed_dict=feed_dict_batch)[1]
175             b_hid2out_present = sess.run(output_logits, feed_dict=feed_dict_batch)[2]
176             b_in2hid_present = sess.run(fc1, feed_dict=feed_dict_batch)[2]
177             sess.run(W_hid2out_form, feed_dict={W_hid2out_form: W_hid2out_present})
178             sess.run(W_in2hid_form, feed_dict={W_in2hid_form: W_in2hid_present})
179             sess.run(b_hid2out_form, feed_dict={b_hid2out_form: b_hid2out_present})
180             sess.run(b_in2hid_form, feed_dict={b_in2hid_form: b_in2hid_present})
181
182             if (iteration + 1) \% display_freq == 0:
183                 loss_batch = sess.run(loss, feed_dict=feed_dict_batch)
184                 print("Iteration: {0:3d}, Loss: {1:.2f}".format((iteration + 1), loss_batch))
185                 print('Delta W: {0:.2f} Gradient: {1:.4f}'.format(diff_batch, grad_batch))
186
187             x_train = train[:, 0:755]
188             x_test=test[:,0:755]
189             feed_dict_train = {x: x_train, y: x_train}
190             feed_dict_test = {x: x_test, y: x_test}
191             grmse_train.append(sess.run(rmse, feed_dict=feed_dict_train))
192             grmse_test.append(sess.run(rmse, feed_dict=feed_dict_test))
193
194         print("Global Step: {0:3d}".format(global_step))
195
196         hidden,W,b=sess.run(fc1,feed_dict={x: train})
197         _,W_out,b_out=sess.run(output_logits,feed_dict={x: train})
198         hid_test=sess.run(fc1,feed_dict={x:test})[0]
199         hidden=np.nan_to_num(hidden)
200         hid_test = np.nan_to_num(hid_test)
201         hid_pca=PCA(hidden)
202         color = ['b', 'g', 'r', 'c']
203
204         fig = plt.figure()
205         ax = fig.add_subplot(111, projection='3d')
206         for i in range(n_classes):
207             ax.scatter(hid_pca.Y[:, 0], hid_pca.Y[:, 1], hid_pca.Y[:, 2], c=color[i])
208         plt.show()
209
210         hid_max=[]
211         activity=[]
212         class_act=[[[]],[],[],[],[]]
213         for i in range(h):
214             hid_max_index=np.argmax(hidden[:,i])

```

```

215         hid_max.append(np.argwhere(label[hid_max_index,:]))
216         activity.append(np.mean(hidden[:,i]))
217         for j in range(n_classes):
218             class_act[j].append(np.mean(hidden[label_train_index[j],i]))
219     hid_max=np.asarray(hid_max)
220     hid_max_sum=[]
221     for i in range(n_classes):
222         hid_max_sum.append(np.count_nonzero(hid_max==[[i]]))
223     print(hid_max_sum)
224     print(np.sort(activity))
225     plt.plot(np.sort(activity))
226     plt.show()
227
228
229     act_sort=np.argsort(activity)
230     inact_sort=(0.05*len(act_sort))
231     act_sort=act_sort[int(inact_sort):]
232     W_act=np.zeros([n_des,h])
233     b_act=np.zeros([h])
234     W_act[:,act_sort]=W[:,act_sort]
235     b_act[act_sort]=b[act_sort]
236     hid_act=sess.run(tf.sigmoid(np.matmul(train,W_act)+b_act))
237     out_act=np.matmul(hid_act,W_out)+b_out
238     loss_act = sess.run(tf.losses.mean_squared_error(labels=train, predictions=out_act))
239     print("Loss after removing inactive nodes: {}".format(loss_act))
240
241     #plt.plot(grmse_train,c='b')
242     #plt.plot(grmse_test,c='r')
243     #plt.show()
244
245     epochs=10
246     with tf.Session() as sess:
247         sess.run(init_op)
248         num_iter = 40
249         global_step = 0
250         for epoch in range(epochs):
251             train_index=np.arange(len(train))
252             print("Training Epoch: {}".format(epoch + 1))
253             for iteration in range(num_iter):
254                 global_step += 1
255                 x_batch,z_batch,train_index=batch(batch_size,hidden[train_index:],label_train[train_index:])
256                 feed_dict_batch={x1:x_batch,z:z_batch,learning_rate:2 / (epoch+1)}
257                 sess.run(optimizer1,feed_dict=feed_dict_batch)
258                 if (iteration + 1) % display_freq == 0:
259                     loss_batch,accuracy_batch = sess.run([loss1,accuracy], feed_dict=feed_dict_batch)
260                     print("Iteration: {}, Loss: {}, Accuracy: {}".format((iteration + 1), loss_batch, accuracy_batch))
261             con_train = [[], [], [], []]
262             con_test = [[], [], [], []]
263             for i in range(n_classes):
264                 con_train[i] = label_train[:, i] == 1
265                 con_test[i] = label_test[:, i] == 1
266                 for j in range(n_classes):
267                     entry_train[i].append(sess.run(ratio, feed_dict={x1: hidden[con_train[i]], indice: [j]}))
268                     entry_test[i].append(sess.run(ratio, feed_dict={x1: hid_test[con_test[i]], indice: [j]}))
269             print(sum(entry_train[i]), sum(entry_test[i]))
270         print("Training set confusion matrix:")
271         print(entry_train)
272         print("Test set confusion matrix:")
273         print(entry_test)

```