

February 5, 2017

## HOMEWORK TWO

### Problem 1. Solution

Results after we re-shuffle the data are shown as Table 1:

Table 1: Results after randomly re-shuffling the data

$\lambda$	Train Accuracy	Validate Accuracy	Test Accuracy
0	0.751225490196	0.757501530925	0.738518064911
0.01	0.748774509804	0.757501530925	0.738518064911
1	0.729166666667	0.753827311696	0.72933251684
100	0.66237745098	0.681567666871	0.680342927128

Listing 1: Key code for Prob.1

```
1 random.seed(0)
2
3 print "Reading data..."
4 dataFile = urllib.urlopen("http://archive.ics.uci.edu/ml/machine-learning-databases/
                             wine-quality/winequality-white.csv")
5 header = dataFile.readline()
6 fields = ["constant"] + header.strip().replace('"', '').split(',')
7 featureNames = fields[:-1]
8 labelName = fields[-1]
9 lines = [[1.0] + [float(x) for x in l.split(',')]] for l in dataFile]
10 random.shuffle(lines)      ##randomly shuffle
11 X = [l[:-1] for l in lines]
12 y = [l[-1] > 5 for l in lines]
13 print "done"
```

### Problem 2. Solution

Classifiers Evaluation Results are shown as Table 2:

Table 2: Classifiers Evaluation Value

True Positives	1129
True Negatives	145
False Positives	321
False Negatives	38
Balanced Error Rate	0.360701663412

Listing 2: Key code for Prob.2

```

1 def performance2(theta):
2
3     scores_test = [inner(theta, x) for x in X_test]
4
5     predictions_test = [s > 0 for s in scores_test]
6     TP = sum([(a == 1 and b == 1) for (a, b) in zip(predictions_test, y_test)])
7     TN = sum([(a == 0 and b == 0) for (a, b) in zip(predictions_test, y_test)])
8     FP = sum([(a == 1 and b == 0) for (a, b) in zip(predictions_test, y_test)])
9     FN = sum([(a == 0 and b == 1) for (a, b) in zip(predictions_test, y_test)])
10    TPR = TP * 1.0 / (TP + FN)
11    TNR = TN * 1.0 / (TN + FP)
12    BER = 1 - (TPR + TNR) / 2
13
14    return TP, TN, FP, FN, BER
15
16
17 lam = 0.01
18 theta = train(lam)
19 TP, TN, FP, FN, BER = performance2(theta)
20 print("TP = " + str(TP) + "; TN=" + str(TN) + "; FP=" + str(FP) + "; FN=" + str(
21 FN) + "; BER=" + str(BER))

```

### Problem 3. Solution

Results are shown as Table 3:

Table 3: Ranking Evaluation Results

Number of Return	Precision	Recall
10	1.0	0.00856898029135
500	0.956	0.409597257926
1000	0.864	0.740359897172

Listing 3: Key code for Prob.3

```
1 def performance3(theta,numR):
2
3     scores_test = [inner(theta, x) for x in X_test]
4     combine = [[scores_test[i]] + [y_test[i]] for i in range(0, len(y_test))]
5     combine.sort(key=lambda x: x[0], reverse=True)
6     total = sum(l[1]==True for l in combine)
7     rel = sum(combine[i][1]==True for i in range(0,numR))
8     precision = rel * 1.0 / numR
9     recall = rel *1.0 / total
10
11     return precision,recall
12
13 lam = 0.01
14 theta = train(lam)
15 for numReturn in [10,500,1000]:
16     precision,recall = performance3(theta,numReturn)
17     print("numReturn = " + str(numReturn) + ";\tprecision=" + str(precision) + ";
                                                recall=" + str(recall))
```

#### Problem 4. Solution

Plot is shown as Fig.1 below:

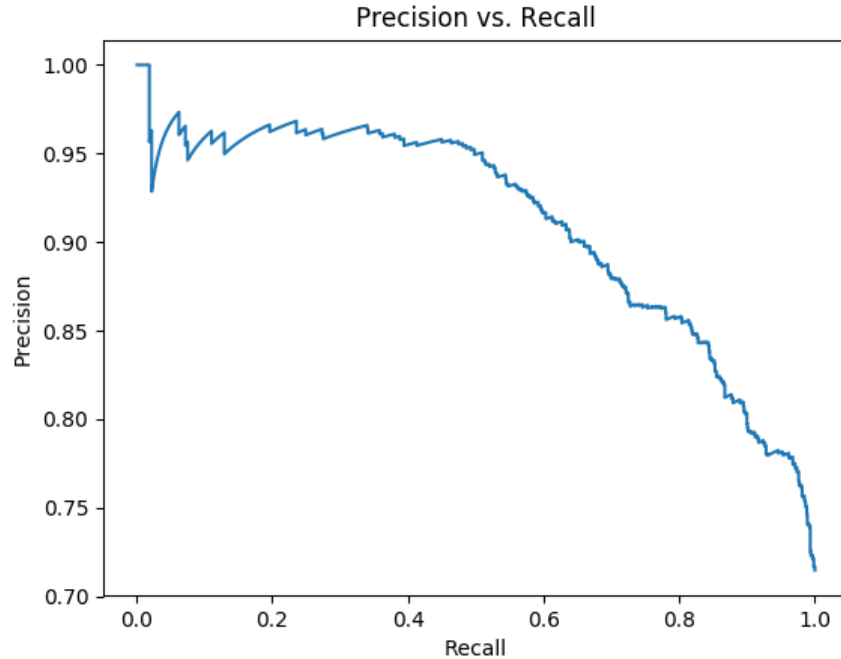


Figure 1: Precision vs. recall as the number of results considered varies

Listing 4: Key code for Prob.4

```
1 def performance3(theta,numR):
2
3     scores_test = [inner(theta, x) for x in X_test]
4     combine = [[scores_test[i]] + [y_test[i]] for i in range(0, len(y_test))]
5     combine.sort(key=lambda x: x[0], reverse=True)
6     total = sum(l[1]==True for l in combine)
7     rel = sum(combine[i][1]==True for i in range(0,numR))
8     precision = rel * 1.0 / numR
9     recall = rel *1.0 / total
10
11     return precision,recall
12
13
14 lam = 0.01
15 theta = train(lam)
16 result_pre = []
17 result_rec = []
18 for numReturn in range(1,len(y_test)+1):
19     precision,recall = performance3(theta,numReturn)
20     result_pre.append(precision)
21     result_rec.append(recall)
22 plt.plot(result_rec,result_pre)
23 plt.xlabel('Recall')
24 plt.ylabel('Precision')
25 plt.title('Precision vs. Recall')
26 plt.show()
```

### Problem 5. Solution

The 'reconstruction error' of the naive 'compression' is: 3675818.61688

Listing 5: Key code for Prob.5

```
1 X_mean=numpy.mean(X_train,axis=0)
2 diff=X_train-X_mean
3 re=sum(sum(numpy.multiply(diff,diff)).T)
4 print re
```

### Problem 6. Solution

PCA components of winequality dataset are shown as Fig.2 below:

```

[[ -3.23636346e-04  1.42201752e-04  3.17030713e-04  5.36390435e-02
   9.30284526e-05  2.54030965e-01  9.65655009e-01  3.19990241e-05
  -2.95831396e-04  3.84043646e-04 -1.00526693e-02]
 [ -7.57985623e-03 -1.66366340e-03  1.04742899e-03  5.21677266e-02
   4.49425600e-05  9.65020304e-01 -2.56793964e-01  7.90089050e-06
   5.24900596e-04 -1.09699394e-03 -2.89827657e-03]
 [  1.82124420e-02  2.54680710e-03  3.31838657e-03  9.93221259e-01
  -1.51888372e-04 -6.42297821e-02 -3.91682592e-02  4.30929482e-04
  -6.93199060e-03 -2.85216045e-03 -8.62920933e-02]
 [  1.56811999e-01  3.28220652e-03  1.66866136e-02  8.28549640e-02
  -6.91822288e-03  1.13029682e-03  5.39110108e-03 -9.49080503e-04
   2.68027305e-03  1.30498102e-03  9.83955205e-01]
 [  9.81360642e-01 -1.45890108e-02  5.92643662e-02 -3.17546064e-02
   5.07483182e-04  8.43759364e-03 -1.77578042e-03  6.03725221e-04
  -9.05011239e-02 -9.35630845e-03 -1.54417839e-01]
 [  7.76578401e-02 -2.37665885e-01  2.23406619e-02  5.04113878e-03
  -1.43564098e-02 -2.14210997e-04 -2.22913844e-04  3.36617054e-03
   8.77254205e-01  4.08570175e-01 -1.54145486e-02]
 [ -7.36289612e-02 -2.61563804e-01  9.43067566e-01 -2.14514264e-03
   1.19104298e-02 -1.68808905e-03  1.42294158e-04 -1.17203197e-04
  -1.45895558e-01  1.23868963e-01 -2.88797236e-03]
 [ -1.37617196e-02  2.11129619e-01 -1.16514121e-01  5.30670319e-04
   1.05181628e-02  1.36446528e-03 -8.21179429e-04  3.09221855e-04
  -3.58358431e-01  9.01728510e-01  3.27758247e-03]
 [  1.74575775e-02  9.10890084e-01  3.04081497e-01 -2.89763923e-03
   2.34615054e-02  1.17406025e-03 -3.85957239e-04  1.23176271e-03
   2.68927937e-01 -6.70756658e-02 -1.12101920e-02]
 [  2.31513441e-03 -2.38717789e-02 -1.67445603e-02  8.92206499e-04
   9.99462734e-01 -9.81109101e-05 -3.32812875e-05  4.14235255e-03
   1.18483756e-02 -3.51543098e-03  6.92344110e-03]
 [  7.48312160e-04  3.08204153e-04  2.55232500e-04  3.49846801e-04
   4.12943179e-03 -6.96565372e-06  4.16951216e-06 -9.99984215e-01
   3.17948604e-03  1.53436134e-03 -1.10029138e-03]]

```

Figure 2: PCA components of winequality dataset

Listing 6: Key code for Prob.6

```

1 print "Reading data..."
2 dataFile = urllib.urlopen("http://archive.ics.uci.edu/ml/machine-learning-databases/
   wine-quality/winequality-white.csv")
3 header = dataFile.readline()
4 fields = ["constant"] + header.strip().replace('"', '').split(';')
5 featureNames = fields[:-1]
6 labelName = fields[-1]
7 lines = [[float(x) for x in l.split(';')] for l in dataFile]

```

```

8 X = [l[:-1] for l in lines]
9 y = [l[-1] > 5 for l in lines]
10 print "done"
11
12 X_train = X[:int(len(X) / 3)]
13 y_train = y[:int(len(y) / 3)]
14 X_validate = X[int(len(X) / 3):int(2 * len(X) / 3)]
15 y_validate = y[int(len(y) / 3):int(2 * len(y) / 3)]
16 X_test = X[int(2 * len(X) / 3):]
17 y_test = y[int(2 * len(X) / 3):]
18
19 pca = PCA(n_components=11)
20 pca.fit(X_train)
21 print pca.components_

```

### Problem 7. Solution

The reconstruction error is: 1345.47557409

We can use two ways to compute the reconstruction error, one is to use the projection matrix we get to reconstruct original data and then compute the error, the other is use the sum of the variance we drop to compute the reconstruction error.

Listing 7: Key code for Prob.7

```

1  pca = PCA(n_components=4)
2  pca.fit(X_train)
3
4
5  X_PCA=pca.fit_transform(X_train)      ##Way 1
6  x_rec2=pca.inverse_transform(X_PCA)
7  x_rec2=numpy.matrix(x_rec2)
8  error2=sum(sum(numpy.multiply(x_rec2-X_train,x_rec2-X_train)).T).tolist()[0][0]
9  print error2
10
11 print pca.noise_variance_*len(x_rec2)  ##Way 2

```

### Problem 8. Solution

Plot is shown as Fig.3 below. We can see from the plot that both training and test set MSE decrease with the increment of the number of dimensions in general. And the test set sometimes gets higher MSE with a higher number of dimensions.

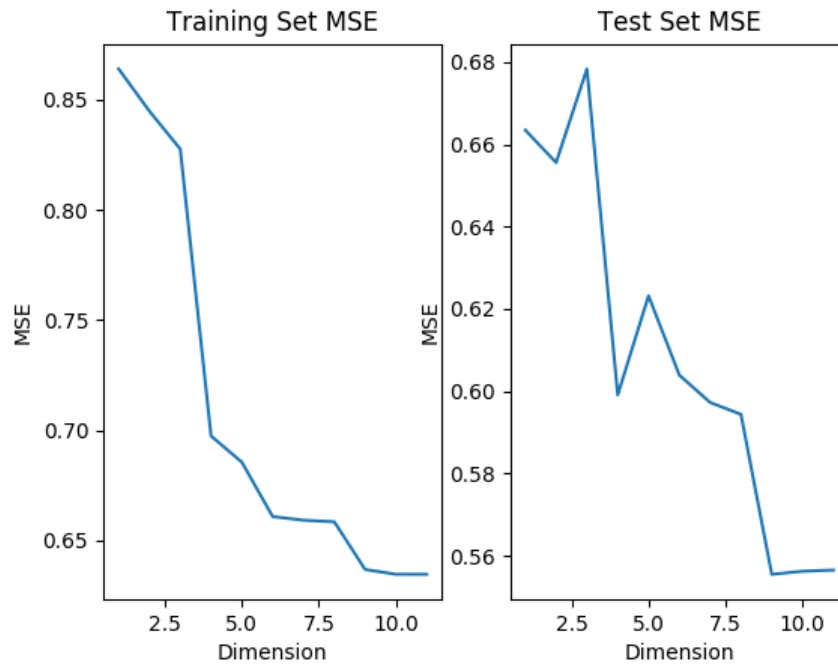


Figure 3: MSE vs. dimensions on the train and test sets

Listing 8: Key code for Prob.8

```

1 train_mse=[]
2 test_mse=[]
3
4 for dim in range(1,12):
5     pca = PCA(n_components=dim)
6     X_PCA = pca.fit_transform(X_train).tolist()
7     X_PCA = [[1] + elem for elem in X_PCA]
8     theta, l, info = scipy.optimize.fmin_l_bfgs_b(f, [0]*(dim+1), fprime, args = (X_PCA
9                                                    , y_train, 0)) #lambda
10    #print "Final log likelihood =", -l
11    X_test_PCA = pca.transform(X_test).tolist()
12    X_test_PCA = [[1] + elem for elem in X_test_PCA]
13    X_test_PCA=numpy.matrix(X_test_PCA)
14    theta=numpy.matrix(theta)
15    y_test=numpy.matrix(y_test)
16    y=X_test_PCA*theta.T
17    k=y-y_test.T
18    mse2=(k.T*k/len(y)).tolist()[0][0]
19    test_mse.append(mse2)
20
21    X_PCA = numpy.matrix(X_PCA)
22    y_train = numpy.matrix(y_train)

```

```

22     y = X_PCA * theta.T
23     k = y - y_train.T
24     mse1 = (k.T * k / len(y)).tolist()[0][0]
25     train_mse.append(mse1)
26 XX=range(1,12)
27 plt.figure()
28 ax1 = plt.subplot(121)
29 ax2 = plt.subplot(122)
30 plt.sca(ax1)
31 plt.plot(XX,train_mse)
32 plt.title('Training Set MSE')
33 plt.xlabel('Dimension')
34 plt.ylabel('MSE')
35 plt.sca(ax2)
36 plt.plot(XX,test_mse)
37 plt.xlabel('Dimension')
38 plt.ylabel('MSE')
39 plt.title('Test Set MSE')
40 plt.show()

```

Submitted by Zhang, Jinhan PID A53211930 on February 5, 2017.