

January 22, 2017

## HOMEWORK ONE

### Problem 1. Solution

$$\theta_0 = -39.1707489, \theta_1 = 0.0214379786$$

Listing 1: Key code for Prob.1

```
1 def feature(datum):
2     feat=[1]
3     feat.append(datum['review/timeStruct']['year'])
4     return feat
5
6 X = [feature(d) for d in data]
7 y = [d['review/overall'] for d in data]
8 theta,residuals,rank,s = numpy.linalg.lstsq(X,y)
9 print theta
```

### Problem 2. Solution

Since the range of the years in dataset is all between 1999 and 2012, we can assume that all year labels are between 1999 and 2012, thus we can use 13 boolean variables to represent year variable(since if the states of all 13 years are False, it must be the remaining year, so we can use 13 extra variable to represent 14 years.) The equation for that is:

$$review/overall \simeq \theta_0 + \theta_1 \times is1999 + \theta_2 \times is2000 + \theta_3 \times is2001 + \theta_4 \times is2002 + \theta_5 \times is2003 + \theta_6 \times is2004 + \theta_7 \times is2005 + \theta_8 \times is2006 + \theta_9 \times is2007 + \theta_{10} \times is2008 + \theta_{11} \times is2009 + \theta_{12} \times is2010 + \theta_{13} \times is2011$$

The MSE for the model in Problem 1 is: 0.49004382

The MSE for our model is: 0.4891519

We see that MSE has decreased when compared to that in Problem 1.

Listing 2: Key code for Prob.2

```
1 MSE1=residuals/len(y) #residuals,X and y are what we get from code in Listing 1
2 print MSE1
3 tmp=numpy.matrix(X)
4 minV=tmp[:,1].min()
5 maxV=tmp[:,1].max()
```

```

7 def feature2(datum):
8     feat=[1]
9     tmp=datum['review/timeStruct']['year']
10    l=[0 for x in range(13)]
11    if tmp!=maxV:
12        l[tmp-minV]=1
13    feat.extend(l)
14    return feat
15
16 X = [feature2(d) for d in data]
17 y = [d['review/overall'] for d in data]
18
19 theta,residuals,rank,s = numpy.linalg.lstsq(X,y)
20 MSE2=residuals/len(y)
21 print MSE2

```

### Problem 3. Solution

Fitted Coefficients:

$\theta_0 = 256.420279$   
 $\theta_1 = 0.135421303$   
 $\theta_2 = -1.72994866$   
 $\theta_3 = 0.102651152$   
 $\theta_4 = 0.109038568$   
 $\theta_5 = -0.276775146$   
 $\theta_6 = 0.00634332168$   
 $\theta_7 = 0.00003.85023977$   
 $\theta_8 = -258.652809$   
 $\theta_9 = 1.19540566$   
 $\theta_{10} = 0.833006285$   
 $\theta_{11} = 0.09.79304353$

MSE on the train data : 0.602307502903

MSE on the test data : 0.562457130315

The code below stores data in dictionary, so the key field is sorted alphabetically in the code, but I've mapped the coefficients to the original order as the assignment required manually. The same applies for all of the rest.

Listing 3: Key code for Prob.3

```

1 l=len(data)      #split the data into 2 groups
2 train=data[0:l/2]
3 test=data[l/2:]
4
5 keys=["alcohol","chlorides","citric acid","density","fixed acidity",
6       "free sulfur dioxide","pH","residual sugar","sulphates",
7       "total sulfur dioxide","volatile acidity"]
8 def feature(datum):
9     feat=[1]

```

```

10     for key in keys:
11         feat.append(eval(datum[key]))
12     return feat
13 X=[feature(d) for d in train]
14 y=[eval(d["quality"]) for d in train]
15 theta,residuals,rank,s = numpy.linalg.lstsq(X,y)
16 print "Theta:",theta
17 print "Training Set MSE:",residuals[0]/len(train)
18
19 X_t=[feature(d) for d in test]
20 y_t=[eval(d["quality"]) for d in test]
21
22 X_t=numpy.matrix(X_t)
23 theta=numpy.matrix(theta)
24 y_t=numpy.matrix(y_t)
25 diff=theta*X_t.T-y_t
26 mse=diff*diff.T
27 k = mse.tolist()[0]
28 print "Test Set MSE:", k[0] / len(test)

```

#### Problem 4. Solution

- (a) MSEs (on the test set) of all 11 ablation experiments:

<i>Removed_Field : MSE</i>	
"fixed acidity"	: 0.559113414376
"volatile acidity"	: 0.596384850161
"citric acid"	: 0.562221702812
"residual sugar"	: 0.553625063967
"chlorides"	: 0.562629266481
"free sulfur dioxide"	: 0.55614081793
"total sulfur dioxide"	: 0.562429005469
"density"	: 0.544726553466
"pH"	: 0.559566626382
"sulphates"	: 0.557346349988
"alcohol"	: 0.573214743558

- (b) Conclusions:

From the result above, we see the MSE is highest when we remove "volatile acidity" field, so we can conclude that "volatile acidity" feature provides the most additional information for the data. And the MSE is lowest when we remove "density" field, which means "density" feature provides the least additional information.

Listing 4: Key code for Prob.4

```

1 l=len(data)      #split the data into 2 groups
2 train=data[0:l/2]

```

```

3 test=data[1/2:]
4
5 keys=["alcohol","chlorides","citric acid","density","fixed acidity",
6 "free sulfur dioxide","pH","residual sugar","sulphates",
7 "total sulfur dioxide","volatile acidity"]
8 def feature(datum):
9     feat=[1]
10    for key in keys:
11        feat.append(eval(datum[key]))
12    return feat
13 X=[feature(d) for d in train]
14 y=[eval(d["quality"]) for d in train]
15
16 for i in range(1,12):
17     X_ab = [x[:i] + x[i+1:] for x in X]
18     theta, residuals, rank, s = numpy.linalg.lstsq(X_ab, y)
19     print "Training Set MSE:",residuals[0]/len(train)
20     X_t = [feature(d) for d in test]
21     y_t = [eval(d["quality"]) for d in test]
22
23     X_ab_t = [x[:i] + x[i+1:] for x in X_t]
24     X_ab_t = numpy.matrix(X_ab_t)
25     theta = numpy.matrix(theta)
26     y_t = numpy.matrix(y_t)
27     diff = theta * X_ab_t.T - y_t
28     mse = diff * diff.T
29     k=mse.tolist()[0]
30     print "Test Set MSE:",k[0]/len(test)
31     print

```

### Problem 5. Solution

When we use the default "penalty parameter" (C=1000), the train accuracy is 100% , while the test accuracy is 66.803%. This result is mainly because we have a large "penalty parameter", which means we can hardly tolerate any mis-classification in training dataset. And that leads to overfitting, so we have a relatively low accuracy in test dataset. And when I change the C value into 0.8, the train accuracy decreases to 90.527% while the test accuracy increases to 69.906%, which means the overfitting has been weakened.

Listing 5: Key code for Prob.5

```

1 keys=["alcohol","chlorides","citric acid","density","fixed acidity",
2 "free sulfur dioxide","pH","residual sugar","sulphates",
3 "total sulfur dioxide","volatile acidity"]
4 def feature(datum):
5     feat=[]
6     for key in keys:
7         feat.append(eval(datum[key]))
8     return feat
9 X=[feature(d) for d in data]
10 y=[eval(d["quality"])>5 for d in data]
11
12 l=len(data)      #split the data into 2 groups
13 X_train = X[:l/2]

```

```

14 y_train = y[:1/2]
15
16 X_test = X[1/2:]
17 y_test = y[1/2:]
18
19 clf = svm.SVC(C=1000)
20 clf.fit(X_train, y_train)
21
22 train_predictions = clf.predict(X_train)
23 test_predictions = clf.predict(X_test)
24
25 accuracy_train = sum([z[0]==z[1] for z in zip(train_predictions,
26 y_train)])*1.0/len(train_predictions)
27 accuracy_test = sum([z[0]==z[1] for z in zip(test_predictions,
28 y_test)])*1.0/len(test_predictions)
29
30
31 print "The train accuracy is: ",accuracy_train
32 print "The test accuracy is: ",accuracy_test

```

### Problem 6. Solution

- (a) Code stub for the derivative (fprime):

```

1 def fprime(theta, X, y, lam):
2     dl = [0.0]*len(theta)
3     for j in range(len(theta)):
4         for i in range(len(X)):
5             logit = inner(X[i], theta)
6             dl[j] += X[i][j] * (1 - sigmoid(logit))
7             if not y[i]:
8                 dl[j] -= X[i][j]
9             dl[j] -= 2 * lam * theta[j]
10            # Negate the return value since we're doing gradient *ascent*
11            # print "dl =", dl
12    return numpy.array([-x for x in dl])

```

- (b) Log-likelihood and accuracy

The log-likelihood after convergence is  $-1388.69589879$

The accuracy is 76.929%

*Submitted by Zhang, Jinhan PID A53211930 on January 22, 2017.*