# Data Mining Assignment 2

Zhengyuan Zhu

CSE5334 - Data Mining

UNIVERSITY OF TEXAS AT ARLINGTON

April 4, 2020

## Problem 1: Non-parameteric density estimation

### Solution for 1.1

First, we set up all the parameters for generating the Gaussian random data where the $\mu_1 = 5$ and $\sigma_1 = 1$

```
data_5_1 = np.random.normal(5, 1, 1000)
```

Then, we write a function p = mykde(X, h, x) that performs kernel density estimation on data X with bandwidth h for 1-Dimensional data. The function take $X$ as the samples that calculate the density of gaussian distribution on $x$. It traverses all the samples and add each sample's density one by one.

```
def mykde(X, h, x):
    density = np.zeros(len(X))
    for i in range(len(X)):
        density += np.exp(-(x-X[i])**2/(2*h**2))/(h*np.sqrt(2*np.pi))/len(
                                                  X)
    return density
```

The figure below is the visualization of the data distribution:
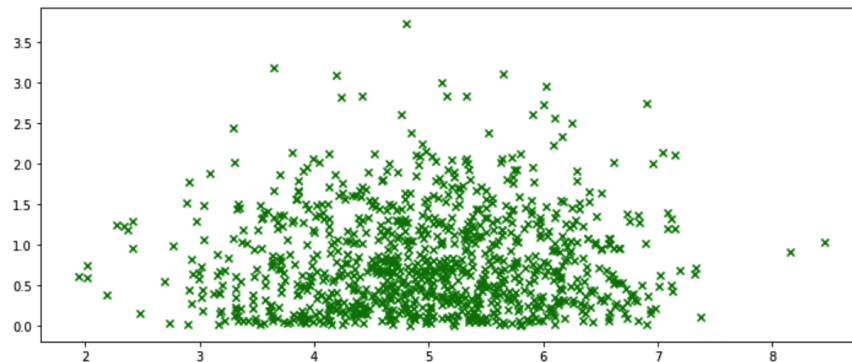


Figure 1: Fake data with $\mu_1, \sigma_1$

Then we do some visualizations of the probability density function. The figure below shows that the different densities with different bandwidth. The red line represents the ground truth distribution and the histogram represents the samples distribution. It shows that the distribution curve becomes more flat when bandwidths get larger and larger.
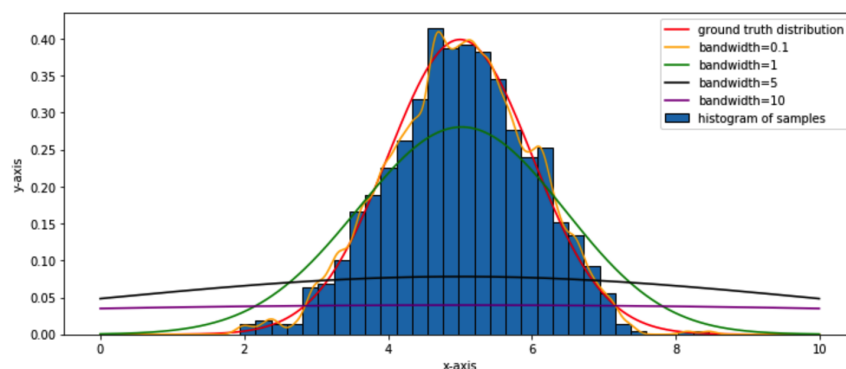


Figure 2: Histogram of X along with the different estimated probability density function

## Solution for 1.2

For the problem 1.2, the methodology essentially is the same as problem 1.1. The only different is that we add a new random normal distribution data with $\mu_2 = 0$ and $\sigma = 0, 2$. We use the codes below for the visualization of results. And the conclusion is pretty much the same as problem 1.1. The larger bandwidth we got, the more flat probability density function in the figure.

```
fig = plt.figure(figsize=(12, 5))
ax = fig.add_subplot(111)
plt.hist(data_5_1, bins=30, edgecolor='black', density=True, label='
                                    histogram of samples')
# Plot kernel function
x = np.arange(0, 10, 0.01)
y = normal(x, 5, 1)
d1 = mykde(data_5_1, 0.1, x)
d2 = mykde(data_5_1, 1, x)
d3 = mykde(data_5_1, 5, x)
d4 = mykde(data_5_1, 10, x)
# Plot the ground truth distribution and estimations
plt.plot(x, y, color='red', label='ground truth distribution')
ax.plot(x, d1, label='bandwidth=0.1', color='orange')
ax.plot(x, d2, label='bandwidth=1', color='green')
ax.plot(x, d3, label='bandwidth=5', color='black')
ax.plot(x, d4, label='bandwidth=10', color='purple')

plt.legend()
plt.xlabel('x-axis')
plt.ylabel('y-axis')
```
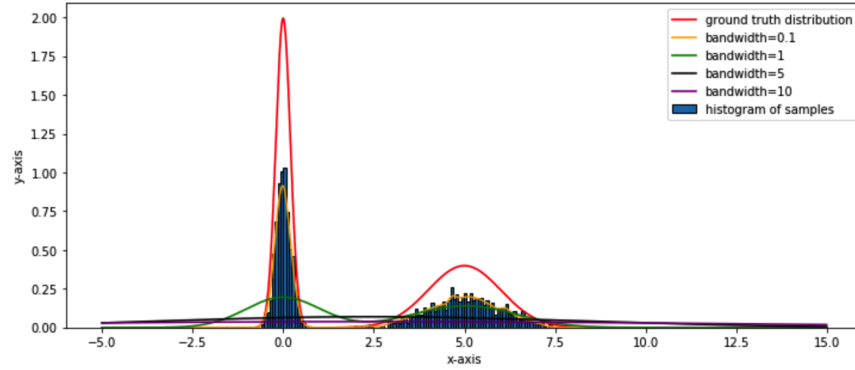
2

Figure 3: Histogram of X along with the different estimated probability density function

## Solution for 1.3

When we have 2 sets of 2-D Gaussian random data, we need to modify mykde function that can fit the 2-D gaussian distribution. I have tried a lot but failed, so that I decide to modify the sklearn vanilla Kernel Density and stack the data we generated to a grid form. The source codes of mykde2D is simple and easy to understand.

```python
from sklearn.neighbors import KernelDensity

def mykde2D(x, y, bandwidth, xbins=100j, ybins=100j, **kwargs):
    # create grid of sample locations (default: 100x100)
    xx, yy = np.mgrid[x.min():x.max():xbins,
                      y.min():y.max():ybins]

    xy_sample = np.vstack([yy.ravel(), xx.ravel()]).T
    xy_train  = np.vstack([y, x]).T

    kde_skl = KernelDensity(bandwidth=bandwidth, **kwargs)
    kde_skl.fit(xy_train)

    # score_samples() returns the log-likelihood of the samples
    z = np.exp(kde_skl.score_samples(xy_sample))
    return xx, yy, np.reshape(z, xx.shape)
```

The visualization of 2D gaussian distribution is a little bit tricky. The white dots represent the 2D gaussian distribution's samples. And I use a colored mesh to create a pseudocolor plot with a non-regular rectangular grid. The shade that cover the figure shows the $\mu$ and $\sigma$ as results of mykde2D. The results shows that when bandwidth is small(such as 0.1), the estimated distribution seems like to be "overfitting". And vice versa the estimated distribution become "underfitting" when bandwidth is large.
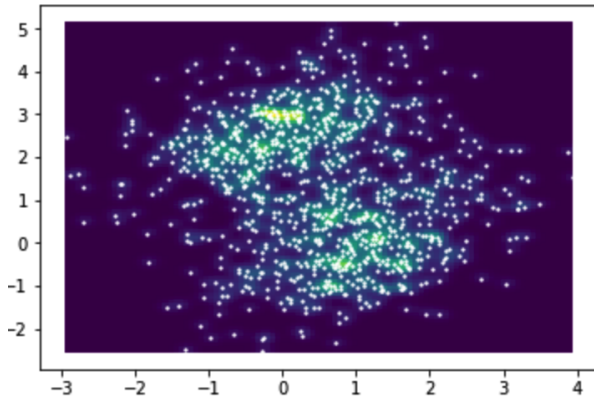
3

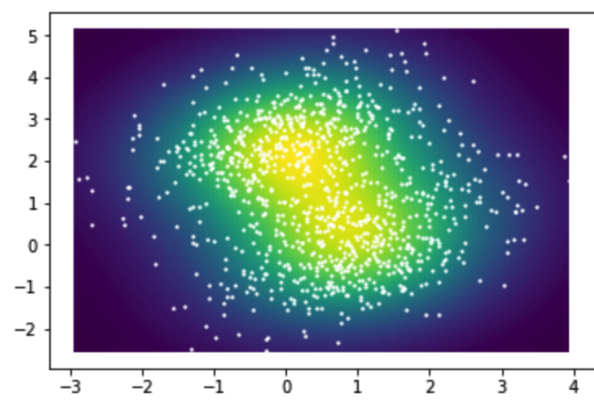Figure 4: Test your function mykde on this data with h = 0.1



Figure 5: Test your function mykde on this data with h = 1



Figure 6: Test your function mykde on this data with h = 5



Figure 7: Test your function mykde on this data with h = 10

# Problem 2 Solution

## Solution 2.1

To solve the problem, the first thing is to generate 1000 training instances in two different classes (500 in each) from multi-variate normal distribution using the following parameters for each class. The source codes are below with samples.

```
mul1, sigma1 = [1, 0], [[1, 0.75], [0.75, 1]]
mul2, sigma2 = [0, 1], [[1, 0.75], [0.75, 1]]
size = 500

train_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size
                                                =size)
train_data2D_1_label = np.zeros((size, 1))
print(train_data2D_1_label.shape)

```

4

```python
train_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size
                                                =size)
train_data2D_2_label = np.ones((size, 1))
X_train = np.vstack([train_data2D_1, train_data2D_2])
y_train = np.vstack([train_data2D_1_label, train_data2D_2_label])
print(X_train[:5], X_train[-5:])
print(y_train[:5], y_train[-5:])

# build test set
mul1, sigma1 = [1, 0], [[1, 0.75], [0.75, 1]]
mul2, sigma2 = [0, 1], [[1, 0.75], [0.75, 1]]
size = 500

test_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size=
                                              size)
test_data2D_1_label = np.zeros((size, 1))
print(test_data2D_1_label.shape)

test_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size=
                                              size)
test_data2D_2_label = np.ones((size, 1))
X_test = np.vstack([test_data2D_1, test_data2D_2])
y_test = np.vstack([test_data2D_1_label, test_data2D_2_label])
print(X_test[:5], X_test[-5:])
print(y_test[:5], y_test[-5:])
---------------------------------------------------------------------------
(500, 1)
[[-0.19002878 -0.14795651]
 [-0.27783423 -1.46171299]
 [ 0.78327396  0.39920896]
 [ 0.57195261 -0.04425206]
 [ 1.0097563  -1.48342926]] [[ 1.0917887   2.14877639]
 [ 0.20435192  1.65422036]
 [-1.49104023  0.3028843 ]
 [-0.48226042  0.62071859]
 [-0.26311759 -0.31374102]]
[[0.]
 [0.]
 [0.]
 [0.]
 [0.]] [[1.]
 [1.]
 [1.]
 [1.]
 [1.]]

 (500, 1)
[[ 1.54609753 -0.01456578]
 [ 1.49064799  0.83087897]
 [-0.59881474 -1.20925996]
 [ 0.34422204 -0.58984347]
 [ 0.42700309 -1.39169054]] [[-0.21408786  0.96185428]
 [ 0.56305549  1.13239575]
 [ 2.30235643  3.27293605]
```

```
60    [ 1.5609519    1.07252668]
61    [-1.14839655   0.92135342]]
62 [[0.]
63    [0.]
64    [0.]
65    [0.]
66    [0.]] [[1.]
67    [1.]
68    [1.]
69    [1.]
70    [1.]]
```

Then I implement your Naive Bayes Classifier [**pred, posterior, err**] = **myNB(X,Y,X test,Y test)** whose inputs are the training data **X**, labels **Y** for **X**, testing data **X_test** and labels **Y_test** for **X_test** and returns predicted labels **pred**, posterior probability **posterior** with which the prediction was made and error rate **err**.

```python
1  import pandas as pd
2  import numpy as np
3  print(pd.__version__)
4
5  def myNB(X, y, X_test, y_test):
6      def norm(x, mean, std):
7          return 1/(np.sqrt(2*np.pi)*std)*np.exp(-(x-mean)**2/2*(std**2))
8
9      X_train, y_train = pd.DataFrame(X, columns=['Feature1', 'Feature2']),
                                       pd.DataFrame(y, columns=['Class']
                                       )
10     X_test, y_test = pd.DataFrame(X_test, columns=['Feature1', 'Feature2']
                                     ), pd.DataFrame(y_test, columns=[
                                     'Class'])
11     y_unique = y_train['Class'].unique()
12     train_set = pd.concat([X_train, y_train], axis=1)
13     prior = np.zeros(len(y_unique))
14     conditional = np.reshape(np.zeros(len(y_unique)*len(X_train.columns)*2
                               ), (len(y_unique), len(X_train.
                               columns), 2))
15
16
17     for i in range(0,len(y_unique)):
18         prior[i]=(sum(y_train['Class']==y_unique[i])+1)/(len(y_train['
                                                 Class'])+len(y_unique))
19 #     print("The prior probability of each class is: ", prior)
20
21     for i, h in enumerate(X_train.columns.values.tolist()):
22         for j in range(0,len(y_unique)):
23             class_feature = train_set[h].loc[(train_set['Class']==y_unique
                                               [j])]
24             mean = np.mean(class_feature)
25             var = np.std(class_feature)
26             conditional[i][j] = [mean, var]
27 #             print("mean and standard variance of current feature is: ",
                                               h, j, mean, var)
28 #     print(conditional)
```

6

```python
29 #     print("Prior distribution is: ", prior)
30     pred_probs = []
31     for idx, row in X_test.iterrows():
32         probs = []
33         for cIdx, pri in enumerate(prior): # class 0 1
34             for fIdx, feat in enumerate(row): # feature 0 1
35                 pri *= norm(feat, conditional[fIdx][cIdx][0], conditional[
                                                 fIdx][cIdx][1])
36             probs.append(pri)
37         pred_probs.append(probs)
38     pred = [np.argmax(p) for p in pred_probs]
39     # calculate error rate
40 #     print(np.array(pred).shape)
41 #     print(y_test.to_numpy().shape)
42     err = np.mean(np.array(pred)!=y_test.to_numpy().squeeze())
43 #     print(err)
44
45     return pred, conditional, err
46
47 pred, posterior, err = myNB(X_train, y_train, X_test, y_test)
```

For the experiments part, I perform the experiments 10 times and take an average of error rate. The source codes and results are shown below.

```python
1 # perform the experiments 10 times
2 avg_err, run_time = 0, 10
3 for i in range(run_time):
4     train_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1,
                                             size=size)
5     train_data2D_1_label = np.zeros((size, 1))
6
7     train_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2,
                                             size=size)
8     train_data2D_2_label = np.ones((size, 1))
9     X_train = np.vstack([train_data2D_1, train_data2D_2])
10    y_train = np.vstack([train_data2D_1_label, train_data2D_2_label])
11
12    test_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1,
                                             size=size)
13    test_data2D_1_label = np.zeros((size, 1))
14
15    test_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2,
                                             size=size)
16    test_data2D_2_label = np.ones((size, 1))
17    X_test = np.vstack([test_data2D_1, test_data2D_2])
18    y_test = np.vstack([test_data2D_1_label, test_data2D_2_label])
19
20    pred, posterior, err = myNB(X_train, y_train, X_test, y_test)
21    print("current time error: ", err)
22    avg_err += err/run_time
23 print("The average error: ", avg_err)
24 --------------------------------------------------------------------
25 current time error:  0.117
26 current time error:  0.057
```

```
27 current time error:   0.101
28 current time error:   0.061
29 current time error:   0.086
30 current time error:   0.08
31 current time error:   0.079
32 current time error:   0.08
33 current time error:   0.071
34 current time error:   0.09
35 The average error:   0.08219999999999998
```

I have performed prediction on the testing data with source code. The accuracy, precision, recall, confusion matrix as well as a scatter plot of data points whose labels are color coded are shown below.

```python
1  from sklearn import metrics
2  import matplotlib.pyplot as plt
3
4  train_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size
                                                    =size)
5  train_data2D_1_label = np.zeros((size, 1))
6
7  train_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size
                                                    =size)
8  train_data2D_2_label = np.ones((size, 1))
9  X_train = np.vstack([train_data2D_1, train_data2D_2])
10 y_train = np.vstack([train_data2D_1_label, train_data2D_2_label])
11
12 test_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size=
                                                   size)
13 test_data2D_1_label = np.zeros((size, 1))
14
15 test_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size=
                                                   size)
16 test_data2D_2_label = np.ones((size, 1))
17 X_test = np.vstack([test_data2D_1, test_data2D_2])
18 y_test = np.vstack([test_data2D_1_label, test_data2D_2_label])
19
20 pred, posterior, err = myNB(X_train, y_train, X_test, y_test)
21
22 # report the metrics
23 acc = metrics.accuracy_score(np.array(pred), y_test.squeeze())
24 print("Accuracy = %.3f" % (sum([p==a for p, a in zip(np.array(pred),
                                    y_test.squeeze())])/len(np.array(pred
                                    ))))
25 precision = metrics.precision_score(np.array(pred), y_test.squeeze())
26 print("Precision = %.3f" % precision)
27 recall = metrics.recall_score(np.array(pred), y_test.squeeze())
28 print("Metrics = %.3f" % recall)
29 cm = metrics.confusion_matrix(np.array(pred), y_test.squeeze())
30 print("Confusion Matrix = {}".format(cm))
31
32 # draw data points
33 plt.scatter(train_data2D_1[:, 0], train_data2D_1[:, 1], label='$\mu_1, \
                                    sigma_1$')
```

```
34  plt.scatter(train_data2D_2[:, 0], train_data2D_2[:, 1], label='$\mu_2, \
                                       sigma_2$')
35  plt.legend()
36  plt.show()
37  ------------------------------------------------------------------------
38  Accuracy = 0.912
39  Precision = 0.896
40  Metrics = 0.926
41  Confusion Matrix = [[464   52]
42                      [ 36 448]]
```
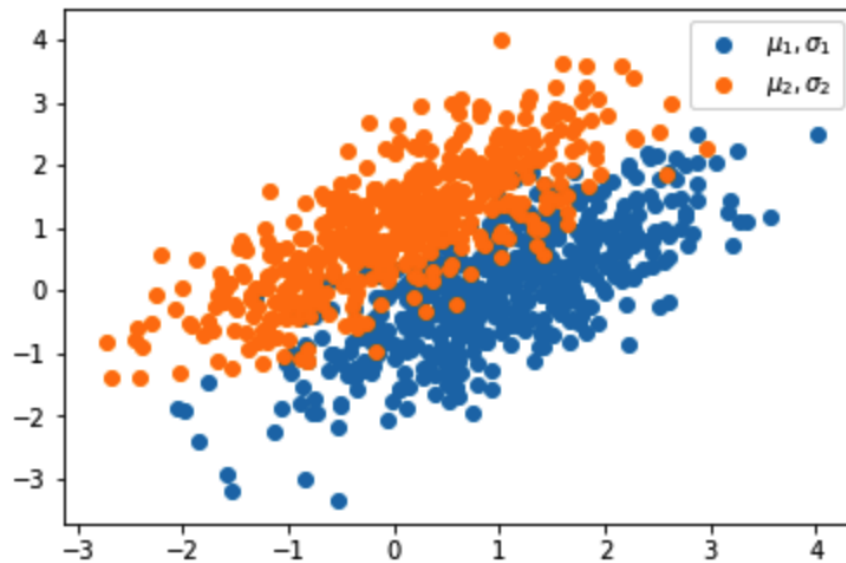


Figure 8: Dataset distribution with two differnet colors

And in my training data, I have changed the number of examples in each class to 10,20,50,100,300,500 and perform prediction on the testing data with my code. The plot of changes of accuracies are shown below. The result shows that the more exmaples we have, more higher accuracy we obtain.

```
1   # change the number of examples
2   accs = []
3   x_axis = [0, 1, 2, 3, 4, 5]
4   for n in [10,20,50,100,300,500]:
5       train_set_1 = train_data2D_1[:n]
6       train_label1 = np.zeros((n, 1))
7       train_set_2 = train_data2D_2[:n]
8       train_label2 = np.ones((n, 1))
9       X_train = np.vstack([train_set_1, train_set_2])
10      y_train = np.vstack([train_label1, train_label2])
11      pred, posterior, err = myNB(X_train, y_train, X_test, y_test)
12      acc = metrics.accuracy_score(np.array(pred), y_test.squeeze())
13      accs.append(acc)
```

```
14
15
16  plt.plot(x_axis, accs, 'ro-', label='Changes of accuracies')
17  plt.xlabel("number of examples")
18  plt.ylabel("accuracy")
19  plt.legend()
20  plt.show()
```
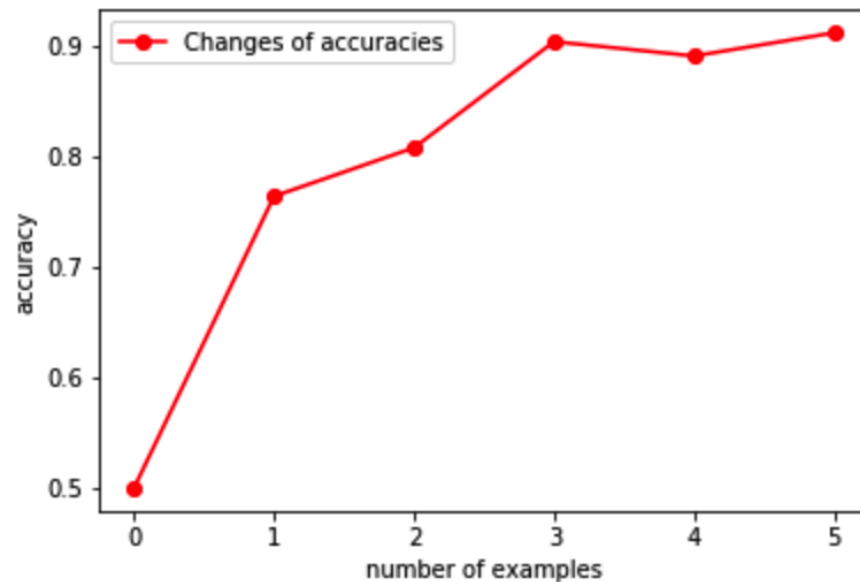


Figure 9: The changes of accuracies along with differnet number of examples

And when in the training data, change the number of examples in class 0 as 700 and the other as 300. Perform prediction on the testing dataset, I found that the accuracy becomes way much lower because of the unbalanced data.

```
1  # change the number of examples
2  train_data2D_1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size
                                                =700)
3  train_data2D_1_label = np.zeros((700, 1))
4
5  train_data2D_2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size
                                                =300)
6  train_data2D_2_label = np.ones((300, 1))
7
8  train_set_1 = train_data2D_1[:700]
9  train_label1 = np.zeros((700, 1))
10 train_set_2 = train_data2D_2[:300]
11 train_label2 = np.ones((300, 1))
12 X_train = np.vstack([train_set_1, train_set_2])
13 y_train = np.vstack([train_label1, train_label2])
14 pred, posterior, err = myNB(X_train, y_train, X_test, y_test)
15 acc = metrics.accuracy_score(np.array(pred), y_test.squeeze())
```

```
16  print("The accuracy when   change  the  number  of  examples  in  class  0  as  700
                                    and  the  other  as  300:", acc)
17  -----------------------------------------------------------------
18  The  accuracy  when   change  the  number  of  examples  in class  0  as  700  and  the
                                    other  as  300: 0.834
```
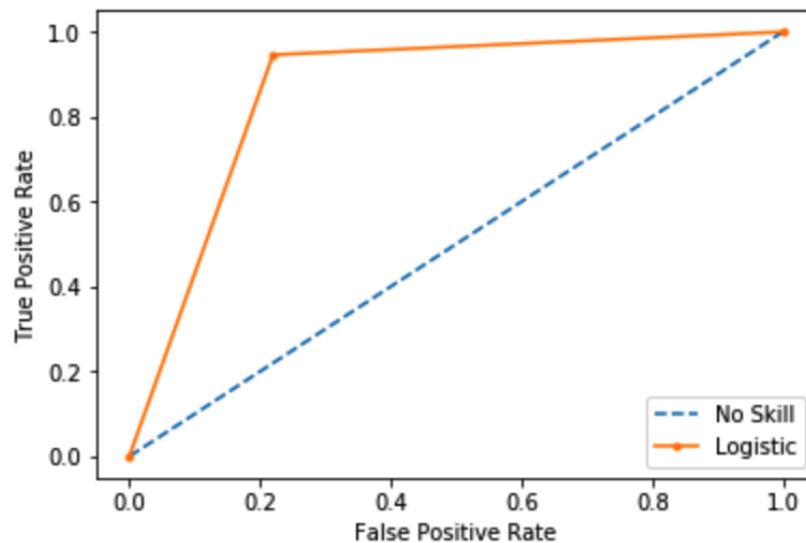
Then the ROC and AUC curve are shown below:



Figure 10: ROC curves from the dataset above

In the last part, I reproduce the tfidf matrix in the homework1 and print some of the reviews.

```
1   # Generate tfidf matrix
2   !pip install nltk
3   import nltk
4   nltk.download('wordnet')
5   nltk.download('stopwords')
6   import pandas as pd
7
8   review=pd.read_csv('./data/Amazon_Reviews.csv')
9   review['Label']=review['Label'].map({'__label__2 ':1,'__label__1 ':0})
10  docs_num = review.shape[0]
11
12  # Preprocessing for tokenize, stopwords removal, lemmatization
13  from nltk.tokenize import RegexpTokenizer
14  from nltk.stem import WordNetLemmatizer
15  from nltk.corpus import stopwords
16
17  tokenizer=RegexpTokenizer(r'\w+')
18  lemmatizer=WordNetLemmatizer()
19
```

```
20  for idx, row in review.iterrows():
21      tmp = []
22      tokens = tokenizer.tokenize(row['Review'])
23      for w in tokens:
24          if w.lower() not in stopwords.words('english'):
25              tmp.append(lemmatizer.lemmatize(w.lower()))
26      review.iloc[idx, 0] = ' '.join(tmp)
27
28  DF = {}
29  for idx, row in review.iterrows():
30      for w in row['Review'].split():
31          try:
32              DF[w].add(idx)
33          except:
34              DF[w] = {idx}
35  DF = {k: len(v) for k, v in DF.items()}
36  vocab = [k for k, v in DF.items()]
37  vocab_len = len(vocab)
38
39  from collections import Counter
40  import numpy as np
41
42  TFIDF = {}
43  for idx, row in review.iterrows():
44      rw = row['Review'].split()
45      cnt = Counter(rw)
46      for w in np.unique(rw):
47          tf = cnt[w]/len(rw)
48          df = DF[w]
49          idf = np.log((docs_num+1)/(df+1))
50          TFIDF[(idx, w)] = tf*idf
51
52  n = 0
53  for k, v in TFIDF.items():
54      if n<5:
55          print(k, v)
56          n += 1
57      else:
58          break
59
60  """
61  P: like(54) great(38) good(35) love(33) best(20)
62  N: however(15) disappointed(11) waste(10) poor(10) hard(8)
63  """
64  def gen_vector(idx, review):
65      Q = np.zeros(10)
66      P = np.zeros(10)
67      ten_words = ['like', 'great', 'good', 'love', 'best',
68                   'however', 'disappointed', 'waste', 'poor', 'hard']
69      for w in ten_words:
70          for k, v in TFIDF.items():
71              if k[0]==idx and k[1]==w:
72                  Q[ten_words.index(w)] = v
73          for word in review:
```

```
74              if word==w:
75                  P[ten_words.index(w)] += 1
76
77      return P, Q
78
79 tfidf_vectors = []
80 count_vectors = []
81 for idx, row in review.iterrows():
82      rw = row['Review'].split()
83      P, Q = gen_vector(idx, rw)
84      tfidf_vectors.append([Q, row['Label']])
85      count_vectors.append(P)
86
87 for i in range(10):
88      print("tf-idf vector for document {} is {}".format(i, tfidf_vectors[i]
                                          ))
89 ----------------------------------------------------------------------------
                                      ------
90 tf-idf vector for document 0 is [array([0.          , 0.          , 0.
                                      , 0.          , 0.05008433,
91      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
92 tf-idf vector for document 1 is [array([0.          , 0.          , 0.
                                      , 0.          , 0.09799108,
93      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
94 tf-idf vector for document 2 is [array([0.02129965, 0.          , 0.
                                      , 0.          , 0.03266369,
95      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
96 tf-idf vector for document 3 is [array([0.05878704, 0.          , 0.
                                      , 0.          , 0.0300506 ,
97      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
98 tf-idf vector for document 4 is [array([0.          , 0.          , 0.03727823
                                      , 0.          , 0.04899554,
99      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
100 tf-idf vector for document 5 is [array([0.          , 0.          , 0.
                                       , 0.          , 0.06091338,
101      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
102 tf-idf vector for document 6 is [array([0., 0., 0., 0., 0., 0., 0., 0., 0
                                       ., 0.]), 0]
103 tf-idf vector for document 7 is [array([0., 0., 0., 0., 0., 0., 0., 0., 0
                                       ., 0.]), 1]
104 tf-idf vector for document 8 is [array([0.          , 0.          , 0.
                                       , 0.03937682, 0.          ,
105      0.          , 0.          , 0.          , 0.          , 0.          ]), 1]
106 tf-idf vector for document 9 is [array([0., 0., 0., 0., 0., 0., 0., 0., 0
                                       ., 0.]), 1]
```

The I use tf-idf weight matrix as features and perform 5-fold cross-validation with the Naive Bayes classifier. The result (average accuracy, precision and reacall across all folds) are shown below.

```
1 from sklearn.model_selection import KFold
2 from sklearn.naive_bayes import GaussianNB
3 from sklearn import metrics
4
```

```python
5  dataset = [[v[0], v[1]] for v in tfidf_vectors]
6  kf = KFold(n_splits=5)
7  avg_acc, avg_precision, avg_recall = [], [], []
8  for train_idx, test_idx in kf.split(dataset):
9      gnb = GaussianNB()
10     X_train, y_train = [dataset[i][0] for i in train_idx], [dataset[i][1]
                                           for i in train_idx]
11     X_test, y_test = [dataset[i][0] for i in test_idx], [dataset[i][1] for
                                           i in test_idx]
12     pred = gnb.fit(X_train, y_train).predict(X_test)
13
14     acc = metrics.accuracy_score(np.array(pred), y_test)
15     print("Accuracy = %.3f" % (sum([p==a for p, a in zip(np.array(pred),
                                           y_test)])/len(np.array(pred))))
16     precision = metrics.precision_score(np.array(pred), y_test)
17     print("Precision = %.3f" % precision)
18     recall = metrics.recall_score(np.array(pred), y_test)
19     print("Recall = %.3f" % recall)
20     print('-----'*4)
21     avg_acc.append(acc)
22     avg_precision.append(precision)
23     avg_recall.append(recall)
24
25 print("Average Accuracy is: ", sum(avg_acc)/5)
26 print("Precision Accuracy is: ", sum(avg_precision)/5)
27 print("Recall Accuracy is: ", sum(avg_recall)/5)
28 --------------------------------------------------------------------------
                                           ------
29
30
31 Accuracy = 0.675
32 Precision = 0.955
33 Recall = 0.636
34 --------------------
35 Accuracy = 0.625
36 Precision = 1.000
37 Recall = 0.583
38 --------------------
39 Accuracy = 0.700
40 Precision = 0.913
41 Recall = 0.677
42 --------------------
43 Accuracy = 0.725
44 Precision = 1.000
45 Recall = 0.667
46 --------------------
47 Accuracy = 0.641
48 Precision = 1.000
49 Recall = 0.622
50 --------------------
51 Average Accuracy is:  0.6732051282051282
52 Precision Accuracy is:  0.9735177865612649
53 Recall Accuracy is:  0.6370809225647934
```

The result shows that cross validation can reduce bias and every data points get to be tested exactly once and is used in training k-1 times. The variance of the resulting estimate is reduced as k increases.

# Reference

- https://jakevdp.github.io/PythonDataScienceHandbook/05.13-kernel-density-estimation.html

- https://stats.stackexchange.com/questions/206963/adaptive-variable-bandwidth-in-2d-gaussian-kernel-density-estimator-to-account