

Non-parametric Classification with Dimensionality Reduction

Kevin De Angeli
kevindeangeli@utk.edu
COSC 522 - Machine Learning
University of Tennessee, Knoxville

October 10, 2019

Abstract

This paper explores the application of K-nearest neighbors (Knn) for classification tasks on a data set that describes diabetes in prime Indian heritage living near Phoenix Arizona. I used Principal Component Analysis (PCA) and Fisher's Linear Discriminant (FLD) to reduce the dimensions of the data set. I have developed a extensive analysis of the performance of Knn with the three different data sets obtained that include accuracy, confusion matrices and running time. Finally, I compared the performance of Knn with three Bayesian classifier under different assumptions.

1 Introduction

K_n Nearest Neighbors (kNN) is one of the core methods for non-parametric classification tasks [1]. This method makes no assumptions about the nature of a population, and makes decisions by merely comparing test points with the training data. This characteristic makes kNN computationally demanding, because each testing point is compared with the entire training data set. The kNN method has been successfully applied to problems in diverse domains. For example, in [2] Sadegh Bafandeh applied kNN to predict economic events, text categorization [3], [4], and visual recognition [5]. In this project, I developed a kNN algorithm from the foundations (using just numpy) and apply it to a data set that contains information about diabetes in prima indians. The algorithms and methods presented in this paper can be easily applied to other data set. Additionally, I applied previously developed Gaussian classifiers and compare the performance. The Python scripts created for this project are provided in the Appendix.

2 Methods

2.1 Data set

The Diabetes Prima Indians data set comes from [1]. It contains data from a population of women who were 21 years old and older. They were women of Prima Indian heritage living in Phoenix, Arizona. Each woman in the data set was tested for diabetes according to World Health Organization criteria. The actual data was collected by the US National Institute of Diabetes and Digestive and Kidney Diseases. The training data set contains 200 rows and the test set has 332 rows. Both have 8 columns, where the last column contains "yes" or "no" values which tells if the person has diabetes. These values were replaced by 0s and 1s at the beginning of the program. 1 shows the specific information contained in the data set.

2.2 Data Pre-processing

2.2.1 Normalization

The first step was to normalize the data so that the values in the data set are scaled similarly. I normalized the data by simply subtracting the mean and dividing by the standard deviation of the specific class:

$$x_i = \frac{x_i - \mu_i}{\sigma_i}$$

I used the normalized data as a data set on its own, and I called it nX .

Column Name	Description
npreg	number of pregnancies
glu	plasma glucose concentration in an oral glucose tolerance test
bp	diastolic blood pressure (mm Hg)
skin	triceps skin fold thickness (mm)
ins	serum insulin (micro U/ml)
bmi	body mass index (weight in $kg/(heightinm)^2$)
ped	diabetes pedigree function
age	in years
type	No / Yes

Table 1

2.2.2 Principal Component Analysis (PCA)

PCA is a unsupervised method for data representation that allows us to reduce the complexity of the data. In order to find the principal components of the data set, I first computed the 7×7 Variance-Covariance matrix of the training data set and the corresponding eigenvalues and eigenvectors:

$$\Sigma = \begin{bmatrix} 1.005025 & 0.171382 & 0.253328 & 0.109597 & 0.058629 & -0.120073 & 0.601932 \\ 0.171382 & 1.005025 & 0.270735 & 0.21869 & 0.217879 & 0.061015 & 0.345133 \\ 0.253328 & 0.270735 & 1.005025 & 0.266295 & 0.240021 & -0.047638 & 0.393039 \\ 0.109597 & 0.21869 & 0.266295 & 1.005025 & 0.662347 & 0.095882 & 0.253192 \\ 0.058629 & 0.217879 & 0.240021 & 0.662347 & 1.005025 & 0.191508 & 0.132582 \\ -0.120073 & 0.061015 & -0.047638 & 0.095882 & 0.191508 & 1.005025 & -0.071768 \\ 0.601932 & 0.345133 & 0.393039 & 0.253192 & 0.132582 & -0.071768 & 1.005025 \end{bmatrix}$$

$$\lambda = [2.42136801 \quad 1.50396726 \quad 0.91644987 \quad 0.80399111 \quad 0.69360322 \quad 0.39177751 \quad 0.3040189]$$

Figure 1 shows the information expressed by each of the eigenvalues. Running some analysis of these values leads to the conclusion that using 5 of the 7 eigenvectors leads to a error of less than 10%. The corresponding eigenvectors have not been added here, but they can be easily visualized by running the python program. The 5 eigenvectors associated with the 5 greatest eigenvalues were put into a matrix of 5×7 dimensions and they were used to reduce the data set into 5 dimensions.

2.2.3 Fisher Linear Discriminant (FLD)

Unlike, PCA Fisher Linear Discriminant is a supervised technique of dimensional reduction. In other words, this method takes the label of the data into consideration to calculate the necessary statistics. Using FLD, I projected the data set into one dimension. In order to find the projection vector, I calculated the following statistics:

$$\begin{aligned} S_1 &= (n_1 - 1)\Sigma_1, & S_2 &= (n_2 - 1)\Sigma_2 \\ S_w &= S_1 + S_2 \\ S_w^{-1} &= inv(S_w) \\ V &= S_w^{-1}(\mu_1 - \mu_2) \end{aligned}$$

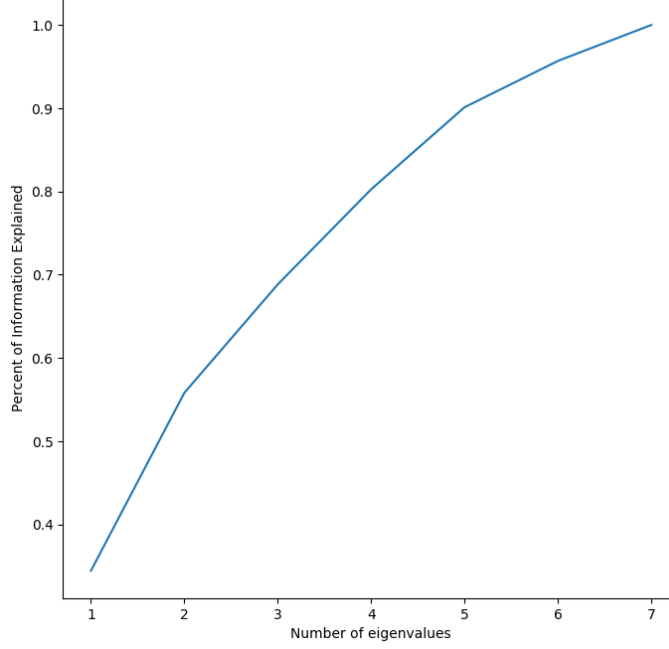


Figure 1: Cumulative information captured by number of eigenvalues. The eigenvalues were ordered from greatest to smallest and then they were added one by one. When $x=7$, all the eigenvalues are used and there is no dimensional reduction.

The resulting vector V has dimensions 1×7 and is used to transform the data into just one dimension.

For this project, the resulting vector obtained was:

$$V = [-0.00138917, -0.00423824, -0.00021323, 0.00039388, -0.00220392, -0.00186505, -0.00192072]$$

2.3 kNN

First, I want to point out that kNN can be understood in terms of posterior probabilities:

$$p(w_i|x) = \frac{p(x|w_i)p(w_i)}{p(x)}$$

$$p(w_i|x) = \frac{\frac{k_i}{n_i V} \frac{n_i}{n}}{\frac{k}{nV}} = \frac{k_i}{k}$$

In terms of algorithms, I have written 3 functions that work together to classify data based on the k closest Neighbors. The function *knn* accepts the training and testing data set, and an integer value for k . This function loops through the testing data set and calculates the accuracy. The function *euclidiandsitance* computes the distance between the test point that it receives and every other point in the training data set. Then, it sends the results to *guessLabel*, which sorts the array of data and makes a guess based on the k closest training data points.

2.4 Performance Evaluation

Different performance matrices has been used to analyze the models. Additionally, the computational time of each model was computed.

In terms of accuracy, I computed the ratio:

$$\frac{\text{Number of labels guessed correctly}}{\text{Number of total rows}}$$

For kNN the accuracy was calculated through a range of K values from 1 to 100.

In order to further evaluate the classifier's accuracy, I have also compared class-wise accuracy. I created bar plots that compares the True Positive (TP), True Negative (TN), False Positive (FP), and False Negative (FN) values. Assuming that class 1 is "positive" and 0 is "negative", I calculated the Sensitivity or True Positive Rate (TPR):

$$\frac{TP}{TP+FN}$$

and the Specificity or True Negative Rate (TNR):

$$\frac{TN}{TN+FP}$$

Sensitivity and Specificity were plotted against different values of prior probabilities.

2.5 Prior Probabilities

In order to run a prior probability analysis with the kNN model, I introduced two constant α_1 and α_2 . Note that based on the data, the prior probability for class 0 is 132/200. In order to find α_1 and α_2 I solved the following equation:

$$\alpha_1 \frac{132}{200} + \alpha_2 \frac{68}{200} = 1$$

We can first check the range of values that α_1 can take:

$$0 \leq \alpha_1 \frac{132}{200} \leq 1$$

So,

$$0 \leq \alpha_1 \leq \frac{200}{132}$$

And therefore,

$$\alpha_2 = \frac{200}{68} - \frac{132}{68} \alpha_1$$

3 Results

3.1 Performance Curves

3.1.1 Normalized Data - nX

Figure 2a shows the performance of kNN using the normalized data set. Running the algorithm with k=1 lead to an accuracy of 71.3% and it took 0.958 seconds. The maximum accuracy is obtained when k=17, and that leads to an accuracy of 80.1%.

3.1.2 Transformed Data set: PCA - pX

In the case of the data set resulting from the PCA transformation, I obtained an accuracy of 60.2% when k=1, and a total accuracy of 73.1% when k=11. The performance curve is displayed in Figure 2b. Running kNN with k=1 took 0.803 seconds.

3.1.3 Transformed Data set: FLD - fX

The FLD data set lead to the highest accuracy (Figure 2c). It obtained a total accuracy of 80.4% when $k=21$. Additionally, the FLD data set also provided the greatest accuracy with $k=1$, since it classified 79.2% of the entries correctly.

3.1.4 Comparison Table

Table 2 provides a direct comparison of the performance curve analysis presented in this section.

Performance Curves Comparison			
Data Set	Accuracy when K=1	K When Maxi- mum Accuracy	Maximum Accu- racy
nX	71.3	17	80.1
pX	60.2	11	73.1
fX	79.2	21	79.2

Table 2

3.2 Performance Comparison

3.2.1 Normalized Data - nX

Graph 3 shows class-wise comparison of kNN vs. each of the 3 Gaussian classifiers. Here, I chose $k=1$ as the kNN parameter. The total accuracy for Case 1 with the nX data set was 75.0%. Plot 3a shows that kNN produce a large amount of false negatives, but much fewer false positives than the Case I Gaussian classifier. Case II of the Gaussian classifier obtained a total accuracy of 77%, and Case III lead to a 70.5% accuracy.

3.2.2 Transformed Data set: PCA - pX

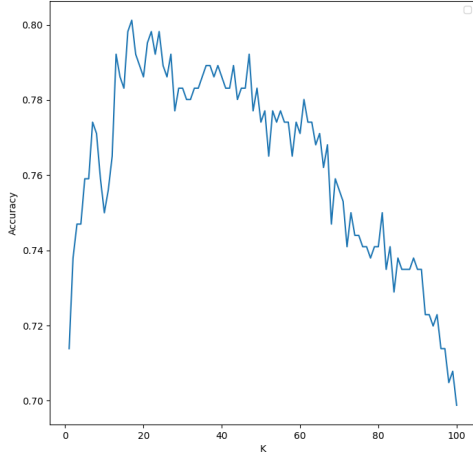
When running the pX data set through the Gaussian cases, I obtained 67.5%, 70.8%, and 60.8% for case I,II,and III respectively (Figure 4a). Here, I also run kNN with $k=1$. One interesting aspect of these results is how poorly case III performed. Even though, it's the graph that makes the least assumptions about the data.

3.2.3 Transformed Data set: FLD - fX

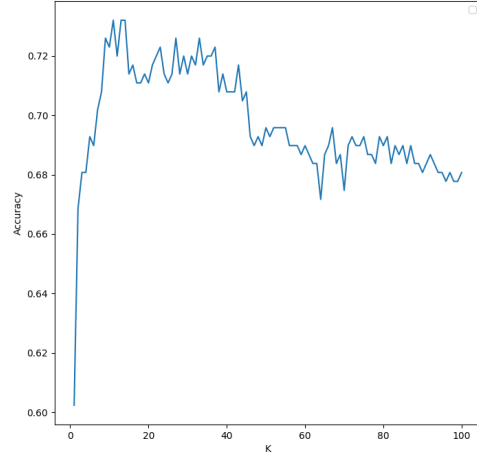
I was not able to come up with results for Case I, Case II, and Case III in the case of the fX data set. The main reason seems to be problems related to calculating statistics such as the covariance matrix, which don't work when the data is in one dimension. I attempted to just use standard deviation, but the program became convoluted enough that it was not possible to provide a solution on time.

3.3 Prior Probability Analysis

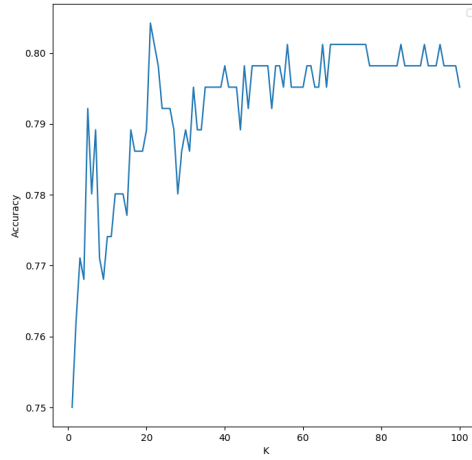
I have run the 3 Gaussian cases together and check the class-wise accuracy. Based on the model output, I created sensitivity and specificity plots.



(a) nX Data set



(b) pX Data Set



(c) fX Data Set

Figure 2: The plots display accuracy vs. K values for the three transformed data sets.

3.3.1 Normalized Data - nX

Figure 5 shows the class-wise accuracy for the three Gaussian classifiers used in this paper. All the classifier were trained using the normalized data. From Figure 5a, Case I seems to provide similar results than kNN.

3.3.2 Transformed Data set: PCA - pX

Similarly as in the previous section, I also plotted the class-wise accuracy of the three Gaussian classifiers and kNN for the pX Data set (Figure 6). Unlike the previous picture, here kNN seems to lead to a more similar result to Case II than the rest.

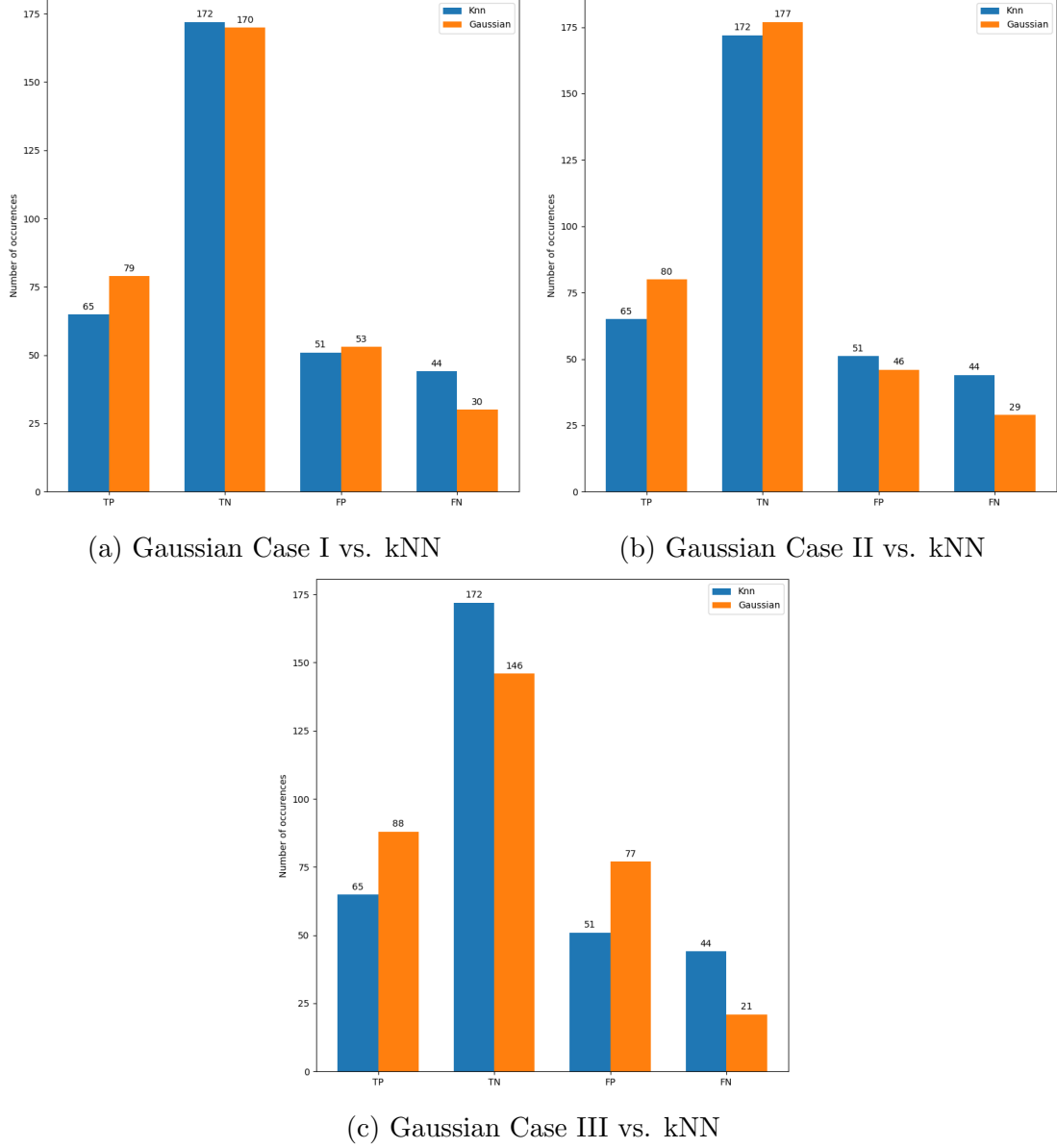


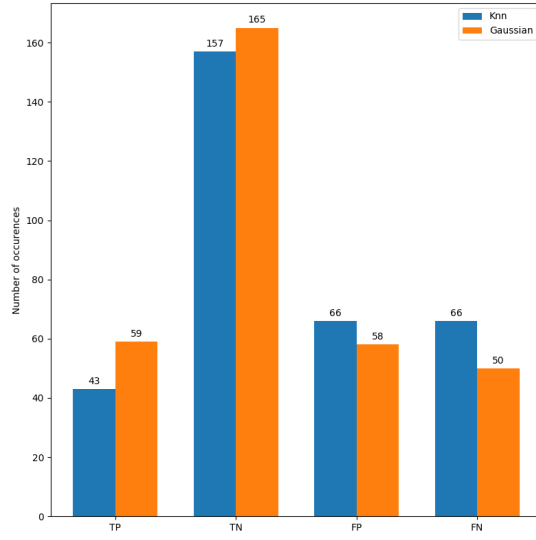
Figure 3: Class-wise accuracy comparison for the Gaussian cases vs. kNN on the nX data set

3.3.3 Transformed Data set: FLD - fX

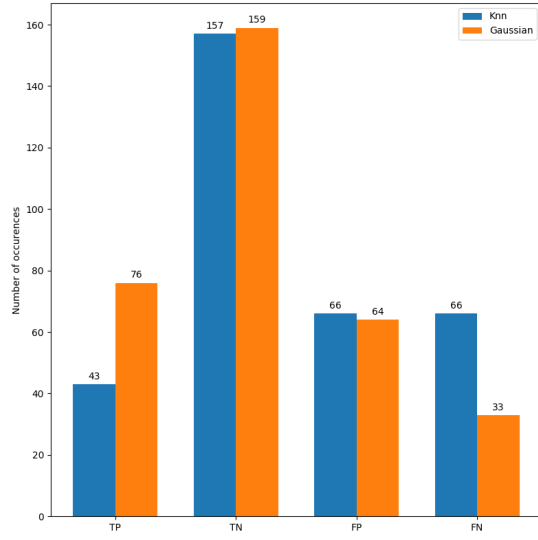
Lastly, the class-wise accuracy for the fX dataset (Figure 7) provide results in similar patterns with Case II and kNN providing almost identical results, and Case III and kNN showing the most differences.

3.4 Eigenvalues vs. Sensitivity and Specificity

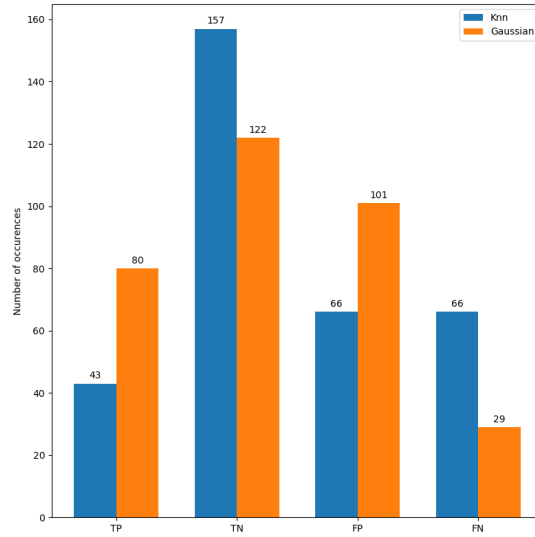
Figure 8 shows the resulting graphs after transforming the data through different numbers of eigenvectors. I first ordered the eigenvalues from greatest to smallest, and I added their eigenvectors to a transformation matrix one by one to compare the results. Each time I added one of the vectors, I transformed the datasets and check the results.



(a) Gaussian Case I vs. kNN

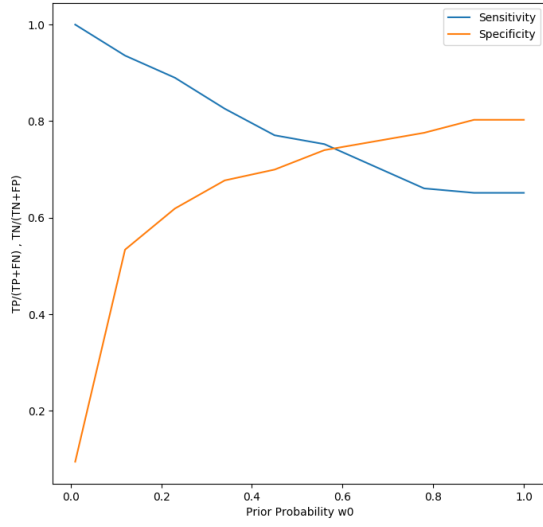


(b) Gaussian Case II vs. kNN

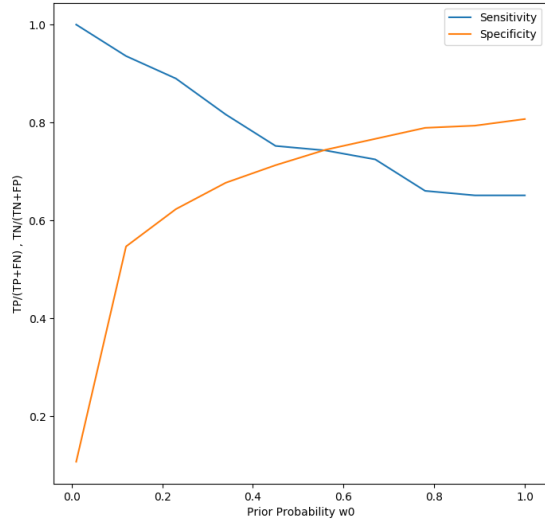


(c) Gaussian Case III vs. kNN

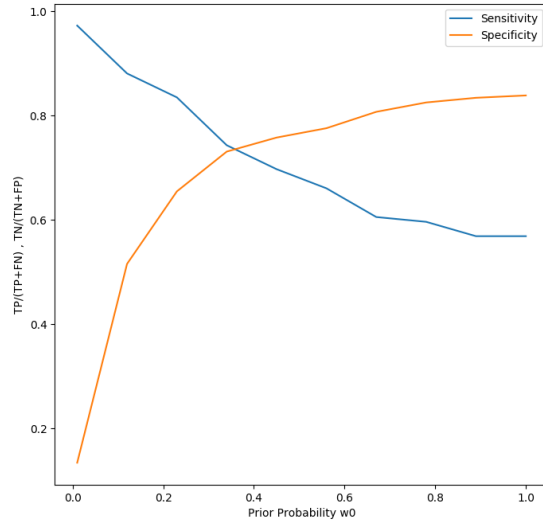
Figure 4: Class-wise accuracy comparison for the Gaussian cases vs. kNN on the pX data set



(a) Gaussian - Case 1

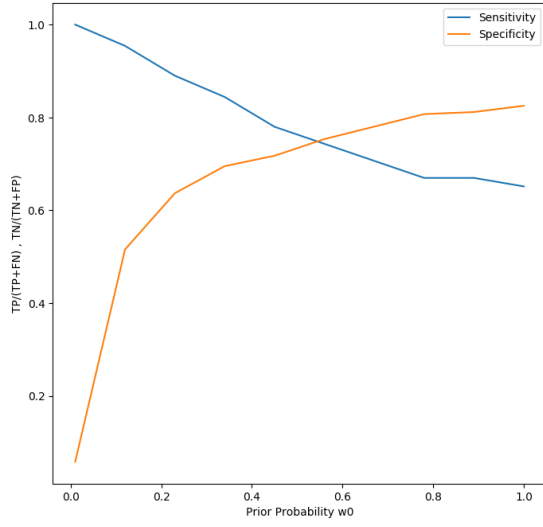


(b) Gaussian - Case 2

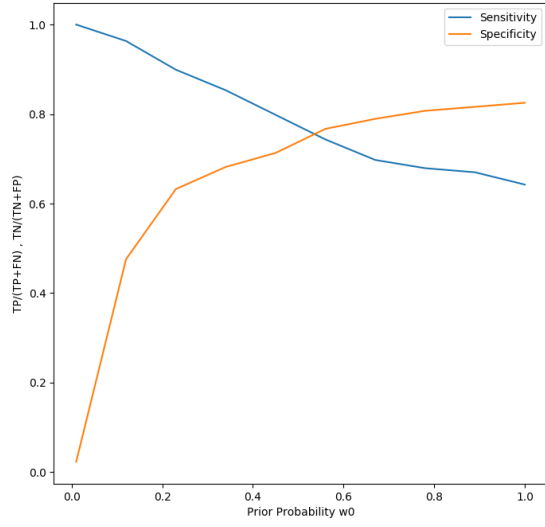


(c) Gaussian - Case 3

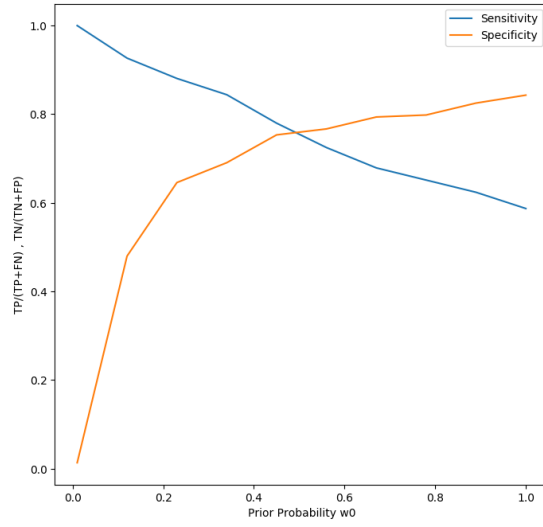
Figure 5: Sensitivity and Specificity curves for the normalized data (nX)



(a) Gaussian - Case 1

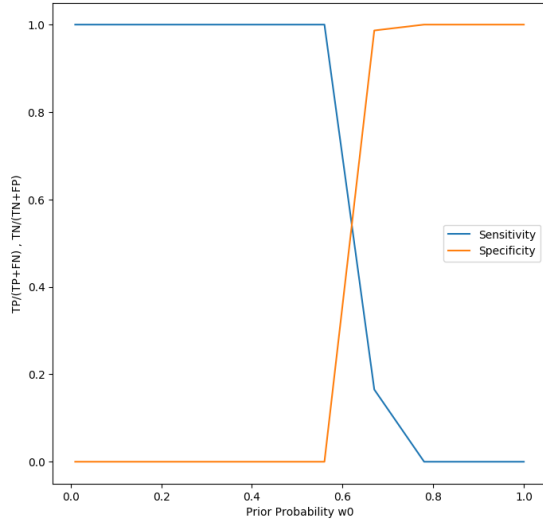


(b) Gaussian - Case 2

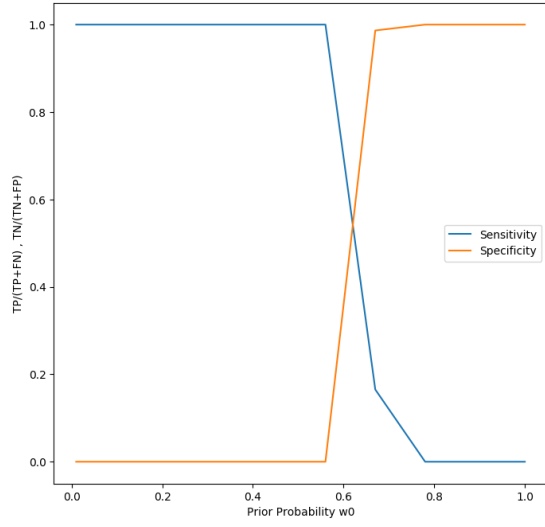


(c) Gaussian - Case 3

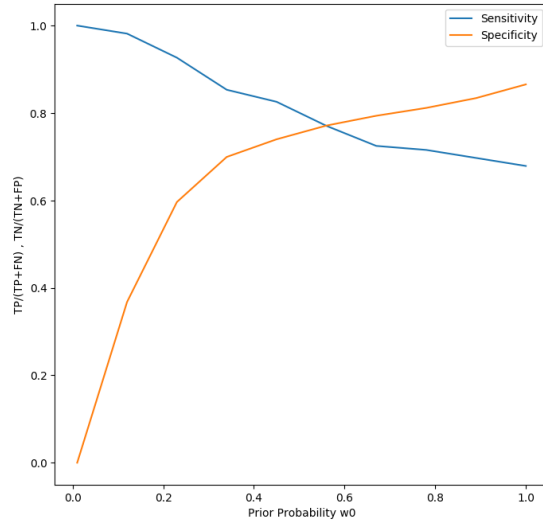
Figure 6: Sensitivity and Specificity curves for the PCA transformed data set (pX)



(a) Gaussian - Case 1

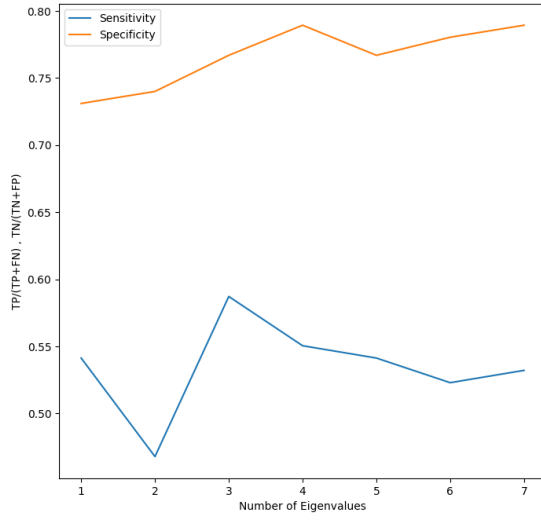


(b) Gaussian - Case 2

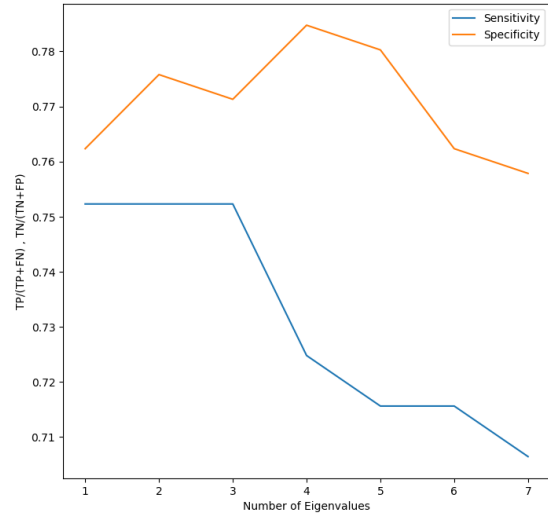


(c) Gaussian - Case 3

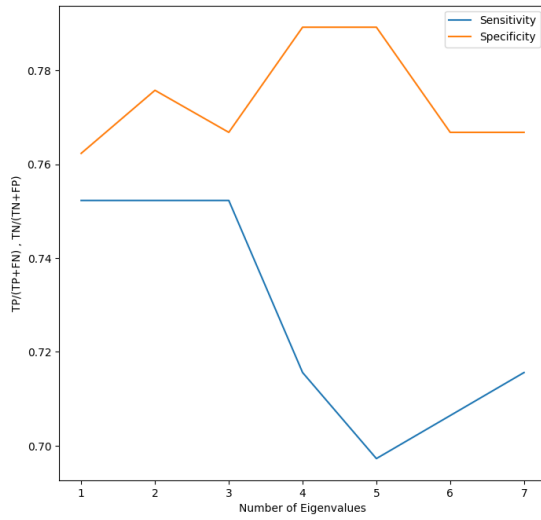
Figure 7: Sensitivity and Specificity curves for the FLD transformed data set (fX)



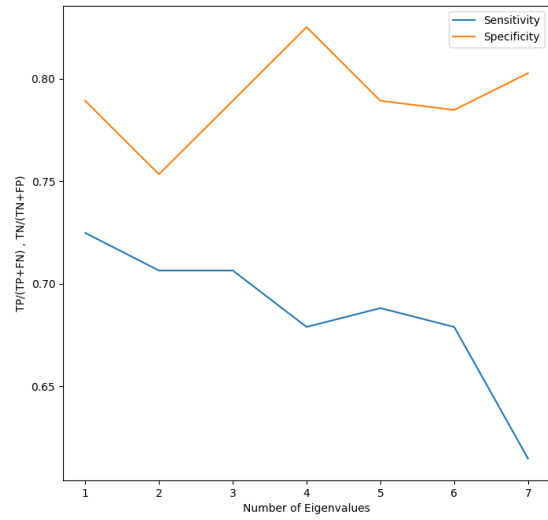
(a)



(b)



(c)



(d)

Figure 8: Cumulative eigenvectors plotted against Sensitivity and Specificity.

4 Discussion

This paper provides an example application of the k-nearest neighbors algorithm that used real world data. In addition to the algorithm itself, I presented the theoretical background behind why kNN actually provides the maximum posterior probability. This project also developed the ideas behind pre-processing data such as normalization and dimension reduction. As expected of a supervised method, Fisher Linear Discriminant method lead to better results than Principal Component Analysis transformation in general. Overall, kNN performed above 70% with all three data sets. That's a great lower bound considering that it has also reached a 80% accuracy when $k=17$ in the nX data set. Additionally, I have shown that the Gaussian classifiers have succeeded in terms of accuracy. One of the non-intuitive aspects that I found out is that Case II often outperformed Case I and III. I used to think that making fewer assumptions about the data is in general better, but this project proved me somehow wrong. The python program available in the appendix can be easily applied to other data set and used as tools to not just classify data, but also to compare Gaussian classifiers with kNN and see contrast plots.

In a personal note, working on Project 1 and 2 have served as great experience in becoming a more efficient and professional programmer while merging statistical theory with applications. I've enjoyed the programming aspect extensively, and I tried to make my code as efficient and generalizable as possible.

References

- [1] Richard O. Duda, Peter E. Hart, David G. Stork. *Pattern Classification*. Second Edition. pdf.
- [2] Sadegh B. Imandoust And Mohammad Bolandraftar *Application of K-Nearest Neighbor (KNN) Approach for Predicting Economic Events: Theoretical Background*. pdf.
- [3] Gongde GuoHui WangDavid BellYaxin BiKieran Greer *An kNN Model-Based Approach and Its Application in Text Categorization*. pdf.
- [4] Bijalwan V., Kumar V., Kumari P., Pascual P. *KNN based Machine Learning Approach for Text and Document Mining*. pdf.
- [5] Liu Q. , Liu C *A Novel Locally Linear KNN Method With Applications to Visual Recognition*. pdf.

5 Appendix

5.1 Python Script

```
1  import numpy as np
2  import pandas as pd
3  import math #Used for Pi and log()
4  import sympy as sym
5  import matplotlib.pyplot as plt
6  from mpl_toolkits.mplot3d import Axes3D
7  from numpy import ones,vstack
8  import time
9
10
11 def euc2(x, y):
12     # calculate squared Euclidean distance
13
14     # check dimension
15     assert x.shape == y.shape
16
17     diff = x - y
18
19     return np.dot(diff, diff)
20
21
22 def mah2(x, y, Sigma):
23     # calculate squared Mahalanobis distance
24
25     # check dimension
26     assert x.shape == y.shape and max(x.shape) == max(Sigma.shape)
27
28     diff = x - y
29
30     return np.dot(np.dot(diff, np.linalg.inv(Sigma)), diff)
31
32
33
34 def load_data(f):
35     #f is the route of the data file.
36     x = pd.read_csv(f, delim_whitespace=1, header=None)
37     x[x.shape[1] - 1] = x[x.shape[1] - 1].map({'Yes': 1, 'No': 0})
38     #Could return x if I wanted to keep it as the first project.
39     X = x.loc[:,0:x.shape[1]-2]
40     Y = x.loc[:,x.shape[1]-1]
41     Y=Y.to_numpy()
42     X=X.to_numpy()
43     X=normalization(X)
```

```

44     return X, Y
45
46 def normalization(X):
47     meanArr = np.mean(X, axis=0)
48     varArr = np.std(X, axis=0)
49     nX = X[:]
50     for i in range(X.shape[1]):
51         nX[:, i] = (X[:, i] - meanArr[i]) / varArr[i]
52     return nX
53
54 def pca(nX,percentError=.9,showGraph=False):
55
56     nX_Cov = np.cov(nX.T) #Note that covariannce in pd is calculated differently
57     nX_eig, nX_eigV = np.linalg.eig(nX_Cov)
58     ordered_eigs = -np.sort(-nX_eig)
59     totalSum = np.sum(ordered_eigs)
60     # Store the indexes of the ordered eigenvalues
61     # So you can create a matrix of eigenvectors
62     # it the correct order:
63     order_index_eigs = []
64     for i in range(ordered_eigs.shape[0]):
65         order_index_eigs.append(np.where(nX_eig == ordered_eigs[i])[0].item())
66
67
68     # Store the indexes of the ordered eigenvalues
69     # So you can create a matrix of eigenvectors
70     # it the correct order:
71     order_index_ordered_eigs = []
72     for i in range(ordered_eigs.shape[0]):
73         order_index_ordered_eigs.append(np.where(nX_eig ==
74             → ordered_eigs[i])[0].item())
75
76 def eigenValErrorAnalysis(ordered_eigs):
77     plt.figure(num=None, figsize=(8, 8), dpi=100, facecolor='w', edgecolor='k')
78     x = np.linspace(1, ordered_eigs.shape[0], ordered_eigs.shape[0])
79     # Cumulative:
80     eg = []
81     totalEig = 0
82     for i in range(ordered_eigs.shape[0]):
83         totalEig += ordered_eigs[i]
84         eg.append(totalEig / totalSum)
85     plt.plot(x, eg)
86     plt.xlabel(xlabel='Number of eigenvalues')
87     plt.ylabel(ylabel='Percent of Information Explained')
88     plt.show()

```



```

89
90 if showGraph == True:
91     eigenValErrorAnalysis(ordered_eigs)
92
93 def numberOfEigValsToUse(eigs,percentError):
94     tmp = 0
95     ix = -1
96     P = []
97     for i in range(eigs.shape[0]):
98         tmp += eigs[i] / totalSum
99         if tmp > percentError:
100             ix = i + 1 # Number of eigenvectors to use is 5
101             for k in range(ix):
102                 a = nX_eigV[order_index_eigs[k]]
103                 P.append(a)
104             P = np.array(P)
105             return P
106
107 P = numberOfEigValsToUse(ordered_eigs,percentError)
108 ans=np.dot(nX, np.transpose(P))
109 return ans
110
111 def fld(nX,y=0, training= True):
112     if training==True:
113         #split the data int two classes:
114         #key = y[:, 0] == 0
115         key0 = y==0
116         y0Values = nX[key0]
117         key1 = y == 1
118         y1Values= nX[key1]
119
120         y0ValuesMean = np.mean(y0Values,axis=0)
121         y1ValuesMean = np.mean(y1Values,axis=0)
122
123
124         y0Cov = np.cov(y0Values.T)
125         y1Cov = np.cov(y1Values.T)
126
127         S_0 = (y0Values.shape[0] - 1) * y0Cov
128         S_1 = (y0Values.shape[0] - 1) * y1Cov
129
130         S_w = S_0 + S_1
131         S_w_inv = np.linalg.inv(np.array(S_w))
132
133         fld.v = np.dot(S_w_inv, (np.transpose(y0ValuesMean) -

```

```

134
135
136     y0Values=np.dot(y0Values, fld.v)
137     y1Values=np.dot(y1Values, fld.v)
138
139
140     nX[:, 0][key0] = y0Values
141     nX[:, 0][key1] = y1Values
142     nX=np.delete(nX, np.s_[1:nX.shape[1]], axis=1)
143 else:
144     nX = np.dot(nX, fld.v)
145 return nX
146
147
148
149
150 class Knn:
151     def __init__(self):
152         self.nX = []
153         self.pX = []
154         self.fX = []
155         self.predictionArr =[]
156         self.totalTime= -1
157
158     def showTime(self):
159         print("Time in seconds: ", self.totalTime)
160
161     def fit(self, X, y):
162         #self.nX = normalization(X)
163         self.nX = X
164         self.pX = pca(self.nX)
165         self.fX = fld(self.nX,y)
166         self.y = y
167
168     #x here is just a point
169     def euclidian_dsitanceList(self, x, X):
170         # x is the test point, X is the dataset
171         # Using Euclidian Distance
172         distancesArr = []
173         # For each row in the data set:
174         dist = -1
175         index=0 # used to associate a distance with a y.
176         for row in X:
177             testPoint = row[0:X.shape[1]]
178             dist = np.linalg.norm(testPoint - x)
179             labelIndex = self.y[index]

```

```

180         distanceAndLabel = (dist, labelIndex)
181         distancesArr.append(distanceAndLabel)
182         index +=1
183     return distancesArr
184
185
186     #x here is just a point
187 def guessLabel(self,x, X, k):
188     label0 = 0
189     guessClass = -1
190     label1 = 0
191     ks = []
192     distancesArr = self.euclidian_distanceList(x, X)
193     distancesArr.sort()
194     for p in range(int(k)):
195         minimum = min(distancesArr)
196         ks.append((minimum[0], minimum[1]))
197         distancesArr.remove(minimum)
198     for q in range(len(ks)):
199         if ks[q][1] == 0:
200             label0 += 1
201         else:
202             label1 += 1
203     if label0 >= label1:
204         guessClass = 0
205     else:
206         guessClass = 1
207     return guessClass
208
209
210 def knn(self,XTest, X, k):
211     #X is the training data
212     guessList =[]
213     guessedLabel = -1
214     #print(XTest)
215     for row in XTest:
216         guessedLabel = self.guessLabel(row, X, k)
217         guessList.append(guessedLabel)
218     return guessList
219
220
221
222 def predict(self, XTest, k=1, data="nX"):
223     start_time = time.time()
224     predictionArr =[]
225     #print(XTest)

```

```

226         if data == "nX":
227             self.predictionArr = self.knn(XTest, self.nX, k)
228         elif data == "pX":
229             XTest = pca(XTest)
230             self.predictionArr = self.knn(XTest, self.pX, k)
231         else:
232             XTest=fld(XTest, training=False)
233             self.predictionArr = self.knn(XTest, self.fX, k)
234             #print(self.predictionArr)
235         self.totalTime= time.time() - start_time
236         return self.predictionArr
237
238     def performanceCurve(self, X_test,ytest, K, data="fX"):
239         k = np.arange(1, K+1, 1).tolist()
240         accuracy_list = []
241         for i in k:
242             array_prediction=self.predict(X_test, i, data)
243
244             ↪ accuracy_list.append(accuracy_score(ytest,array_prediction,showResults=False))
245             #print(accuracy_list)
246             print("Maximum Accuracy is obtained when K=", k[np.argmax(accuracy_list)])
247             print("The accuracy associated with that K = ", np.amax(accuracy_list))
248             plt.figure(num=None, figsize=(8, 8), dpi=100, facecolor='w', edgecolor='k')
249             plt.plot(k, accuracy_list)
250             plt.legend(loc='best')
251             plt.xlabel('K')
252             plt.ylabel('Accuracy')
253             plt.show()
254
255
256
257
258
259     def accuracy_score(yTest, y_model, showResults=True):
260         TP, TN, FP, FN = 0, 0, 0, 0
261         index=0
262         for rightLabel in yTest:
263             guessLabel=y_model[index]
264             index+=1
265             if guessLabel == rightLabel:
266                 if rightLabel == 1:
267                     TP += 1
268                 else:
269                     TN += 1
270             else:

```

```

271         if rightLabel == 1:
272             FN += 1
273         else:
274             FP += 1
275     # print(TP)
276     totalRowsInData = yTest.shape[0]
277     confusion_matrix = [['TP', TP], ["TN", TN], ['FP', FP], ['FN', FN],
278         → ['Accuracy', (TN + TP) / totalRowsInData]]
279     confusionarr = [TP, TN, FP, FN]
280     accuracy = confusion_matrix[4][1]
281     if showResults:
282         print(confusion_matrix)
283     return accuracy, confusionarr
284
285
286 class mpp:
287     def __init__(self, case=1):
288         # init prior probability, equal distribution
289         # self.classsn = len(self.classes)
290         # self.pw = np.full(self.classsn, 1/self.classsn)
291
292         # self.covs, self.means, self.covavg, self.varavg = \
293         #     self.train(self.train_data, self.classes)
294         self.case_ = case
295         self.pw_ = None
296
297     def fit(self, Tr, y, fX=False):
298         # derive the model
299         self.covs_, self.means_ = {}, {}
300         self.covsum_ = None
301
302         self.classes_ = np.unique(y) # get unique labels as dictionary items
303         self.classsn_ = len(self.classes_)
304
305         for c in self.classes_:
306             arr = Tr[y == c]
307             self.covs_[c] = np.cov(np.transpose(arr))
308             self.means_[c] = np.mean(arr, axis=0) # mean along rows
309             if self.covsum_ is None:
310                 self.covsum_ = self.covs_[c]
311             else:
312                 self.covsum_ += self.covs_[c]
313
314         if fX==False:
315             # used by case II

```

```

316         self.covavg_ = self.covsum_ / self.classn_
317
318         # used by case I
319         self.varavg_ = np.sum(np.diagonal(self.covavg_)) / len(self.classes_)
320     else:
321         self.covavg_ = np.std(Tr)
322         self.varavg_ = np.var(Tr)
323
324     def predict(self, T):
325         # eval all data
326         y = []
327         disc = np.zeros(self.classn_)
328         nr, _ = T.shape
329
330         if self.pw_ is None:
331             self.pw_ = np.full(self.classn_, 1 / self.classn_)
332
333         for i in range(nr):
334             for c in self.classes_:
335                 if self.case_ == 1:
336                     edist2 = euc2(self.means_[c], T[i])
337                     disc[c] = -edist2 / (2 * self.varavg_) + np.log(self.pw_[c])
338                 elif self.case_ == 2:
339                     mdist2 = mah2(self.means_[c], T[i], self.covavg_)
340                     disc[c] = -mdist2 / 2 + np.log(self.pw_[c])
341                 elif self.case_ == 3:
342                     mdist2 = mah2(self.means_[c], T[i], self.covs_[c])
343                     disc[c] = -mdist2 / 2 - np.log(np.linalg.det(self.covs_[c])) /
344                         ↪ 2 \
345                         + np.log(self.pw_[c])
346                 else:
347                     print("Can only handle case numbers 1, 2, 3.")
348                     sys.exit(1)
349             y.append(disc.argmax())
350
351         return y
352
353     def plotBarsAccuracy(confusionM1, confusionM2):
354         plt.figure(num=None, figsize=(8, 8), dpi=100, facecolor='w', edgecolor='k')
355         labels = ['TP', 'TN', 'FP', 'FN']
356         x = np.arange(len(labels)) # the label locations
357         width = 0.35 # the width of the bars
358
359         fig, ax = plt.subplots(num=None, figsize=(8, 8), dpi=100, facecolor='w',
360                               ↪ edgecolor='k')

```

```

360 rects1 = ax.bar(x - width/2, confusionM1, width, label='Knn')
361 rects2 = ax.bar(x + width/2, confusionM2, width, label='Gaussian')
362
363 # Add some text for labels, title and custom x-axis tick labels, etc.
364 ax.set_ylabel('Number of occurrences')
365 #ax.set_title('Scores by group and gender')
366 ax.set_xticks(x)
367 ax.set_xticklabels(labels)
368 ax.legend()
369
370
371 def autolabel(rects):
372     """Attach a text label above each bar in *rects*, displaying its
373     ↪ height."""
374     for rect in rects:
375         height = rect.get_height()
376         ax.annotate('{}'.format(height),
377                     xy=(rect.get_x() + rect.get_width() / 2, height),
378                     xytext=(0, 3), # 3 points vertical offset
379                     textcoords="offset points",
380                     ha='center', va='bottom')
381
382 autolabel(rects1)
383 autolabel(rects2)
384
385 fig.tight_layout()
386 plt.show()
387
388
389
390
391 def main():
392     trainingData=
393     ↪ '/Users/kevindeangeli/Desktop/Fall2019/COSC522/Project2/Dataset_PimaTr.txt'
394     testingData=
395     ↪ '/Users/kevindeangeli/Desktop/Fall2019/COSC522/Project2/Dataset_PimaTest.txt'
396     Xtrain,ytrain = load_data(trainingData)
397     Xtest,ytest = load_data(testingData)
398     model = Knn()
399     model.fit(Xtrain, ytrain)
400
401     #Predicts accepts values "nX", "pX", "fX". "nX" is by default.
402     y_model = model.predict(Xtest,k=1, data="pX")
403     accuracy1, confusion_matrix1 = accuracy_score(ytest, y_model)
404     #model.showTime()

```

```

403     #model.performanceCurve(Xtest,ytest, K=100, data="fX")
404
405
406
407
408     Xtrain2, ytrain2 = load_data(trainingData)
409     Xtrain2=pca(Xtrain2)
410     Xtest2, ytest2 = load_data(testingData)
411     Xtest2=pca(Xtest2)
412     model = mpp(3)
413     model.fit(Xtrain2, ytrain2,fX=False)
414     y_model2 = model.predict(Xtest2)
415     accuracy2, confusion_matrix2 = accuracy_score(ytest2, y_model2)
416
417
418
419
420     plotBarsAccuracy(confusion_matrix1,confusion_matrix2)
421
422
423
424     if __name__ == "__main__":
425         main()

```