

Aprendizagem 2022  
Homework II – Group 58

**Part I:** Pen and paper

1.

Homework II								
(I)	1)	Distance Matrix						
	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$
positive	$\frac{1}{2}$	$\frac{5}{2}$	$\frac{3}{2}$	$\frac{1}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{5}{2}$
		$\frac{1}{2}$	$\frac{3}{2}$	$\frac{5}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{7}{2}$
			$\frac{1}{2}$	$\frac{3}{2}$	$\frac{5}{2}$	$\frac{5}{2}$	$\frac{1}{2}$	$\frac{3}{2}$
				$\frac{1}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{3}{2}$	$\frac{5}{2}$
negative					$\frac{1}{2}$	$\frac{5}{2}$	$\frac{3}{2}$	$\frac{3}{2}$
						$\frac{1}{2}$	$\frac{5}{2}$	$\frac{3}{2}$
							$\frac{1}{2}$	$\frac{3}{2}$
							$\frac{1}{2}$	$\frac{7}{2}$

$S_{nn} \rightarrow 5$  nearest neighbours

$S_{nn}(x_1) = [x_4, x_3, x_5, x_6, x_7]$      $S_{nn}(x_5) = [x_6, x_1, x_2, x_4, x_8]$

$S_{nn}(x_2) = [x_8, x_3, x_5, x_6, x_7]$      $S_{nn}(x_6) = [x_5, x_1, x_2, x_4, x_8]$

$S_{nn}(x_3) = [x_7, x_1, x_2, x_4, x_8]$      $S_{nn}(x_7) = [x_3, x_1, x_2, x_4, x_8]$

$S_{nn}(x_4) = [x_1, x_3, x_5, x_6, x_7]$      $S_{nn}(x_8) = [x_2, x_3, x_5, x_6, x_7]$

weight =  $\frac{1}{\text{distance}}$

		Weight Matrix							
		$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$
positive	$x_1$	2	$\frac{2}{5}$	$\frac{2}{3}$	2	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{3}{5}$	$\frac{2}{5}$
	$x_2$	$\frac{2}{5}$	2	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	2
	$x_3$	$\frac{2}{3}$	$\frac{2}{3}$	2	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{5}$	2	$\frac{2}{3}$
	$x_4$	2	$\frac{2}{5}$	$\frac{2}{3}$	2	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{5}$
negative	$x_5$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{3}$	2	2	$\frac{2}{5}$	$\frac{2}{3}$
	$x_6$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{3}$	2	2	$\frac{2}{5}$	$\frac{2}{3}$
	$x_7$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{5}$	2	$\frac{2}{3}$
	$x_8$	$\frac{2}{5}$	2	$\frac{2}{3}$	$\frac{2}{5}$	$\frac{2}{3}$	$\frac{2}{3}$	$\frac{2}{3}$	2

$$S_{nn}(x_1) = [x_4, x_3, x_5, x_6, x_7]$$

$$WA(x_1) = \left(\frac{2}{3} + 2\right) \times (+1) + \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (-1) = \frac{8}{3} - \frac{6}{3} = \frac{2}{3} > 0$$

$$\text{Class}(x_1) = P \quad \text{Class Predicted}(x_1) = P$$

$$S_{nn}(x_2) = [x_8, x_3, x_5, x_6, x_7]$$

$$WA(x_2) = \frac{2}{3} \times (+1) + \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3} + 2\right) \times (-1) = \frac{2}{3} - \frac{12}{3} = -\frac{10}{3} < 0$$

$$\text{Class}(x_2) = N \quad \text{Class Predicted}(x_2) = N$$

$$S_{nn}(x_3) = [x_7, x_1, x_2, x_4, x_8]$$

$$WA(x_3) = \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (+1) + \left(2 + \frac{2}{3}\right) \times (-1) = \frac{6}{3} - \frac{8}{3} = -\frac{2}{3} < 0$$

$$\text{Class}(x_3) = P \quad \text{Class Predicted}(x_3) = N$$

$$S_{nn}(x_4) = [x_1, x_3, x_5, x_6, x_7]$$

$$WA(x_4) = \left(2 + \frac{2}{3}\right) \times (+1) + \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (-1) = \frac{8}{3} - \frac{6}{3} = \frac{2}{3} > 0$$

$$\text{Class}(x_4) = P \quad \text{Class Predicted}(x_4) = P$$

$$S_{nn}(x_5) = [x_6, x_1, x_2, x_4, x_8]$$

$$WA(x_5) = \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (+1) + \left(2 + \frac{2}{3}\right) \times (-1) = \frac{6}{3} - \frac{8}{3} = -\frac{2}{3} < 0$$

$$\text{Class}(x_5) = N \quad \text{Class Predicted}(x_5) = N$$

$$S_{nn}(x_6) = [x_5, x_1, x_2, x_4, x_8]$$

$$WA(x_6) = \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (+1) + \left(2 + \frac{2}{3}\right) \times (-1) = \frac{6}{3} - \frac{8}{3} = -\frac{2}{3} < 0$$

$$\text{Class}(x_6) = N \quad \text{Class Predicted}(x_6) = N$$

$$S_{nn}(x_7) = [x_3, x_1, x_2, x_4, x_8]$$

$$WA(x_7) = \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (+1) + \frac{2}{3} \times (-1) = \frac{8}{3} - \frac{2}{3} = 2 > 0$$

$$\text{Class}(x_7) = N \quad \text{Class Predicted}(x_7) = P$$

$$S_{nn}(x_8) = [x_2, x_3, x_5, x_6, x_7]$$

$$WA(x_8) = \left(2 + \frac{2}{3}\right) \times (+1) + \left(\frac{2}{3} + \frac{2}{3} + \frac{2}{3}\right) \times (-1) = \frac{8}{3} - \frac{6}{3} = \frac{2}{3} > 0$$

$$\text{Class}(x_8) = N \quad \text{Class Predicted}(x_8) = P$$

## Leave One Out Execution

$$D_{train} = \{x_2, \dots, x_8\} \quad D_{test} = \{x_1\}$$

target prediction = P

target true = P

$$D_{train} = \{x_1, x_2, x_4, \dots, x_8\} \quad D_{test} = \{x_3\}$$

target prediction = N

target true = P

$$D_{train} = \{x_1, \dots, x_4, x_6, \dots, x_8\} \quad D_{test} = \{x_5\}$$

target prediction = N

target true = N

$$D_{train} = \{x_1, x_3, \dots, x_8\} \quad D_{test} = \{x_2\}$$

target prediction = N

target true = P

$$D_{train} = \{x_1, \dots, x_5, x_7, x_8\} \quad D_{test} = \{x_6\}$$

target prediction = P

target true = P

$$D_{train} = \{x_1, \dots, x_5, x_7, x_8\} \quad D_{test} = \{x_6\}$$

target prediction = N

target true = N

$$D_{train} = \{x_1, \dots, x_6, x_8\} \quad D_{test} = \{x_7\}$$

target prediction = P

target true = N

predictions = [P, N, N, P, N, N, P, P]

truths = [P, P, P, P, N, N, N, N]

Confusion Matrix

		truth	
		P	N
predicted	P	2	2
	N	2	2

$$TP = 2$$

$$FP = 2$$

$$TF = 2$$

$$FN = 2$$

$$\text{recall} = \frac{TP}{TP+FN} = \frac{2}{2+2} = \frac{1}{2} = 0.5$$

2.

(2)

$$P(y|x) = \frac{P(x|y)P(y)}{P(x)}$$

↙  
Posterior      ↙  
all taken from data directly

	$y_1$	$y_2$	Class		$y_3$	Class
$x_1$	A	0	P		1,2	P
$x_2$	B	1	P		0,8	P
$x_3$	A	1	P		0,5	P
$x_4$	A	0	P		0,9	P
$x_5$	B	0	N		1	N
$x_6$	B	0	N		0,9	N
$x_7$	A	1	N		1,2	N
$x_8$	B	1	N		0,8	N
$x_9$	B	0	P		0,8	P

→  $y_1$  and  $y_2$  are dependent

→  $y_3$  is independent from  $\{y_1, y_2\}$

→  $y_3$  is normally distributed

P → Positive    N → Negative

$$P(P) = \frac{\#(\text{Class} = P)}{\#\text{total}} = \frac{5}{9} \quad P(N) = \frac{\#(\text{Class} = N)}{\#\text{total}} = \frac{4}{9}$$

$$P(y_1=A, y_2=0 | P) = \frac{\#(y_1=A, y_2=0, \text{Class} = P)}{\#(\text{Class} = P)} = \frac{2}{5}$$

$$P(y_1=A, y_2=1 | P) = \frac{\#(y_1=A, y_2=1, \text{Class} = P)}{\#(\text{Class} = P)} = \frac{1}{5}$$

$$P(y_1=B, y_2=0 | P) = \frac{\#(y_1=B, y_2=0, \text{Class} = P)}{\#(\text{Class} = P)} = \frac{1}{5}$$

$$P(y_1=B, y_2=1 | P) = \frac{\#(y_1=B, y_2=1, \text{Class}=P)}{\#(\text{Class}=P)} = \frac{1}{5}$$

$$P(y_1=A, y_2=0 | N) = \frac{\#(y_1=A, y_2=0, \text{Class}=N)}{\#(\text{Class}=N)} = 0$$

$$P(y_1=A, y_2=1 | N) = \frac{\#(y_1=A, y_2=1, \text{Class}=N)}{\#(\text{Class}=N)} = \frac{1}{4}$$

$$P(y_1=B, y_2=0 | N) = \frac{\#(y_1=B, y_2=0, \text{Class}=N)}{\#(\text{Class}=N)} = \frac{2}{4} = \frac{1}{2}$$

$$P(y_1=B, y_2=1 | N) = \frac{\#(y_1=B, y_2=1, \text{Class}=N)}{\#(\text{Class}=N)} = \frac{1}{4}$$

$$P(x | \mu, \sigma^2, p) \quad \mu_p = \frac{1,2 + 0,8 + 0,5 + 0,9 + 0,8}{5} = \frac{4,2}{5} = 0,84$$

$$P(x | \mu, \sigma^2, N) \quad \mu_N = \frac{1 + 0,9 + 1,2 + 0,8}{4} = \frac{3,9}{4} = 0,975$$

$$\sigma_p^2 = \frac{1}{N-1} \sum_{i=0}^N (y_{pi} - \mu_p)^2 = \frac{1}{5-1} \left[ (1.2 - 0.84)^2 + (0.8 - 0.84)^2 + (0.5 - 0.84)^2 + (0.9 - 0.84)^2 + (0.8 - 0.84)^2 \right] =$$

$$= \frac{1}{4} (0.1296 + 0.0016 + 0.0156 + 0.0036 + 0.0016) = \frac{1}{4} (0.252) = 0.063$$

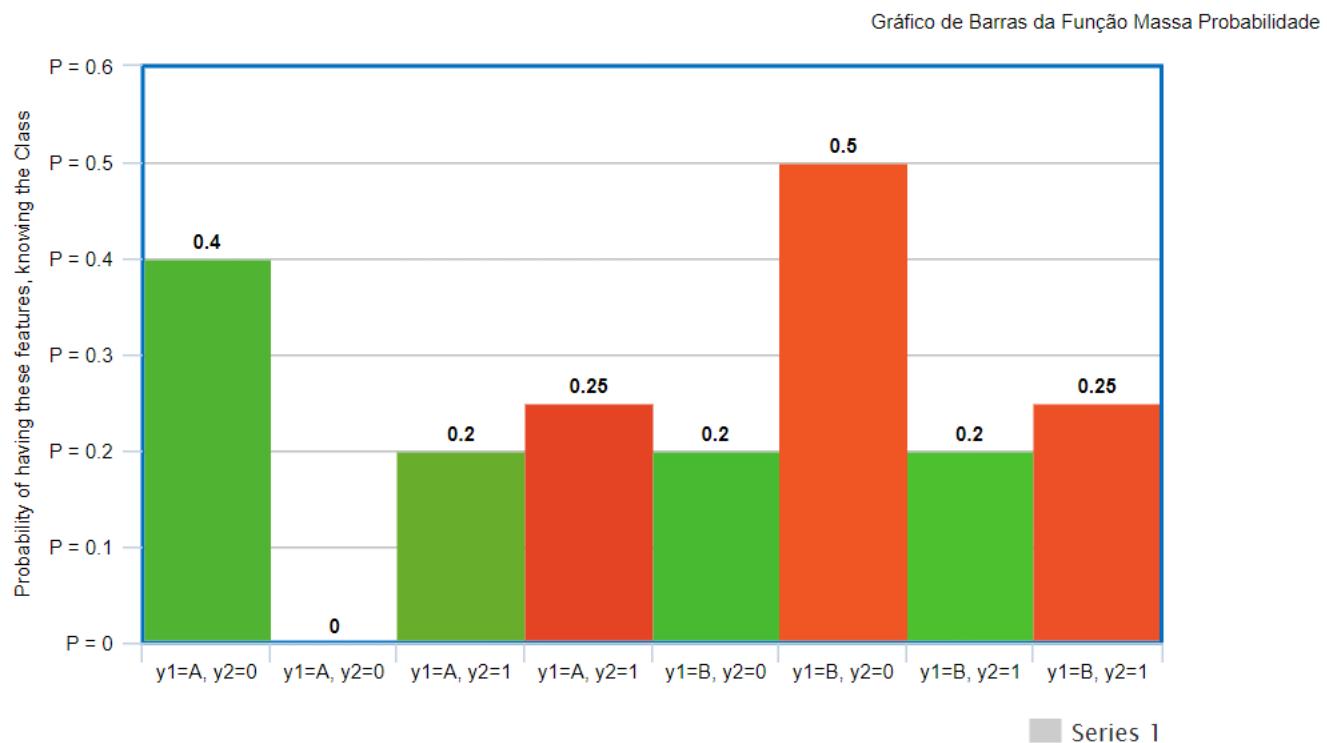
$$\sigma_N^2 = \frac{1}{N-1} \sum_{i=0}^N (y_{Ni} - \mu_N)^2 = \frac{1}{4-1} \left[ (1 - 0.975)^2 + (0.9 - 0.975)^2 + (1.2 - 0.975)^2 + (0.8 - 0.975)^2 \right]$$

$$= \frac{1}{3} [0.000625 + 0.005625 + 0.050625 + 0.030625] =$$

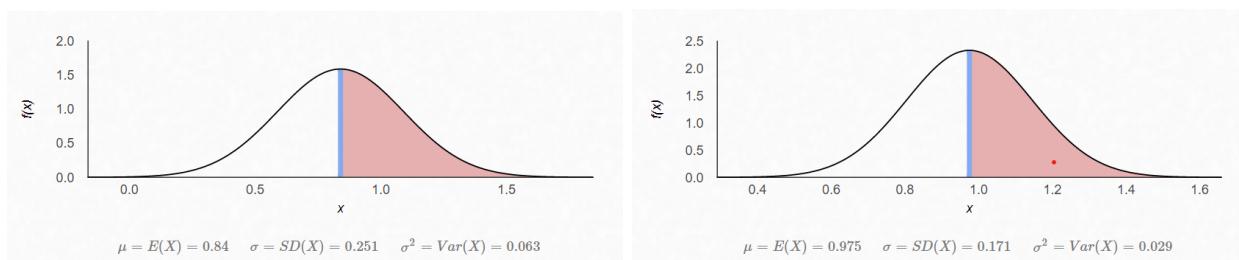
$$= \frac{1}{3} \times 0.0875 \approx 0.0292$$

$$P(Y_3=x | \mu=0.84, \sigma^2=0.063, P) = \frac{1}{\sqrt{2\pi \times 0.063}} e^{-\frac{(x-0.84)^2}{2 \times 0.063}}$$
$$P(Y_3=x | \mu=0.975, \sigma^2=0.0292, N) = \frac{1}{\sqrt{2\pi \times 0.0292}} e^{-\frac{(x-0.975)^2}{2 \times 0.0292}}$$

The probability mass function for the Bayesian classifier, the red boxes mean  $P(y_1, y_2|N)$  and the green mean  $P(y_1, y_2|P)$ , for  $y_1$  and  $y_2$  are:



The normal distributions of  $P(y_3|P)$  and  $P(y_3|N)$  are, respectively:



3.

$$\textcircled{3} \quad P(P|Y_1, Y_2, Y_3) = \frac{P(Y_1, Y_2, Y_3 | P) \times P(P)}{P(Y_1, Y_2, Y_3)} = \frac{P(Y_1, Y_2 | P) \times P(Y_3 | P) \times P(P)}{P(Y_1, Y_2, Y_3)}$$

observação 1: (A, 1, 0.8)

$$P(\text{Positive} | Y_1=A, Y_2=1, Y_3=0.8) = P(\text{Positive}) \times \frac{P(Y_1=A, Y_2=1, Y_3=0.8 | P)}{P(Y_1=A, Y_2=1, Y_3=0.8)}$$

$$P(\text{Positive}) = \frac{5}{9}$$

$$P(Y_1=A, Y_2=1 | P) = 0.2 \quad P(Y_1=A, Y_2=1 | N) = 0.25$$

$$P(Y_1, Y_2, Y_3) = P(Y_1, Y_2 | P) \times P(Y_3 | P) \times P(P) + P(Y_1, Y_2 | N) \times P(Y_3 | N) \times P(N)$$

$$= 0.2 \times 1.5694 \times \frac{5}{9} + 0.25 \times 1.3819 \times \frac{4}{9} = 0.19438 + 0.15354 = 0.32792$$

$$P(P | Y_1, Y_2, Y_3) = \frac{0.2 \times 0.8719}{0.32792} = 0.53178$$

observação 2: (B, 1, 1)

$$P(\text{Positive} | Y_1=B, Y_2=1, Y_3=1) = \frac{P(Y_1, Y_2, Y_3 | P) \times P(P)}{P(Y_1, Y_2, Y_3)}$$

$$P(\text{Positive}) = \frac{5}{9}$$

$$P(Y_1=B, Y_2=1 | P) = 0.2 \quad P(Y_1=B, Y_2=1 | N) = 0.25$$

$$P(Y_1, Y_2, Y_3) = P(Y_1, Y_2 | P) \times P(Y_3 | P) \times P(P) + P(Y_1, Y_2 | N) \times P(Y_3 | N) \times P(N)$$

$$= 0.2 \times 1.2972 \times \frac{5}{9} + 0.25 \times 2.3095 \times \frac{4}{9} = 0.14413 + 0.25661 = 0.40074$$

$$P(P | Y_1, Y_2, Y_3) = \frac{0.2 \times 1.2972 \times \frac{5}{9}}{0.40074} = 0.3597$$

observación 3: (B, 0, 0.9)

$$P(P|y_1, y_2, y_3) = \frac{P(y_1, y_2, y_3 | P) \times P(P)}{P(y_1, y_2, y_3)}$$

$$P(y_1=B, y_2=0 | P) = 0.2 \quad P(y_1=B, y_2=0 | N) = 0.5$$

$$\begin{aligned} P(y_1, y_2, y_3) &= P(y_1, y_2 | P) \times P(y_3 | P) \times P(P) + P(y_1, y_2 | N) \times P(y_3 | N) \times P(N) \\ &= 0.2 \times 1.5446 \times \frac{5}{9} + 0.5 \times 2.1201 \times \frac{4}{9} = \\ &= 0.17162 + 0.47113 = 0.64275 \end{aligned}$$

$$P(P|y_1, y_2, y_3) = \frac{P(y_1, y_2 | P) \times P(y_3 | P) \times P(P)}{P(y_1, y_2, y_3)} = \frac{0.2 \times 1.5446 \times \frac{5}{9}}{0.64275} = 0.2670$$

4.

(4)

$$1 - 0.53178$$

values of  $P(\text{Positive} | x)$  for observations: 2 - 0.3597

$$3 - 0.2670$$

Classification of observations w/ different thresholds:

	$\theta = 0.3$	$\theta = 0.5$	$\theta = 0.7$	actual Class
obs. 1	P	P	N	P
obs. 2	P	N	N	P
obs. 3	N	N	N	N

$$\text{accuracy} = \frac{TP + TN}{\text{All}}$$

$$\theta = 0.3$$

$$TP = 2 \quad TN = 1 \quad accuracy = \frac{2+1}{3} = 100\% \\ All = 3$$

$$\theta = 0.5$$

$$TP = 1 \quad TN = 1 \quad FN = 1 \quad accuracy = \frac{1+1}{3} = 66,7\% \\ All = 3$$

$$\theta = 0.7$$

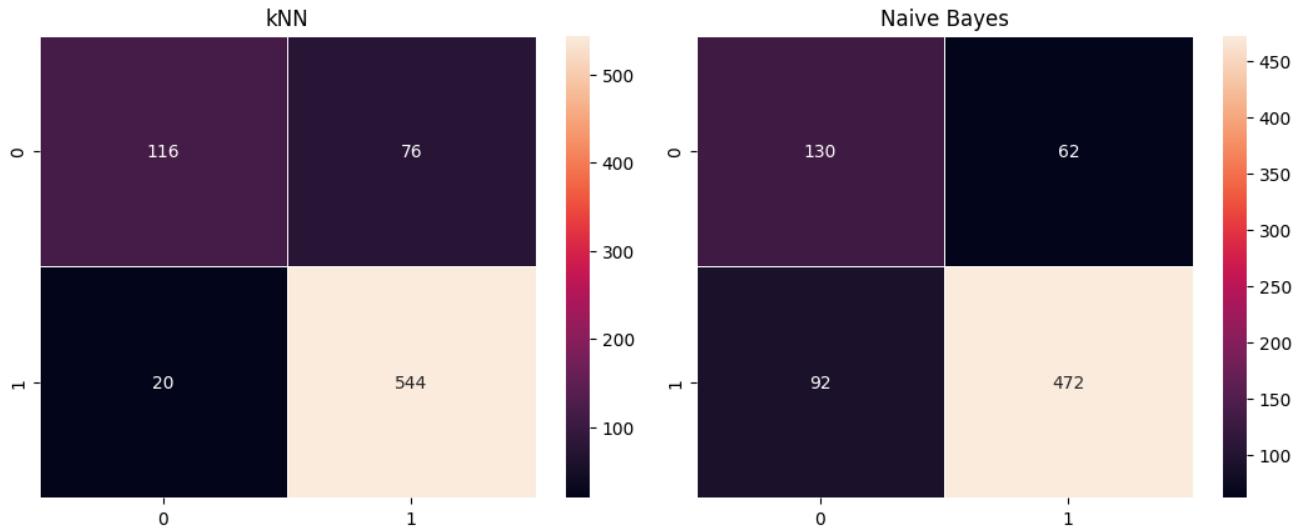
$$FN = 2 \quad TN = 1 \quad accuracy = \frac{1}{3} = 33,3\%$$

$$accuracy_{\theta=0.3} > accuracy_{\theta=0.5} > accuracy_{\theta=0.7}$$

The accuracy for  $\theta = 0.3$  is the biggest, 100%, for our previous observations, the threshold that optimizes accuracy is  $\theta = 0.3$

## Part II: Programming

5.



6.  $knn_{acc} > gnb_{acc}$ ?  $pval = 0.001316817828490826$

$knn_{acc} < gnb_{acc}$ ?  $pval = 0.9986831821715092$

$knn_{acc} \neq gnb_{acc}$ ?  $pval = 0.002633635656981652$

Through the  $pval$  we can observe that the accuracy of  $knn$  is indeed statistically superior to Naive Bayes' accuracy, with a  $pval$  of less than 0.05 ( $0.001 < 0.05$ ), which indicates strong statistical significance, so, our hypothesis is true.

7. Since, Naive Bayes is a parametric algorithm, it assumes a functional form, Gaussian distribution in this case, our data set might be non-Gaussian, hence, the poorer performance compared to  $knn$ . On a problem that focuses on similarity between observations,  $knn$  has an overall better performance due to its nature to locally optimize, although, outliers can worsen its performance significantly, hence, the importance of adjusting  $k$  to maximise performance.

## Appendix

```
1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 from scipy import stats
5 from scipy.io.arff import loadarff
6 from sklearn.model_selection import StratifiedKFold
7 from sklearn.neighbors import KNeighborsClassifier
8 from sklearn.metrics import accuracy_score, confusion_matrix
9 from sklearn.preprocessing import StandardScaler
10 from sklearn.naive_bayes import GaussianNB
11
12 cm_sum_knn = [[0, 0], [0, 0]]
13 cm_sum_gnb = [[0, 0], [0, 0]]
14 knn_acc = []
15 gnb_acc = []
16 data = loadarff('pd_speech.arff')
17 df = pd.DataFrame(data[0])
18 df['class'] = df['class'].str.decode('utf-8')
19 X = df.drop('class', axis=1)
20 y = df['class']
21
22 skf = StratifiedKFold(n_splits = 10, shuffle = True, random_state = 0)
23 knn = KNeighborsClassifier()
24 gnb = GaussianNB()
25
26 for train_index, test_index in skf.split(X, y):
27     X_train, X_test, y_train, y_test = X.iloc[train_index], X.iloc[test_index],\
28                                         y.iloc[train_index], y.iloc[test_index]
29     scaler = StandardScaler().fit(X_train)
30     X_train, X_test = scaler.transform(X_train), scaler.transform(X_test)
31
32
33     knn.fit(X_train, y_train)
34     y_pred_knn = knn.predict(X_test)
35     cm_knn = confusion_matrix(y_test, y_pred_knn)
36     knn_acc.append(accuracy_score(y_test, y_pred_knn))
37
38     for i in range(2):
39         for j in range(2):
40             cm_sum_knn[i][j] += cm_knn[i][j]
41
42     gnb.fit(X_train, y_train)
43     y_pred_gnb = gnb.predict(X_test)
44     cm_gnb = confusion_matrix(y_test, y_pred_gnb)
45     gnb_acc.append(accuracy_score(y_test, y_pred_gnb))
46
47     for i in range(2):
48         for j in range(2):
49             cm_sum_gnb[i][j] += cm_gnb[i][j]
50
51 f, ax = plt.subplots()
52 sns.heatmap(cm_sum_knn, annot = True, linewidths= 0.5, fmt="g")
53 plt.title('kNN')
54 plt.show()
```

```
55
56 sns.heatmap(cm_sum_gnb ,annot = True, linewidths= 0.5, fmt="g")
57 plt.title('Naive Bayes')
58 plt.show()
59
60 res = stats.ttest_rel(knn_acc, gnb_acc, alternative='greater')
61 print("knn>gnb? pval=",res.pvalue)
62
63 res = stats.ttest_rel(knn_acc, gnb_acc, alternative='less')
64 print("knn<gnb? pval=",res.pvalue)
65
66 res = stats.ttest_rel(knn_acc, gnb_acc, alternative='two-sided')
67 print("knn!=gnb? pval=",res.pvalue)
```