

I Pen-and-paper

1.

	y_1	y_2	Class	Classification
x_1	A	0	P	P
x_2	B	1	P	N
x_3	A	1	P	N
x_4	A	0	P	P
x_5	B	0	N	-
x_6	B	0	N	-
x_7	A	1	N	-
x_8	B	1	N	-

x_1)

x_i	x_2	x_3	x_4	x_5	x_6	x_7	x_8
$d(x_1, x_i)$	2.5	1.5	0.5	1.5	1.5	1.5	2.5
Class	P	P	P	N	N	N	N

neighbours = $\{x_3, x_4, x_5, x_6, x_7\}$
P: $\frac{1}{1.5} + \frac{1}{0.5} = 2, \bar{6}$
N: $\frac{1}{1.5} + \frac{1}{1.5} + \frac{1}{1.5} = 2$
 $\hat{z}_{x_1} = \text{P} \implies x_1$ is a true positive.

x_2)

x_i	x_1	x_3	x_4	x_5	x_6	x_7	x_8
$d(x_2, x_i)$	2.5	1.5	2.5	1.5	1.5	1.5	0.5
Class	P	P	P	N	N	N	N

neighbours = $\{x_3, x_5, x_6, x_7, x_8\}$
P: $\frac{1}{1.5} = 0, \bar{6}$
N: $\frac{1}{1.5} + \frac{1}{1.5} + \frac{1}{1.5} + \frac{1}{0.5} = 4$
 $\hat{z}_{x_2} = \text{N} \implies x_2$ is a false negative.

x_3)

x_i	x_1	x_2	x_4	x_5	x_6	x_7	x_8
$d(x_3, x_i)$	1.5	1.5	1.5	2.5	2.5	0.5	1.5
Class	P	P	P	N	N	N	N

neighbours = $\{x_1, x_2, x_4, x_7, x_8\}$
P: $\frac{1}{1.5} + \frac{1}{1.5} + \frac{1}{1.5} = 2$
N: $\frac{1}{0.5} + \frac{1}{1.5} = 2, \bar{6}$
 $\hat{z}_{x_3} = \text{N} \implies x_3$ is a false negative.

x_4)

x_i	x_1	x_2	x_3	x_5	x_6	x_7	x_8
$d(x_4, x_i)$	0.5	2.5	1.5	1.5	1.5	1.5	2.5
Class	P	P	P	N	N	N	N

neighbours = $\{x_1, x_3, x_5, x_6, x_7\}$
P: $\frac{1}{0.5} + \frac{1}{1.5} = 2, \bar{6}$
N: $\frac{1}{1.5} + \frac{1}{1.5} + \frac{1}{1.5} = 2$
 $\hat{z}_{x_4} = \text{P} \implies x_4$ is a true positive.

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

2. $y_1 \mid P = \{A, A, A, B, B\}$
 $y_1 \mid N = \{A, B, B, B\}$
 $y_2 \mid P = \{0, 0, 0, 1, 1\}$
 $y_2 \mid N = \{0, 0, 1, 1\}$
 $y_3 \mid P = \{1.2, 0.8, 0.5, 0.9, 0.8\}$
 $y_3 \mid N = \{1, 0.9, 1.2, 0.8\}$

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

$$P(\text{class} = N) = 1 - P(\text{class} = P) = 1 - \frac{5}{9} = \frac{4}{9}$$

$$P(y_1 = A, y_2 = 0) = \frac{2}{9}$$

$$P(y_1 = A, y_2 = 1) = \frac{2}{9}$$

$$P(y_1 = B, y_2 = 0) = \frac{3}{9}$$

$$P(y_1 = B, y_2 = 1) = \frac{2}{9}$$

$$P(y_1 = A, y_2 = 0 \mid \text{class} = P) = \frac{2}{5}$$

$$P(y_1 = A, y_2 = 1 \mid \text{class} = P) = \frac{1}{5}$$

$$P(y_1 = B, y_2 = 0 \mid \text{class} = P) = \frac{1}{5}$$

$$P(y_1 = B, y_2 = 1 \mid \text{class} = P) = \frac{1}{5}$$

$$P(y_1 = A, y_2 = 0 \mid \text{class} = N) = \frac{0}{4}$$

$$P(y_1 = A, y_2 = 1 \mid \text{class} = N) = \frac{1}{4}$$

$$P(y_1 = B, y_2 = 0 \mid \text{class} = N) = \frac{2}{4}$$

$$P(y_1 = B, y_2 = 1 \mid \text{class} = N) = \frac{1}{4}$$

$$P(y_3) = \mathcal{N}(y_3 \mid \mu, \sigma^2), \text{ where } \mu = 0.9 \text{ and } \sigma^2 = 0.0475$$

$$P(y_3 \mid \text{class} = P) = \mathcal{N}(y_3 \mid \mu_P, \sigma_P^2), \text{ where } \mu_P = 0.84 \text{ and } \sigma_P^2 = 0.063$$

$$P(y_3 \mid \text{class} = N) = \mathcal{N}(y_3 \mid \mu_N, \sigma_N^2), \text{ where } \mu_N = 0.975 \text{ and } \sigma_N^2 = 0.02916667$$

$$\mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \times e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$$\mu = \frac{\sum_{i=1}^n x_i}{n}$$

$$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

3.

$$\begin{aligned}
P(\text{class} = P \mid y_1 = A, y_2 = 1, y_3 = 0.8) &= \\
&= \frac{P(\text{class} = P) \times P(y_1 = A, y_2 = 1, y_3 = 0.8 \mid \text{class} = P)}{P(y_1 = A, y_2 = 1, y_3 = 0.8)} = \\
&= \frac{P(\text{class} = P) \times P(y_1 = A, y_2 = 1 \mid \text{class} = P) \times P(y_3 = 0.8 \mid \text{class} = P)}{P(y_1 = A, y_2 = 1) \times P(y_3 = 0.8)} = \\
&= \frac{\frac{5}{9} \times \frac{1}{5} \times 1.569369}{\frac{2}{9} \times 1.647586} = 0.4762631
\end{aligned}$$

$$\begin{aligned}
P(\text{class} = P \mid y_1 = B, y_2 = 1, y_3 = 1) &= \\
&= \frac{P(\text{class} = P) \times P(y_1 = B, y_2 = 1, y_3 = 1 \mid \text{class} = P)}{P(y_1 = B, y_2 = 1, y_3 = 1)} = \\
&= \frac{P(\text{class} = P) \times P(y_1 = B, y_2 = 1 \mid \text{class} = P) \times P(y_3 = 1 \mid \text{class} = P)}{P(y_1 = B, y_2 = 1) \times P(y_3 = 1)} = \\
&= \frac{\frac{5}{9} \times \frac{1}{5} \times 1.297186}{\frac{2}{9} \times 1.647586} = 0.3936626
\end{aligned}$$

$$\begin{aligned}
P(\text{class} = P \mid y_1 = B, y_2 = 0, y_3 = 0.9) &= \\
&= \frac{P(\text{class} = P) \times P(y_1 = B, y_2 = 0, y_3 = 0.9 \mid \text{class} = P)}{P(y_1 = B, y_2 = 0, y_3 = 0.9)} = \\
&= \frac{P(\text{class} = P) \times P(y_1 = B, y_2 = 0 \mid \text{class} = P) \times P(y_3 = 0.9 \mid \text{class} = P)}{P(y_1 = B, y_2 = 0) \times P(y_3 = 0.9)} = \\
&= \frac{\frac{5}{9} \times \frac{1}{5} \times 1.544655}{\frac{3}{9} \times 1.830473} = 0.2812852
\end{aligned}$$

4. For $\theta = 0.3$:

$(y_1 = A, y_2 = 1, y_3 = 0.8)$ is classified as positive.

$(y_1 = B, y_2 = 1, y_3 = 1)$ is classified as positive.

$(y_1 = B, y_2 = 0, y_3 = 0.9)$ is classified as negative.

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{All}} = \frac{3}{3} = 1$$

For $\theta = 0.5$:

$(y_1 = A, y_2 = 1, y_3 = 0.8)$ is classified as negative

$(y_1 = B, y_2 = 1, y_3 = 1)$ is classified as negative

$(y_1 = B, y_2 = 0, y_3 = 0.9)$ is classified as negative

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{All}} = \frac{1}{3} = 0.\bar{3}$$

For $\theta = 0.7$:

$(y_1 = A, y_2 = 1, y_3 = 0.8)$ is classified as negative

$(y_1 = B, y_2 = 1, y_3 = 1)$ is classified as negative

$(y_1 = B, y_2 = 0, y_3 = 0.9)$ is classified as negative

$$\text{accuracy} = \frac{\text{TP} + \text{TN}}{\text{All}} = \frac{1}{3} = 0.\bar{3}$$

Between the given values, the $\theta = 0.3$ decision threshold optimizes testing accuracy.

II Programming

5.

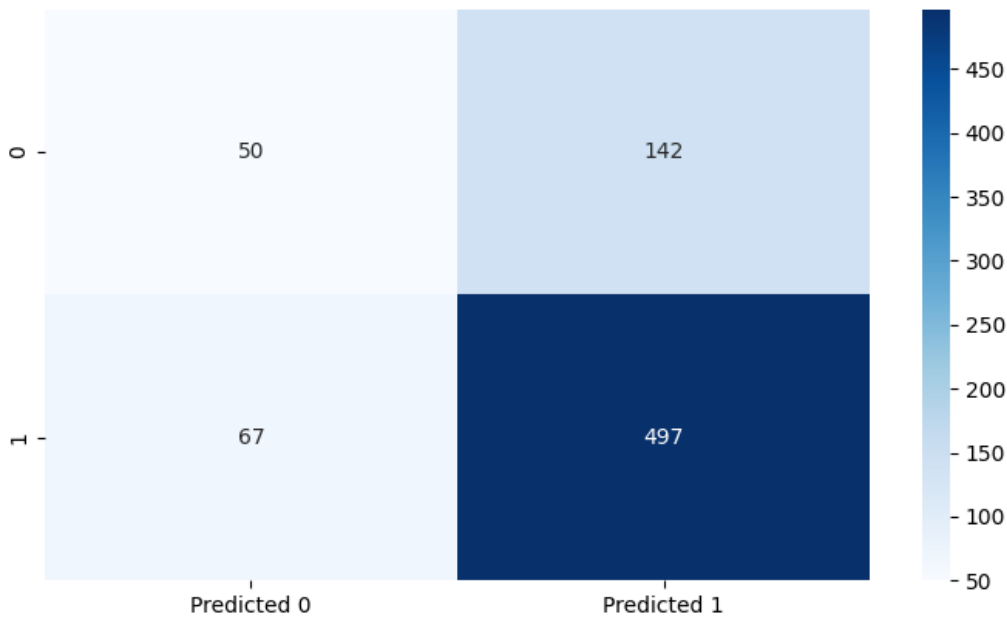


Figure 1: kNN Confusion Matrix

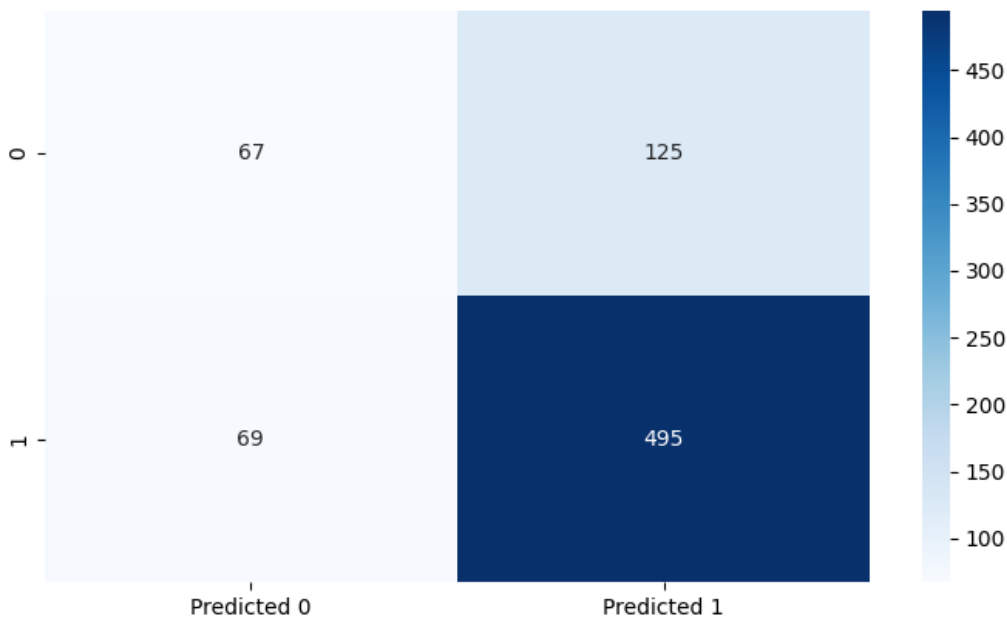


Figure 2: Naïve Bayes Gaussian Confusion Matrix

6. Defining H_0 as “*kNN is statistically superior to Naïve Bayes regarding accuracy*” and H_1 as “*kNN is not statistically superior to Naïve Bayes regarding accuracy*”, the hypothesis test lead to a p-value of 0.9104476998751558. Therefore, we reject H_0 , concluding that kNN is not statistically superior to Naïve Bayes regarding accuracy.
7. The Naïve Bayes classifier may have had better accuracy results due to the small number of neighbours, k , of the kNN classifier, which is known to have its accuracy increased as the value of k rises. The reduced dependence between the variables of the dataset favors the Naïve Bayes classifier, which assumes variable independence, increasing this classifier’s accuracy. Furthermore, the Naïve Bayes classifier is known to have better accuracy, compared to the kNN classifier, when applied to big data.

III Appendix

```

1 from scipy.io.arff import loadarff
2 from sklearn.metrics import confusion_matrix
3 from sklearn.model_selection import StratifiedKFold
4 from sklearn.naive_bayes import GaussianNB
5 from sklearn.neighbors import KNeighborsClassifier
6 import matplotlib.pyplot as plt
7 import numpy as np
8 import pandas as pd
9 import seaborn as sns
10 from scipy import stats
11
12 if __name__ == "__main__":
13     data = loadarff("pd_speech.arff")
14     df = pd.DataFrame(data[0])
15     df["class"] = df["class"].astype(int)
16     X, y = df.drop("class", axis=1), df["class"]
17
18     classifiers = [KNeighborsClassifier(n_neighbors=5), GaussianNB()]
19     matrices = [np.zeros((2, 2)), np.zeros((2, 2))]
20     skf = StratifiedKFold(n_splits=10, random_state=0, shuffle=True)
21     acc = [[], []]
22
23     for train_index, test_index in skf.split(X, y):
24         X_train, X_test = X.iloc[train_index], X.iloc[test_index]
25         y_train, y_test = y.iloc[train_index], y.iloc[test_index]
26
27         for i in range(len(classifiers)):
28             classifiers[i].fit(X_train, y_train)
29             y_pred = classifiers[i].predict(X_test)
30             matrices[i] += confusion_matrix(y_test, y_pred)
31             acc[i].append(classifiers[i].score(X_test, y_test))
32
33     for matrix in matrices:
34         matrix = pd.DataFrame(matrix,
35                               index=["0", "1"],
36                               columns=["Predicted 0", "Predicted 1"])
37         sns.heatmap(matrix, annot=True, fmt="g", cmap="Blues")
38         plt.show()
39
40     print(
41         "kNN is statistically superior to Naive Bayes regarding accuracy is
42         suported by a p-value of {}."
43         .format(stats.ttest_rel(acc[0], acc[1], alternative="greater").pvalue))

```