

I. Pen-and-paper [13v]

1)

	x1	x2	x3	x4	x5	x6	x7	x8
x1	–	2.5	1.5	0.5	1.5	1.5	1.5	2.5
x2		–	1.5	2.5	1.5	1.5	1.5	0.5
x3			–	1.5	2.5	2.5	0.5	1.5
x4				–	1.5	1.5	1.5	0.5
x5					–	0.5	2.5	1.5
x6						–	2.5	1.5
x7							–	1.5
x8								–

$$w_{ij} = \frac{1}{d(x_i, x_j)}$$

$$5NN(x_1) = (x_3, x_4, x_5, x_6, x_7) = \operatorname{argmax}((\frac{2}{3} + 2)P, (3 \cdot \frac{2}{3})N) = P$$

$$5NN(x_2) = (x_3, x_5, x_6, x_7, x_8) = \operatorname{argmax}(\frac{2}{3}P, (3 \cdot \frac{2}{3} + 2)N) = N$$

$$5NN(x_3) = \dots\dots\dots = N$$

$$5NN(x_4) = \dots\dots\dots = P$$

$$5NN(x_5) = \dots\dots\dots = N$$

$$5NN(x_6) = \dots\dots\dots = N$$

$$5NN(x_7) = \dots\dots\dots = P$$

$$5NN(x_8) = \dots\dots\dots = P$$

Reais \ Previstos	P	N
P	2	2
N	2	2

$$Recall_P = \frac{TP}{TP + FN} = \frac{2}{2 + 2} = 0.5 = Recall_N \text{ (porque os valores são todos iguais)}$$

2) A=0, B=1

Para $z = P, y_1 y_2 \{$

$$\mu = \frac{1}{5} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} \right) = \begin{bmatrix} 0.4 \\ 0.4 \end{bmatrix} \text{ (not a fraction, can't do "vertical" vector)}$$

$$\text{var}(y_1) = \frac{1}{4} \sum_{i=1}^5 (y_{1i} - \mu)^2 = \dots = 0.1812 ; \text{var}(y_2) = \dots = 0.1812$$

$$\text{cov}(y_1, y_2) = \frac{1}{4} \sum_{i=1}^5 (y_{1i} - \mu)(y_{2i} - \mu) = \dots = 0.1098$$

$$\Sigma = \begin{bmatrix} 0.1812 & 0.1098 \\ 0.1098 & 0.1812 \end{bmatrix} ; |\Sigma| = 0.02077 ; \Sigma^{-1} = \frac{1}{0.02077} \Sigma = \begin{bmatrix} 8.72 & -5.285 \\ -5.285 & 8.72 \end{bmatrix}$$

$$\mathbf{P}(y_1, y_2 | z=P) = \frac{1}{2\pi\sqrt{0.02077}} \exp\left(-\frac{1}{2} [y_1 - 0.4, y_2 - 0.4] \begin{bmatrix} 8.72 & -5.285 \\ -5.285 & 8.72 \end{bmatrix} \begin{bmatrix} y_1 - 0.4 \\ y_2 - 0.4 \end{bmatrix}\right)$$

$\}$
 $z = N, y_1 y_2 \{$

$$\mu = \frac{1}{4} \left(\begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right) = \begin{bmatrix} 0.75 \\ 0.5 \end{bmatrix}$$

$$\text{var}(y_1) = \frac{1}{3} \sum_{i=1}^4 (y_{1i} - \mu)^2 = \dots = 0.25 ; \text{var}(y_2) = \dots = \frac{1}{3}$$

$$\text{cov}(y_1, y_2) = \frac{1}{3} \sum_{i=1}^4 (y_{1i} - \mu)(y_{2i} - \mu) = \dots = -\frac{1}{6}$$

$$\Sigma = \begin{bmatrix} 0.25 & -\frac{1}{6} \\ -\frac{1}{6} & \frac{1}{3} \end{bmatrix} ; |\Sigma| = \frac{1}{18} ; \Sigma^{-1} = 18 * \Sigma = \begin{bmatrix} 6 & 3 \\ 3 & 4.5 \end{bmatrix}$$

$$\mathbf{P}(y_1, y_2 | z=N) = \frac{1}{2\pi\sqrt{\frac{1}{18}}} \exp\left(-\frac{1}{2} [y_1 - 0.75, y_2 - 0.5] \begin{bmatrix} 6 & 3 \\ 3 & 4.5 \end{bmatrix} \begin{bmatrix} y_1 - 0.75 \\ y_2 - 0.5 \end{bmatrix}\right)$$

$\}$

$z = P, y_3 \{$

$$\mu = \frac{1}{5} (1 + 0.9 + 1.2 + 0.8) = 0.84$$

$$\sigma^2 = \frac{1}{4} \sum_{i=1}^5 (y_{3i} - \mu)^2 = \dots \approx 0.063$$

$$\mathbf{P}(y_3 | z = P) = \frac{1}{\sqrt{2\pi*0.063}} \exp\left(-\frac{1}{2*0.063} (y_3 - 0.84)^2\right)$$

$\}$

$z = N, y_3 \{$

$$\mu = \frac{1}{4} (1 + 0.9 + 1.2 + 0.8) = 0.975$$

$$\sigma^2 = \frac{1}{3} \sum_{i=1}^4 (y_{3i} - \mu)^2 = \dots \approx 0.029$$

$$\mathbf{P}(y_3 | z = N) = \frac{1}{\sqrt{2\pi*0.029}} \exp\left(-\frac{1}{2*0.029} (y_3 - 0.975)^2\right)$$

$\}$

3) (usando as expressões obtidas na alínea anterior)

$$p(c = P | x_{new1}) = \frac{p(y_1=A, y_2=1 | c=P) * p(y_3=0.3 | c=P) * p(c=P)}{p(y_1=A, y_2=1 | c=P) * p(y_3=0.8 | c=P) * p(c=P) + p(y_1=A, y_2=1 | c=N) * p(y_3=0.8 | c=N) * p(c=N)} \approx 0.0618$$

$$p(c = P | x_{new2}) = \frac{p(y_1=B, y_2=1 | c=P) * p(y_3=1 | c=P) * p(c=P)}{p(y_1=B, y_2=1 | c=P) * p(y_3=1 | c=P) * p(c=P) + p(y_1=B, y_2=1 | c=N) * p(y_3=1 | c=N) * p(c=N)} \approx 0.29104$$

$$p(c = P | x_{new3}) = \frac{p(y_1=B, y_2=0 | c=P) * p(y_3=0.9 | c=P) * p(c=P)}{p(y_1=B, y_2=0 | c=P) * p(y_3=0.9 | c=P) * p(c=P) + p(y_1=B, y_2=0 | c=N) * p(y_3=0.9 | c=N) * p(c=N)} \approx 0.00439$$

II. Programming and critical analysis [7v]

```
# Import Wall
import pandas as pd
import numpy as np
from sklearn import metrics, datasets, tree
from sklearn.metrics import confusion_matrix
from sklearn.metrics import accuracy_score
from sklearn.model_selection import StratifiedKFold
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
import matplotlib.pyplot as plt
from scipy.io.arff import loadarff
import seaborn as sns
from scipy import stats

data = loadarff('pd_speech.arff')
df = pd.DataFrame(data[0])
df['class'] = df['class'].str.decode('utf-8')

X = df.iloc[:, :-1]
Y = df.iloc[:, -1]
inputs = df.drop('class', axis=1).values
outputs = df['class'].values

# stratified k-fold cross validation
skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=0)

# predict using Knn
predictorKnn = KNeighborsClassifier(n_neighbors=5, metric='euclidean',
weights='uniform')
# predict using gaussian naive bayes
predictorGnb = GaussianNB()

cm = None
m = []
accuracyKnn = []
accuracyNB = []
for train, test in skf.split(X, Y):
    X_train, X_test = X.iloc[train], X.iloc[test]
    Y_train, Y_test = Y.iloc[train], Y.iloc[test]
    predictorKnn.fit(X_train, Y_train)
```

```

Y_pred = predictorKnn.predict(X_test)
accuracyKnn.append(accuracy_score(Y_test, Y_pred))
if type(cm) == type(None):
    cm = np.array(confusion_matrix(Y_test, Y_pred, labels=['0',
'1']))
else:
    cm += np.array(confusion_matrix(Y_test, Y_pred, labels=['0',
'1']))

result = pd.DataFrame(cm, index=['0', '1'], columns=['Predicted 0',
'Predicted 1'])
m.append(result)

cm = None
for train, test in skf.split(X, Y):
    X_train, X_test = X.iloc[train], X.iloc[test]
    Y_train, Y_test = Y.iloc[train], Y.iloc[test]
    predictorGnb.fit(X_train, Y_train)
    Y_pred = predictorGnb.predict(X_test)
    accuracyNB.append(accuracy_score(Y_test, Y_pred))
    if type(cm) == type(None):
        cm = np.array(confusion_matrix(Y_test, Y_pred, labels=['0',
'1']))
    else:
        cm += np.array(confusion_matrix(Y_test, Y_pred, labels=['0',
'1']))

result = pd.DataFrame(cm, index=['0', '1'], columns=['Predicted 0',
'Predicted 1'])
m.append(result)

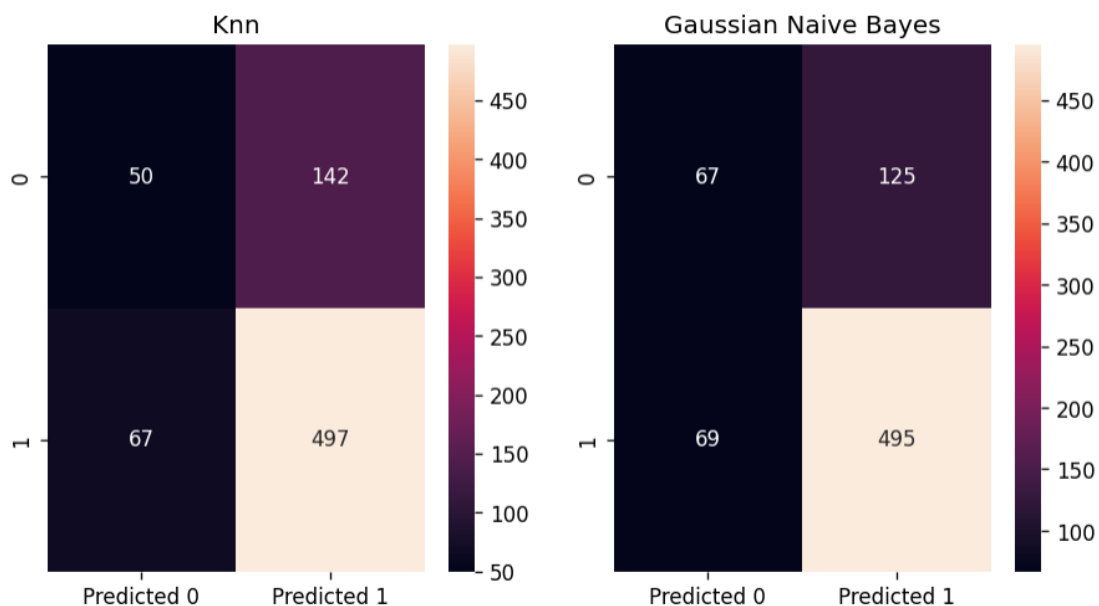
#plot confusion matrix
fig, ax = plt.subplots(1, 2, figsize=(10, 5))
sns.heatmap(m[0], annot=True, ax=ax[0], fmt='d')
sns.heatmap(m[1], annot=True, ax=ax[1], fmt='d')
ax[0].set_title('Knn')
ax[1].set_title('Gaussian Naive Bayes')
plt.show()

# Knn p value > NB p value
res = stats.ttest_rel(accuracyKnn, accuracyNB, alternative='greater')

```

```
print('Knn > Nb -> p-value: ', res.pvalue)
# Knn p value > NB p value
res = stats.ttest_rel(accuracyKnn, accuracyNB, alternative='less')
print('Knn < Nb -> p-value: ', res.pvalue)
```

5)



6)

Applying the T-Student formula for the hypothesis H0: “*k*NN is statistically superior to Naïve Bayes”, and H1 being the null hypothesis, we’ll get the following values:

```
Knn > Nb -> p-value: 0.9104476998751558
Knn < Nb -> p-value: 0.08955230012484414
```

In this way after seeing the values , since the p-value of the hypothesis H1 is lower than 0.10 this indicates strong evidence against the null hypothesis therefore validating the H0 statement.

7) Differences in performance between *k*NN and Naïve Bayes can be caused by:

- Naïve Bayes makes the assumption that all variables are independent, and there is no guarantee that this is the case in this context.
- *k*NN is a simple algorithm when compared to Naïve Bayes, but it is quite accurate when handling datasets of small samples, even more than Naïve Bayes, that requires a larger, more complex set of data to be more accurate.

- kNN has better performance when we normalize the data to the same scale that way the distance is more coherent with the classification.

END