

**I. Pen-and-paper**

1)a)

$$1) a) \quad P(y_6 = A) = \frac{3}{7}$$

$$P(y_6 = B) = \frac{4}{7}$$

$$P(y_3 = 0, y_4 = 0 | A) = 0$$

$$P(y_3 = 0, y_4 = 1 | A) = \frac{1}{3}$$

$$P(y_3 = 1, y_4 = 0 | A) = \frac{1}{3}$$

$$P(y_3 = 1, y_4 = 1 | A) = \frac{1}{3}$$

$$P(y_5 = 0 | A) = \frac{1}{3}$$

$$P(y_5 = 1 | A) = \frac{1}{3}$$

$$P(y_5 = 2 | A) = \frac{1}{3}$$

$$P(y_3 = 0, y_4 = 0 | B) = \frac{1}{2}$$

$$P(y_3 = 0, y_4 = 1 | B) = \frac{1}{4}$$

$$P(y_3 = 1, y_4 = 0 | B) = \frac{1}{4}$$

$$P(y_3 = 1, y_4 = 1 | B) = 0$$

$$P(y_5 = 0 | B) = \frac{1}{4}$$

$$P(y_5 = 1 | B) = \frac{1}{2}$$

$$P(y_5 = 2 | B) = \frac{1}{4}$$

$$\mu(y_1, y_2 | A) = \begin{pmatrix} 0.24 \\ 0.52 \end{pmatrix}$$

$$\Sigma(y_1, y_2 | A) = \begin{bmatrix} 0.0064 & 0.0096 \\ 0.0096 & 0.0336 \end{bmatrix}$$

$$\mu(y_1, y_2 | B) = \begin{pmatrix} 0.5925 \\ 0.3275 \end{pmatrix}$$

$$\Sigma(y_1, y_2 | B) = \begin{bmatrix} 0.0289 & -0.00976 \\ -0.00976 & 0.0315 \end{bmatrix}$$

1)b)

b) Com MAP

$$P(y_6 = A | x_8) = \frac{P(x_8 | A) P(A)}{P(x_8)} \stackrel{\text{MAP}}{=} P(x_8 | A) P(A) =$$

$$= P(y_1 = 0.38, y_2 = 0.52, y_3 = 0, y_4 = 1, y_5 = 0 | A) P(A) =$$

$$= P(y_1 = 0.38, y_2 = 0.52 | A) P(y_3 = 0, y_4 = 1 | A) P(y_5 = 0 | A) P(A) =$$

$$= 0.422 \times \frac{1}{3} \times \frac{1}{3} \times \frac{3}{7} \approx 0.0201$$

$$P(B | x_8) \stackrel{\text{MAP}}{=} \frac{P(x_8 | B) P(B)}{P(x_8)} = P(y_1 = 0.38, y_2 = 0.52 | B) P(y_3 = 0, y_4 = 1 | B) P(y_5 = 0 | B) P(B) =$$

$$= 1.121 \times \frac{1}{4} \times \frac{1}{4} \times \frac{4}{7} \approx 0.04$$

$$\text{Normalização: } P(A | x_8)_{\text{normalizado}} = \frac{P(A | x_8)}{P(A | x_8) + P(B | x_8)} = \frac{0.0201}{0.0201 + 0.04} \approx 0.3344$$

$$P(B | x_8)_{\text{normalizado}} = 1 - P(A | x_8)_{\text{norm.}} = 0.6656$$

Classificamos como **(B)**

$$P(A | x_9) \stackrel{\text{MAP}}{=} \frac{P(x_9 | A) P(A)}{P(x_9)} = P(y_1 = 0.42, y_2 = 0.59 | A) P(y_3 = 0, y_4 = 1 | A) P(y_5 = 1 | A) P(A) =$$

$$= 0.1727 \times \frac{1}{3} \times \frac{1}{3} \times \frac{3}{7} \approx 0.00822$$

$$P(B | x_9) \stackrel{\text{MAP}}{=} \frac{P(x_9 | B) P(B)}{P(x_9)} = P(y_1 = 0.42, y_2 = 0.59 | B) P(y_3 = 0, y_4 = 1 | B) P(y_5 = 1 | B) P(B) =$$

$$= 0.9878 \times \frac{1}{4} \times \frac{1}{2} \times \frac{4}{7} \approx 0.0706$$

$$\text{Normalização: } P(A | x_9)_{\text{norm.}} = \frac{P(A | x_9)}{P(A | x_9) + P(B | x_9)} = \frac{0.00822}{0.00822 + 0.0706} \approx 0.1043$$

$$P(B | x_9)_{\text{norm.}} = 1 - P(A | x_9)_{\text{norm.}} \approx 0.8957$$

Logo, também classificamos  $x_9$  como B



1)c)

$$1)c) \quad \begin{aligned} P(A|x_8) &= 0.3344 \\ P(A|x_9) &= 0.1043 \end{aligned}$$

$$f(x|\theta) = \begin{cases} A, & P(A|x) > \theta \\ B, & \text{caso contrário} \end{cases}$$

$$P(A|x_8)_{\text{norm.}} = \frac{0.3344}{0.3344 + 0.1043} = 0.762$$

$$P(A|x_9)_{\text{norm.}} = 1 - 0.762 = 0.238$$

Se o threshold  $\theta$  for menor que  $P(A|x_8)$  e  $P(A|x_9)$ :

$$\begin{aligned} x_8 &\rightarrow 0.762 > \theta \Rightarrow A \checkmark \\ x_9 &\rightarrow 0.238 > \theta \Rightarrow A \checkmark \end{aligned} \quad \left. \vphantom{\begin{aligned} x_8 &\rightarrow 0.762 > \theta \Rightarrow A \checkmark \\ x_9 &\rightarrow 0.238 > \theta \Rightarrow A \checkmark \end{aligned}} \right\} \text{accuracy} = \frac{1}{2} = 50\%$$

Se o threshold for um valor entre  $P(A|x_8)$  e  $P(A|x_9)$ ,  $0.238 < \theta < 0.762$ :

$$\begin{aligned} x_8 &\rightarrow 0.762 > \theta \Rightarrow A \checkmark \\ x_9 &\rightarrow 0.238 < \theta \Rightarrow B \checkmark \end{aligned} \quad \left. \vphantom{\begin{aligned} x_8 &\rightarrow 0.762 > \theta \Rightarrow A \checkmark \\ x_9 &\rightarrow 0.238 < \theta \Rightarrow B \checkmark \end{aligned}} \right\} \text{accuracy} = 100\%$$

Se o threshold for maior que  $P(A|x_8)$  e  $P(A|x_9)$ :

$$\begin{aligned} x_8 &\rightarrow 0.762 < \theta \Rightarrow B \times \\ x_9 &\rightarrow 0.238 < \theta \Rightarrow B \checkmark \end{aligned} \quad \left. \vphantom{\begin{aligned} x_8 &\rightarrow 0.762 < \theta \Rightarrow B \times \\ x_9 &\rightarrow 0.238 < \theta \Rightarrow B \checkmark \end{aligned}} \right\} \text{accuracy} = 50\%$$

Logo, qualquer threshold entre  $P(A|x_8)$  e  $P(A|x_9)$  atinge a testing accuracy

2)a)

2)a)

Binarização de  $y_2 \rightarrow$  domínio é  $[0;1]$ , logo:

$\text{lim}_0 = [0; 0.5[ \rightarrow y_2 = 0$

$\text{lim}_1 = [0.5; 1] \rightarrow y_2 = 1$

	$y_2$
$x_1$	0
$x_2$	0
$x_3$	1
$x_4$	0
$x_5$	0
$x_6$	0
$x_7$	1
$x_8$	1
$x_9$	1

Fold 1:

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_1$	0.24	0	1	1	0	A
$x_2$	0.16	0	1	0	1	A
$x_3$	0.32	1	0	1	2	A
$x_4$	0.54	0	0	0	1	B
$x_5$	0.66	0	0	0	0	B
$x_6$	0.76	0	1	0	2	B

Fold 2:

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_4$	0.54	0	0	0	1	B
$x_5$	0.66	0	0	0	0	B
$x_6$	0.76	0	1	0	2	B
$x_7$	0.41	1	0	1	1	B
$x_8$	0.38	1	0	1	0	A
$x_9$	0.42	1	0	1	1	B

Fold 3:

	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_1$	0.24	0	1	1	0	A
$x_2$	0.16	0	1	0	1	A
$x_3$	0.32	1	0	1	2	A
$x_7$	0.41	1	0	1	1	B
$x_8$	0.38	1	0	1	0	A
$x_9$	0.42	1	0	1	1	B



2)b)

2)b)  
Vamos Falt 1  $\rightarrow x_1$  a  $x_6$  training observations e  $x_7$  a  $x_9$  testing observations

Distância de Hamming de  $x_7$  a  $x_1$  até  $x_6$ :

$$H(x_7, x_1) = 1+1+0+1+1=4$$

$$H(x_7, x_4) = \boxed{2}$$

$$H(x_7, x_2) = 1+1+1+0+1=4$$

$$H(x_7, x_5) = \boxed{3}$$

$$H(x_7, x_3) = \boxed{2}$$

$$H(x_7, x_6) = 4$$

$$\text{Logo, } \hat{y}_1(x_7) = \frac{\frac{1}{2} \times 0.32 + \frac{1}{2} \times 0.54 + \frac{1}{3} \times 0.66}{\frac{1}{2} + \frac{1}{2} + \frac{1}{3}} = 0.4875$$

Distância de Hamming de  $x_8$  a  $x_1$  até  $x_6$ :

$$H(x_8, x_1) = \boxed{2}$$

$$H(x_8, x_3) = \boxed{1}$$

$$H(x_8, x_5) = \boxed{3}$$

$$H(x_8, x_2) = 4$$

$$H(x_8, x_4) = 4$$

$$H(x_8, x_6) = 5$$

$$\text{Logo, } \hat{y}_1(x_8) = \frac{\frac{1}{2} \times 0.24 + 1 \times 0.32 + \frac{1}{3} \times 0.66}{\frac{1}{2} + 1 + \frac{1}{3}} = 0.36$$

Distância de Hamming de  $x_9$  a  $x_1$  até  $x_6$ :

$$H(x_9, x_1) = 4$$

$$H(x_9, x_3) = \boxed{2}$$

$$H(x_9, x_5) = \boxed{3}$$

$$H(x_9, x_2) = 4$$

$$H(x_9, x_4) = \boxed{2}$$

$$H(x_9, x_6) = 4$$

$$\text{Logo, } \hat{y}_1(x_9) = \frac{\frac{1}{2} \times 0.32 + \frac{1}{2} \times 0.54 + \frac{1}{3} \times 0.66}{\frac{1}{2} + \frac{1}{2} + \frac{1}{3}} = 0.4875$$

$$\Rightarrow \text{MAE} = \frac{1}{n} \sum_i |res_i| = \frac{|0.41 - 0.4875| + |0.38 - 0.36| + |0.42 - 0.4875|}{3} = 0.055$$

## II. Programming and critical analysis

1)a)

```
import numpy as np
import matplotlib.pyplot as plt
from scipy.io import arff
from sklearn.model_selection import StratifiedKFold, cross_val_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from scipy.stats import ttest_rel
from sklearn.preprocessing import LabelEncoder

# Load the ARFF dataset
data, meta = arff.loadarff('column_diagnosis.arff')

# Convert the dataset to a Pandas DataFrame
import pandas as pd
df = pd.DataFrame(data)

# Extract features (X) and encode the target labels (y)
X = df.drop(columns=['class']).values
y = df['class'].str.decode('utf-8') # Decoding binary strings
label_encoder = LabelEncoder()
y = label_encoder.fit_transform(y)

# Set random seed for reproducibility
random_seed = 0

# Define classifiers
knn_classifier = KNeighborsClassifier(n_neighbors=5)
nb_classifier = GaussianNB()

# Define stratified 10-fold cross-validation
cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=random_seed)

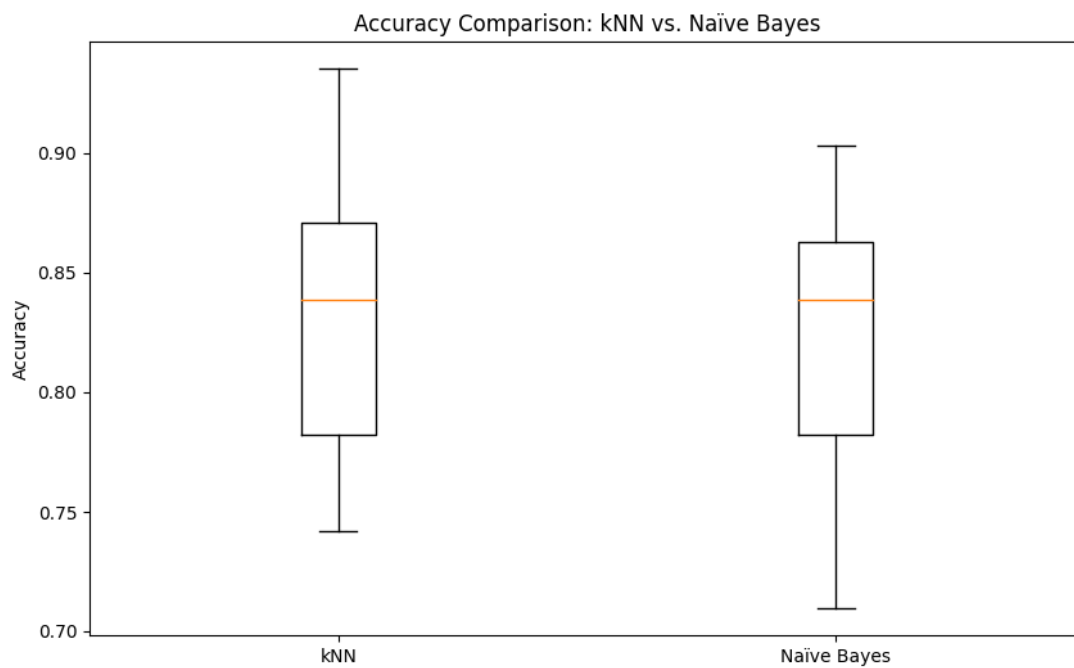
# Perform cross-validation and collect accuracy scores
knn_scores = cross_val_score(knn_classifier, X, y, cv=cv, scoring='accuracy')
nb_scores = cross_val_score(nb_classifier, X, y, cv=cv, scoring='accuracy')

# Plot boxplots of the fold accuracies
plt.figure(figsize=(10, 6))
plt.boxplot([knn_scores, nb_scores], labels=['kNN', 'Naïve Bayes'])
plt.title('Accuracy Comparison: kNN vs. Naïve Bayes')
plt.ylabel('Accuracy')
plt.savefig("Ex1ab")
```

```
# Perform a statistical test to compare the classifiers
t_statistic, p_value = ttest_rel(knn_scores, nb_scores)

# Set the significance level
alpha = 0.05

if p_value < alpha:
    print(f"kNN is statistically superior to Naïve Bayes (p-
value={p_value:.4f})")
else:
    print(f"There is no statistically significant difference between kNN and
Naïve Bayes (p-value={p_value:.4f})")
```



1)b)

There is no statistically significant difference between kNN and Naïve Bayes (p-value=0.3809)

2)

```
import pandas as pd
from scipy.io import arff
from sklearn.model_selection import cross_val_predict, StratifiedKFold
from sklearn.preprocessing import LabelEncoder
from sklearn.neighbors import KNeighborsClassifier
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Load the ARFF dataset
data = arff.loadarff('column_diagnosis.arff')
df = pd.DataFrame(data[0])

# Ensure the 'class' column contains only valid class labels
df['class'] = df['class'].str.decode('utf-8') # Convert bytes to strings

# Encode the 'class' column to numerical values
le = LabelEncoder()
df['class'] = le.fit_transform(df['class'])

# Split the dataset into features (X) and the target variable (y)
X = df.drop('class', axis=1) # Features
y = df['class'] # Target variable

# Create k-NN classifiers with k = 1 and k = 5
knn_1 = KNeighborsClassifier(n_neighbors=1)
knn_5 = KNeighborsClassifier(n_neighbors=5)

# Initialize StratifiedKFold with 10 folds and shuffling
cv = StratifiedKFold(n_splits=10, shuffle=True, random_state=0)

# Perform cross-validation and get predicted labels for each fold for k=1 and k=5
predicted_labels_1 = cross_val_predict(knn_1, X, y, cv=cv)
predicted_labels_5 = cross_val_predict(knn_5, X, y, cv=cv)

# Get the unique class names
class_names = le.classes_

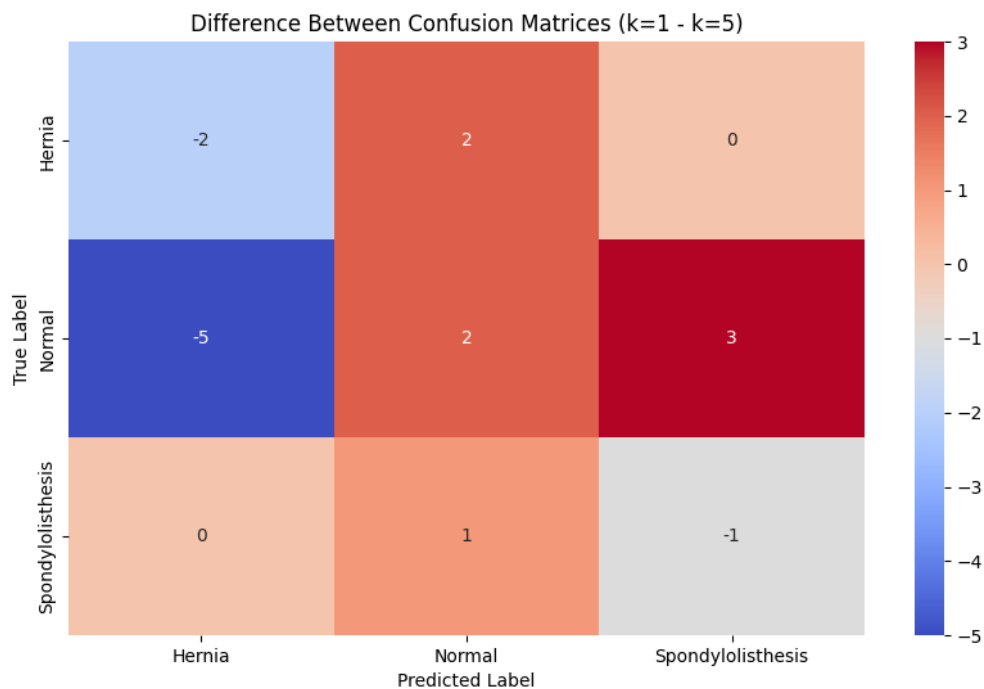
# Calculate confusion matrices for both k-NN classifiers
confusion_matrix_1 = pd.crosstab(index=y, columns=predicted_labels_1,
rownames=['True'], colnames=['Predicted'])
confusion_matrix_5 = pd.crosstab(index=y, columns=predicted_labels_5,
rownames=['True'], colnames=['Predicted'])
```



```
# Calculate the difference between the two confusion matrices
difference_matrix = confusion_matrix_1 - confusion_matrix_5

# Print the difference matrix
print("Matrix k=1:\n")
print(confusion_matrix_1)
print("Matrix k=5:\n")
print(confusion_matrix_5)
print("\nDifference Matrix (k=1 - k=5):")
print(difference_matrix)

# Create a heatmap to visualize the differences
plt.figure(figsize=(10, 6))
sns.heatmap(difference_matrix, annot=True, fmt="d", cmap="coolwarm", cbar=True)
plt.title('Difference Between Confusion Matrices (k=1 - k=5)')
plt.xlabel('Predicted Label')
plt.ylabel('True Label')
plt.xticks(np.arange(len(class_names)) + 0.5, class_names)
plt.yticks(np.arange(len(class_names)) + 0.5, class_names)
plt.savefig("Ex2")
```



Analisando os resultados vemos que quanto mais positivos forem os valores das células, concluímos que um  $k=1$  é melhor que um  $k=5$ , no que toca a estimar a class do atributo em questão. Por outro

lado, se o valor da célula for negativo tiramos uma conclusão semelhante, onde aqui o  $k$  preferível é o  $k=5$ . Nas células cinzentas, neutras sabemos que tanto o  $k=1$  e o  $k=5$  têm a mesma performance, sendo igual escolher um ou outro.

Tendo em conta esta observação concluímos que não existe um  $k$  que seja universalmente ótimo, sendo preciso ter sempre em consideração vários dados, como as características dos dados, a natureza das classes, etc.

3)

Tendo em conta `column_diagnosis`, podemos identificar 3 possíveis dificuldades:

1. **Suposição de Independência de Características:** O Naïve Bayes assume que as características são condicionalmente independentes dadas as classes. No entanto, no conjunto de dados "`column_diagnosis`", algumas características podem não ser inteiramente independentes. Por exemplo, medidas relacionadas à saúde da coluna, como '`pelvic_incidence`' e '`pelvic_tilt`', podem estar correlacionadas. Essa suposição de independência pode não ser verdadeira, o que pode afetar o desempenho do modelo.
2. **Características Contínuas:** O Naïve Bayes foi projetado para lidar com dados discretos e características categóricas. Se o conjunto de dados contém características contínuas, como é comum em conjuntos de dados biomédicos, pode ser necessário discretizá-las para encaixá-las no framework do Naïve Bayes. A escolha do método de discretização pode afetar o desempenho do modelo, e uma discretização inadequada pode resultar na perda de informações.
3. **Limitação na Expressividade do Modelo:** O Naïve Bayes é um classificador simples e linear. Ele pode não capturar relações complexas entre características nos dados. No conjunto de dados "`column_diagnosis`", pode haver padrões não lineares ou intrincados que o Naïve Bayes tem dificuldade em modelar eficazmente. Modelos mais avançados, como árvores de decisão, florestas aleatórias ou máquinas de vetores de suporte, podem ser mais adequados para capturar essas relações.

**END**