

I. Pen-and-paper

1)

Para 1

	$P(Y_1 C=0)$	$P(Y_1 C=1)$
μ	0,25	0,05
σ	0,238	0,288

μ_0 (média de y_1 's para $C=0$) = $\frac{0,6 + 0,1 + 0,2 + 0,1}{4} = 0,25$
 μ_1 (" " " " " $C=1$) = $\frac{0,3 + 0,1 + 0,3 + 0,2 + 0,4 + 0,2}{6} = 0,05$

σ_{y_0} (desvio padrão de y_1 's quando $C=0$) = $\sqrt{\frac{(0,6 - 0,25)^2 + (0,1 - 0,25)^2 + (0,2 - 0,25)^2 + (0,1 - 0,25)^2}{4 - 1}} = 0,238$
 σ_{y_1} (" " " " " $C=1$) = $\sqrt{\frac{(0,3 - 0,05)^2 + \dots}{6 - 1}} = 0,288$

Fica: $C=0: N(x_{muw} | 0,25; 0,238)$; $C=1: N(x_{muw} | 0,05; 0,288)$ para Y_1

$P(C=0) = 4/10$ $P(C=1) = 6/10$
 $P(Y_2=A | C=0) = 2/4$ $P(Y_2=A | C=1) = 1/6$
 $" = B | " = 1/4$ $" = B | " = 2/6$
 $" = C | " = 1/4$ $" = C | " = 3/6$

Para Y_2

	$P(Y_2, Y_4 C=0)$	$P(Y_2, Y_4 C=1)$
μ	$\begin{bmatrix} 0,2 \\ 0,25 \end{bmatrix}$	$\begin{bmatrix} 0,117 \\ 0,083 \end{bmatrix}$
Σ	$\begin{bmatrix} 0,18 & 0,18 \\ 0,18 & 0,25 \end{bmatrix}$	$\begin{bmatrix} 0,1097 & 0,1223 \\ 0,1223 & 0,2137 \end{bmatrix}$

$\mu_0 = \frac{1}{4} \left(\begin{bmatrix} 0,2 \\ 0,4 \end{bmatrix} + \begin{bmatrix} -0,1 \\ -0,4 \end{bmatrix} + \begin{bmatrix} -0,1 \\ 0,2 \end{bmatrix} + \begin{bmatrix} 0,0 \\ 0,8 \end{bmatrix} \right) = \begin{bmatrix} 0,2 \\ 0,25 \end{bmatrix}$
 $\Sigma_{y_2} = \frac{1}{3} \left(\begin{bmatrix} 0,4 - 0,25 \\ -0,4 - 0,25 \end{bmatrix}^2 + \begin{bmatrix} -0,2 - 0,25 \\ 0,2 - 0,25 \end{bmatrix}^2 + \begin{bmatrix} 0,8 - 0,25 \end{bmatrix}^2 \right) = \dots$

$\mu_1 = \frac{1}{6} \left(\begin{bmatrix} 0,1 \\ 0,3 \end{bmatrix} + \dots \right) = \begin{bmatrix} 0,117 \\ 0,083 \end{bmatrix}$
 $\Sigma_{y_4} = \frac{1}{5} \left((0,2 - 0,083)(0,1 - 0,117) + (-0,2 - 0,083)(0,2 - 0,117) + \dots \right) = 0,1223$

Para $C=1$

Fica $N(x_{muw} | \mu, \Sigma)$ para $X_3 \in Y_3$; conforme $C=0$ ou $C=1$, substitua μ e Σ da tabela

$P(C=0 | x_{muw}) = \frac{P(x_{muw} | C=0) P(C=0)}{P(x_{muw})} = \frac{P(C=0) P(Y_1 | C=0) P(Y_2 | C=0) P(Y_3, Y_4 | C=0)}{P(Y_1, Y_2, Y_3, Y_4)}$
 (igual a $C=1$)

Como $P(C=0 | x_{muw}) + P(C=1 | x_{muw}) = 1 \Leftrightarrow P(Y_1, Y_2, Y_3, Y_4) = P(C=0) P(Y_1 | C=0) \dots P(Y_3, Y_4 | C=0) + P(C=1) P(Y_1 | C=1) \dots$

$P(C=1 | x_1) = \dots = 0,081$
 Fazer igual para os restantes

2)

Para x_1 : $P(C=0 | x_1) = \frac{P(C=0) P(Y_1 | C=0) P(Y_2 | C=0) P(Y_3, Y_4 | C=0)}{P(Y_1, Y_2, Y_3, Y_4)} = \frac{4/10 \times N(0,6 | 0,25; 0,238) \times 2/4 \times N\left(\begin{bmatrix} 0,2 \\ 0,4 \end{bmatrix} \middle| \begin{bmatrix} 0,2 \\ 0,25 \end{bmatrix}; \begin{bmatrix} 0,18 & 0,18 \\ 0,18 & 0,25 \end{bmatrix}\right)}{P(Y_1, Y_2, Y_3, Y_4)}$

valores de x_1 \rightarrow não comparem manualmente; \rightarrow decompõem no índice na questão 4) = 0,137

$P(C=1 | x_1) = \dots = 0,081$

Fazer igual para os restantes

Aprendizagem 2021/22
Homework I – Group 96

2-

	$P(C=0 x_i)$	$P(C=1 x_i)$	Prod. Prod.	True
x_1	0,137	0,081	0	0
x_2	0,063	0,196	1	0
x_3	0,232	0,220	0	0
x_4	0,070	0,042	0	0
x_5	0,1252	0,229	0	1
x_6	0,038	0,244	1	1
x_7	0,016	0,120	1	1
x_8	0,237	0,203	0	1
x_9	0,020	0,026	1	1
x_{10}	0,061	0,321	1	1

$= P(C=0)P(x_1/C=0) \dots P(x_9/C=0) + P(C=1)P(x_2/C=1) \dots$

→ uma vez por mês normalizadas arredondadas as milésimas

TP: 4
FP: 1

TN: 3
FN: 2

	True		
	0	1	
Prod.	0 TN	2 TP	5
1 FP	1	4	5
	4	6	10

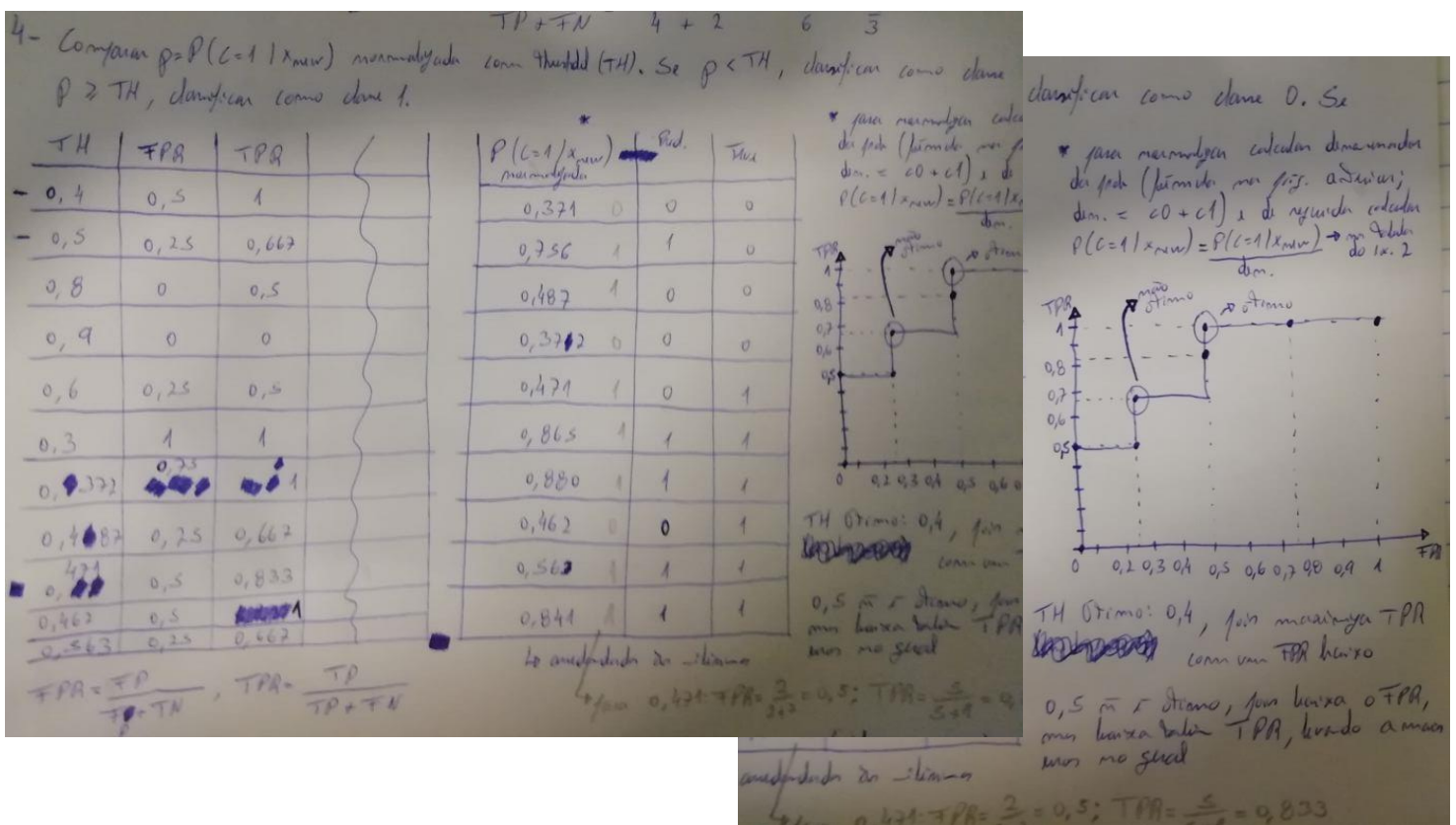
3)

3- $\frac{1}{F1} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right) = \frac{1}{2} \left(\frac{5}{4} + \frac{3}{2} \right) = \frac{1}{2} \left(\frac{11}{4} \right) = \frac{11}{8} \Rightarrow F1 = \frac{8}{11} = 0,727$

$P = \frac{TP}{TP+FP} = \frac{4}{4+1} = \frac{4}{5}$

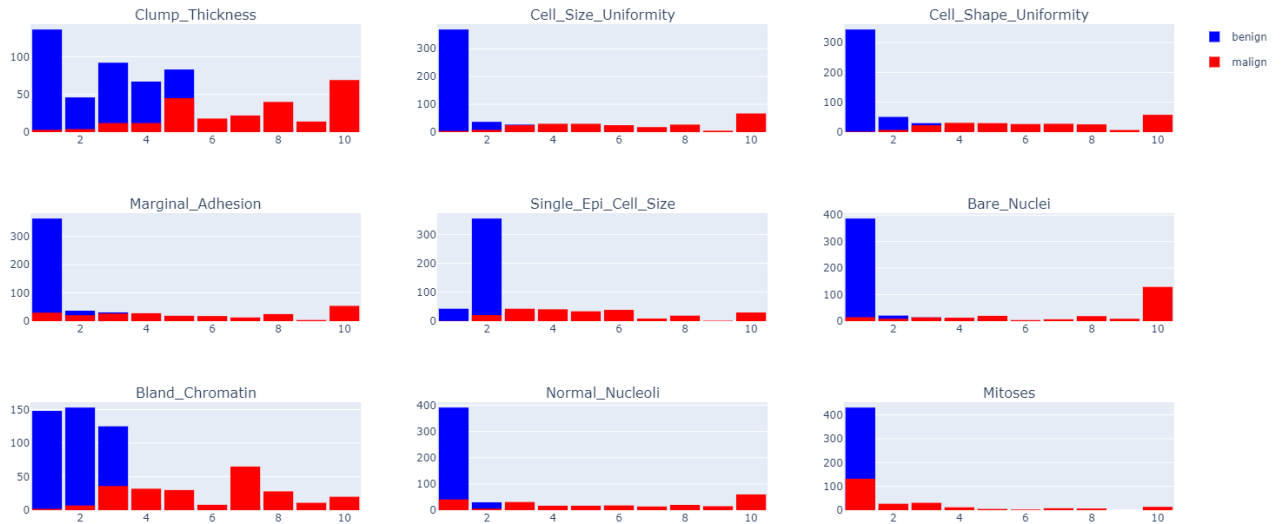
$R = \frac{TP}{TP+FN} = \frac{4}{4+2} = \frac{4}{6} = \frac{2}{3}$

4)



II. Programming and critical analysis

5)



6) 3NN Mean Test Accuracy = 96.6%, 5NN Mean Test Accuracy = 97.1%, 7NN Mean Test Accuracy = 96.9%

O 7NN é menos suscetível ao risco de overfitting, porque considera os 7 vizinhos mais próximos do novo elemento a classificar, levando a uma análise melhor, já que acaba por não ser tão influenciado por dados mal classificados junto do novo elemento. Quanto menor o nº de vizinhos, maior o grau de overfitting. Conforme se aumenta o nº de vizinhos, o overfitting diminui, mas a partir de certa altura começa a haver underfitting.

7) Ttest pvalue = 0.0004326 = 0.04326%. Como o pvalue é menor que o threshold de 1%, 5% e 10%, podemos rejeitar a hipótese nula de accuracies iguais entre o 3NN e o Naïve Bayes (multinomial assumption).

8) Com os dados obtidos anteriormente, concluímos que o KNN tem uma accuracy maior que o Naïve Bayes. Uma razão para isto acontecer é o facto do KNN lidar melhor com o overfitting do que o Naïve Bayes, levando a uma melhor classificação dos dados de teste. Outra razão para isto acontecer é devido ao Naïve Bayes assumir que as features são independentes quando na verdade isso pode não acontecer.

III. APPENDIX

[illegible]

Aprendizagem 2021/22

Homework I – Group 96

```

        "Bare_Nuclei": float,
        "Bland_Chromatin": int,
        "Normal_Nucleoli": int,
        "Mitoses": int,
        "Class": str}, na_values=["?"])

df = df.dropna()

columns = ["Clump_Thickness", "Cell_Size_Uniformity", "Cell_Shape_Uniformity",
"Marginal_Adhesion", "Single_Epi_Cell_Size", "Bare_Nuclei", "Bland_Chromatin",
"Normal_Nucleoli", "Mitoses", "Class"]
features = columns[:-1]

fig = make_subplots(rows=3, cols=3, subplot_titles=("Clump_Thickness",
"Cell_Size_Uniformity", "Cell_Shape_Uniformity", "Marginal_Adhesion",
"Single_Epi_Cell_Size", "Bare_Nuclei", "Bland_Chromatin", "Normal_Nucleoli",
"Mitoses"))

ben = df[df["Class"] == "benign"]
mal = df[df["Class"] == "malignant"]

fig.add_trace(go.Histogram(
    x=ben["Clump_Thickness"],
    bingroup=1, name="benign", legendgroup='group1', marker_color='#0000FF'),
row=1, col=1)

fig.add_trace(go.Histogram(
    x=mal["Clump_Thickness"],
    bingroup=1, name="malign", legendgroup='group2', marker_color="#FF0000"),
row=1, col=1)

fig.add_trace(go.Histogram(
    x=ben["Cell_Size_Uniformity"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
marker_color='#0000FF'), row=1, col=2)

fig.add_trace(go.Histogram(
    x=mal["Cell_Size_Uniformity"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
marker_color="#FF0000"), row=1, col=2)

fig.add_trace(go.Histogram(
    x=ben["Cell_Shape_Uniformity"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
marker_color='#0000FF'), row=1, col=3)

fig.add_trace(go.Histogram(
    x=mal["Cell_Shape_Uniformity"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
marker_color="#FF0000"), row=1, col=3)

fig.add_trace(go.Histogram(
    x=ben["Marginal_Adhesion"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
marker_color='#0000FF'), row=2, col=1)

fig.add_trace(go.Histogram(
    x=mal["Marginal_Adhesion"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
marker_color="#FF0000"), row=2, col=1)

fig.add_trace(go.Histogram(

```

Aprendizagem 2021/22
Homework I – Group 96

```
x=ben["Single_Epi_Cell_Size"],
bingroup=1, name="benign", legendgroup='group1', showlegend=False,
marker_color='#0000FF'), row=2, col=2)

fig.add_trace(go.Histogram(
    x=mal["Single_Epi_Cell_Size"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
    marker_color="#FF0000"), row=2, col=2)

fig.add_trace(go.Histogram(
    x=ben["Bare_Nuclei"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
    marker_color='#0000FF'), row=2, col=3)

fig.add_trace(go.Histogram(
    x=mal["Bare_Nuclei"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
    marker_color="#FF0000"), row=2, col=3)

fig.add_trace(go.Histogram(
    x=ben["Bland_Chromatin"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
    marker_color='#0000FF'), row=3, col=1)

fig.add_trace(go.Histogram(
    x=mal["Bland_Chromatin"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
    marker_color="#FF0000"), row=3, col=1)

fig.add_trace(go.Histogram(
    x=ben["Normal_Nucleoli"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
    marker_color='#0000FF'), row=3, col=2)

fig.add_trace(go.Histogram(
    x=mal["Normal_Nucleoli"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
    marker_color="#FF0000"), row=3, col=2)

fig.add_trace(go.Histogram(
    x=ben["Mitoses"],
    bingroup=1, name="benign", legendgroup='group1', showlegend=False,
    marker_color='#0000FF'), row=3, col=3)

fig.add_trace(go.Histogram(
    x=mal["Mitoses"],
    bingroup=1, name="malign", legendgroup='group2', showlegend=False,
    marker_color="#FF0000"), row=3, col=3)

fig.update_layout(
    barmode="overlay",
    bargap=0.1)

fig.show() # Pergunta 5

Y = df["Class"]
X = df.drop(columns=["Class"])

kf10 = KFold(n_splits=10)
model = neighbors.KNeighborsClassifier(n_neighbors=3, weights='uniform', p=2)
gnb = MultinomialNB()
```

Aprendizagem 2021/22
Homework I – Group 96

```
Y = [0 if x == "benign" else 1 for x in Y]
Y = np.ravel(pd.DataFrame(Y))

acc_knn = []
acc_gnb = []
i = 1
for train_index, test_index in kf10.split(df):
    X_train = X.iloc[train_index].loc[:, features].values
    X_test = X.iloc[test_index][features].values
    y_train = Y[train_index]
    y_test = Y[test_index]

    model.fit(X_train, y_train)
    gnb.fit(X_train, y_train)
    acc_knn += [accuracy_score(y_test, model.predict(X_test))]
    acc_gnb += [accuracy_score(y_test, gnb.predict(X_test))]
    i += 1

print("KNN: ", mean(acc_knn), " GNB: ", mean(acc_gnb)) # Pergunta 6

print(stats.ttest_rel(acc_knn, acc_gnb)) # Pergunta 7
```

END

Nota: Não houve paridade de esforço entre os elementos do grupo. Eu, Afonso Ferreira - 96832, fiz este homework sozinho.