

I. Pen-and-paper

1)

1) Assumindo as variáveis/conjuntos de variáveis independentes, podemos utilizar um classificador "Naïve Bayes" para responder ao problema.

$$\text{class} = \underset{k \in \{0,1\}}{\text{argmax}} p(\text{class}_k) \prod_{i=1}^n p(x_i | \text{class}_k)$$

• $P(Y_1 = 0 | \text{class} = 0)$: Y_1 é distribuída normalmente ($N(\mu, \sigma^2)$)

$$\bar{x} = \mu = \frac{1}{N} \sum_{i=1}^n x_i = \frac{1}{4} \cdot (0,6 + 0,1 + 0,2 + 0,1) = 0,25$$

$$s_{\text{sample}} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \approx 0,238$$

• $P(Y_1 = 0 | \text{class} = 1)$: ($N(\mu, \sigma^2)$)
 $\mu = 0,05$ $s_{\text{sample}} = 0,288$

$$P(Y_2 = A | \text{class} = 0) = \frac{1}{2} ; P(Y_2 = A | \text{class} = 1) = \frac{1}{6}$$

$$P(Y_2 = B | \text{class} = 0) = \frac{1}{4} ; P(Y_2 = B | \text{class} = 1) = \frac{2}{6} = \frac{1}{3}$$

$$P(Y_2 = C | \text{class} = 0) = \frac{1}{4} ; P(Y_2 = C | \text{class} = 1) = \frac{3}{6} = \frac{1}{2}$$

• $P(Y_3 = 0_3 ; Y_4 = 0_4 | \text{class} = 0)$: $\{Y_3, Y_4\}$ tem distribuição normal multivariada (2D) $N(\underline{\mu}, \underline{\Sigma})$

$$\underline{\mu} = [0,2 ; 0,28]^T = [\mu_1 ; \mu_2]^T = [\bar{x}_1 ; \bar{x}_2]$$

$$\underline{\Sigma}_{i,j} = c_{i,j} = \frac{\sum_{k=1}^n (x_{k,i} - \bar{x}_i)(x_{k,j} - \bar{x}_j)}{N-1} = \begin{bmatrix} 0,18 & 0,18 \\ 0,18 & 0,25 \end{bmatrix}$$

• $P(Y_3 = 0_3 ; Y_4 = 0_4 | \text{class} = 1)$: $N(\underline{\mu}, \underline{\Sigma})$

$$\underline{\mu} = \left[\frac{7}{60} ; \frac{1}{12} \right]^T = [\mu_1 ; \mu_2]^T = [\bar{x}_1 ; \bar{x}_2]$$

$$\underline{\Sigma} = \begin{bmatrix} 0,1097 & 0,1223 \\ 0,1223 & 0,2137 \end{bmatrix}$$

Formulas usadas nas distribuições normais:

$$\text{Univariável: } \frac{1}{\sqrt{2\pi} \cdot \sigma} \cdot \exp\left(-\frac{1}{2 \cdot \sigma^2} \cdot (x - \mu)^2\right)$$

$$2 \text{ variáveis: } \frac{1}{2 \cdot \pi} \cdot \frac{1}{|\underline{\Sigma}|^{1/2}} \cdot \exp\left(-\frac{1}{2} \cdot (x - \underline{\mu})^T \underline{\Sigma}^{-1} (x - \underline{\mu})\right)$$

Quando queremos classificar um novo evento, por exemplo, x_{new} , calculamos a sua probabilidade condicionada à classificação = 0 e à classificação = 1. A probabilidade superior dará a sua respetiva classificação respetivamente, 0 ou 1.

2) Answer 2

2) Usando ferramentas computacionais obtivemos os seguintes valores para os eventos x_1, \dots, x_{10}

~~$P(Y_1 = 0 | \text{class} = 1)$~~

$$P(x_1 | \text{class} = 0) = 0,067 \dots$$

$$P(x_1 | \text{class} = 1) = 0,017 \dots$$

$$P(x_2 | \text{class} = 0) = 0,028 \dots$$

$$P(x_2 | \text{class} = 1) = 0,170 \dots$$

$$P(x_3 | \text{class} = 0) = 0,103 \dots$$

$$P(x_3 | \text{class} = 1) = 0,048 \dots$$

$$P(x_4 | \text{class} = 0) = 0,037 \dots$$

$$P(x_4 | \text{class} = 1) = 0,0542 \dots$$

$$P(x_5 | \text{class} = 0) = 0,095 \dots$$

$$P(x_5 | \text{class} = 1) = 0,149 \dots$$

$$P(x_6 | \text{class} = 0) = 0,008 \dots$$

$$P(x_6 | \text{class} = 1) = 0,158 \dots$$

$$P(x_7 | \text{class} = 0) = 0,003 \dots$$

$$P(x_7 | \text{class} = 1) = 0,078 \dots$$

$$P(x_8 | \text{class} = 0) = 0,079 \dots$$

$$P(x_8 | \text{class} = 1) = 0,132 \dots$$

$$P(x_9 | \text{class} = 0) = 0,026 \dots$$

$$P(x_9 | \text{class} = 1) = 0,016 \dots$$

$$P(x_{10} | \text{class} = 0) = 0,0735 \dots$$

$$P(x_{10} | \text{class} = 1) = 0,209 \dots$$

$$\hat{y} = [0, 1, 0, 1, 1, 1, 1, 1, 0, 1] = \text{predicted}$$

$$y = \text{observed} = [0, 0, 0, 0, 1, 1, 1, 1, 1, 1]$$

$\hat{y} \backslash y$	1	0
1	5	1
0	2	2

3) Answer 3

4) Answer 4

3) F1:

$$\frac{1}{F} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)$$

$$P = \text{Precision} = \frac{TP}{TP+FP} = \frac{5}{5+2} = \frac{5}{7}$$

$$R = \text{Recall} = \frac{TP}{P} = \frac{TP}{TP+FN} = \frac{5}{5+1} = \frac{5}{6}$$

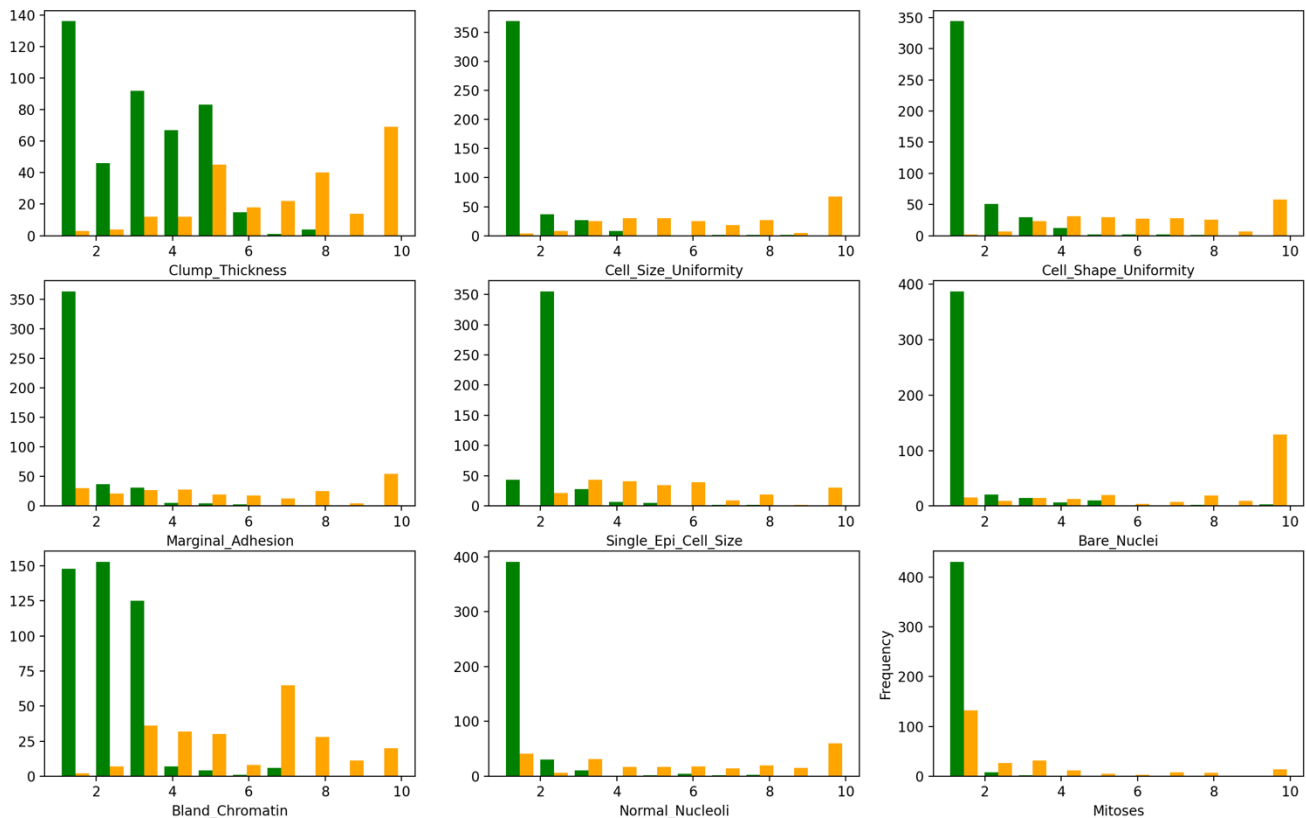
$$\frac{1}{F} = \frac{1}{2} \left(\frac{1}{\frac{5}{7}} + \frac{1}{\frac{5}{6}} \right) \Leftrightarrow \frac{1}{F} = \frac{1}{2} \left(\frac{7}{5} + \frac{6}{5} \right) \Leftrightarrow$$

$$\Leftrightarrow \frac{1}{F} = \frac{1}{2} \cdot \frac{13}{5} \Leftrightarrow \frac{1}{F} = \frac{13}{10} \Leftrightarrow F = \frac{10}{13} \approx 77\%$$

4) Usamos o AOC para analisar o nosso "Bayesian Classifier" e obtemos o seguinte gráfico onde podemos observar que o "decision threshold" que otimiza o treino é no ponto (0.5; 0.83) pois assim temos um "True positive rate" elevado face ao "False Positive rate". A partir desse ponto a "False Positive rate" aumenta bastante face ao "True Positive Rate". O que vai piorar os resultados, pois iremos ter mais falsos positivos.

II. Programming and critical analysis

5) Answer 5



6) Answer 6

When $k = 3$ Distance Accuracy train : 0.9746 test : 0.9883

Uniform Accuracy train : 0.9785 test : 0.9825

When $k = 5$ Distance Accuracy train : 0.9746 test : 0.9649

Uniform Accuracy train : 0.9844 test : 0.9649

When $k = 7$ Distance Accuracy train : 0.9785 test : 0.9649

Uniform Accuracy train : 0.9805 test : 0.9825

A melhor tentativa foi usando $k = 3$ através da distância Euclidiana. O overfitting é mais provável acontecer quando o conjunto de treino tem melhor rating que o conjunto de teste. Para os $K = 3$ é possível verificar que os score do teste > score de treino, sendo benéfico para a amostra, reduzindo o risco de Overfitting. Por outro lado quanto mais Accuracy, mais detalhada vai ser a regressão, o que não é ótimo para a amostra de dados (eventual perda de amostras significativa). Logo quanto menor for o Accuracy melhor será a regressão.

7) Answer 7

Multinomial NB score = 0.9048

$K = 3 = 0.9746$

A hipótese é verificável para esta experiência.

8) Answer 8

O classificador kNN em relação ao classificador de Naives Bayes é o mais fiável para a experiência. Quanto mais a taxa de acerto se aproxima de 1, mais detalhada vai ser a curva de regressão, o que não é fiável para os resultados, levando para situações de overfitting. Ainda assim o classificador kNN é o mais indicado para amostras com muitos atributos.

III. APPENDIX

```

import numpy as np
import pandas as pd
from scipy.io.arff import loadarff
import codecs
import matplotlib.pyplot as plt
from scipy.io.arff.arffread import print_attribute
from sklearn.neighbors import DistanceMetric
from numpy.random.mtrand import RandomState, random_sample
from sklearn.model_selection import KFold
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.model_selection import train_test_split
from sklearn.neighbors import KNeighborsClassifier
from sklearn import neighbors, datasets, preprocessing
from sklearn.metrics import confusion_matrix
from sklearn.naive_bayes import MultinomialNB

if __name__ == '__main__':
    raw_data = loadarff('breast.w.arff')
    df_data = pd.DataFrame(raw_data[0])
    classe = df_data.pop('Class')
    df_data = df_data.astype(int)
    Y = classe.str.decode('utf-8')
def freq(n): #getting frequency by symbol for benign and malignant
    a = []
    b = []
    for val in range(0,683):#iterating columns
        c = df_data.iloc[val,n]
        if Y[val] == 'benign':
            a.append(c)
        else:
            b.append(c)
    return [a,b]
fig, axs = plt.subplots(3,3)
names = ['Clump_Thickness',
'Cell_Size_Uniformity','Cell_Shape_Uniformity','Marginal_Adhesion','Single_Epi_Cell_Size','Bare_Nuc
lei','Bland_Chromatin','Normal_Nucleoli','Mitoses'] #Column names (Attributes)
for i in range(1,10):# plotting by column (attributes)
    plt.subplot(3,3,i)
    plt.hist(freq(i-1),10,range=(1,10), align=('mid'), color=['green', 'orange'], label=['benign',
'malignant'])
    plt.xlabel(names[i-1])
plt.ylabel("Frequency")
plt.show()
#Exercise 6
X = df_data
kf = KFold(10,shuffle=True,random_state=114)
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.25)
knn = KNeighborsClassifier(n_neighbors = 3 ,weights='uniform')#uniform weights
knn = KNeighborsClassifier(n_neighbors = 3 ,metric='wminkowski', p=2,
metric_params={'w': np.random.random(X_train.shape[1])})#euclidean
knn.fit(X,Y)
a = knn.score(X_train, Y_train)
b = knn.score(X_test,Y_test)
print("When k = 3 uniform Accuracy train : %.4f test : %.4f" % (a,b))
#Exercise 7
clf = MultinomialNB()
clf.fit(X, Y)
print("Multinomial NB score = %.4f " %clf.score(X,Y))

```

END