

### Part I: Pen and paper

1. Consider  $x_1$ – $x_7$  to be training observations,  $x_8$ – $x_9$  to be testing observations,  $y_1$  –  $y_5$  to be input variables and  $y_6$  to be the target variable.
- a. Learn a Bayesian classifier assuming: i)  $\{y_1, y_2\}$ ,  $\{y_3, y_4\}$  and  $\{y_5\}$  sets of independent variables (e.g.,  $y_1 \perp y_3$  yet  $y_1 \not\perp y_2$ ), and ii)  $y_1 \times y_2 \in \mathbb{R}^2$  is normally distributed. Show all parameters (distributions and priors for subsequent testing).

priors:  $P(A) = \frac{3}{7}$        $P(B) = \frac{4}{7}$

distributions:

$\{y_1, y_2\}$        $P(y_1, y_2 | A) \sim N(\mu_A, \Sigma_A)$

$$\mu_A = \begin{bmatrix} \frac{0,24 + 0,16 + 0,32}{3} \\ \frac{0,36 + 0,48 + 0,72}{3} \end{bmatrix} = \begin{bmatrix} 0,24 \\ 0,52 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \text{cov}(y_1, y_1) & \text{cov}(y_1, y_2) \\ \text{cov}(y_1, y_2) & \text{cov}(y_2, y_2) \end{bmatrix}$$

$$\text{cov}(y_1, y_1) = \frac{\sum_{i=1}^n (x_{i1} - \bar{y}_1)^2}{n-1} = \frac{(0,24-0,24)^2 + (0,16-0,24)^2 + (0,32-0,24)^2}{2} =$$

porque é uma amostra  $\rightarrow$   $= 6,4 \times 10^{-3} = 0,0064$

$$\text{cov}(y_2, y_2) = \frac{\sum_{i=1}^n (x_{i2} - \bar{y}_2)^2}{n-1} = \frac{(0,36-0,52)^2 + (0,48-0,52)^2 + (0,72-0,52)^2}{2} =$$

$$= 0,0336$$

$$\text{cov}(y_1, y_2) = \frac{\sum_{i=1}^n (x_{i1} - \bar{y}_1)(x_{i2} - \bar{y}_2)}{n-1} =$$

$$= \frac{(0,24-0,24)(0,36-0,52) + (0,16-0,24)(0,48-0,52) + (0,32-0,24)(0,72-0,52)}{2} =$$

$$= 9,6 \times 10^{-3} = 0,0096$$

$$\Sigma_A = \begin{bmatrix} 0,0064 & 0,0096 \\ 0,0096 & 0,0336 \end{bmatrix}$$

$$|\Sigma_A| = 0,0064 \times 0,0336 - 0,0096^2 = 1,2288 \times 10^{-2} = 0,00012288$$

$$\Sigma_A^{-1} = \begin{bmatrix} 273,4375 & -78,125 \\ -78,125 & 52,0833 \end{bmatrix}$$

$$p(x|\mu, \Sigma) = N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{n/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

$$p(x|\mu_A, \Sigma_A) = \frac{1}{2\pi \times 0,011085} \exp\left(-\frac{1}{2}\left(x - \begin{bmatrix} 0,24 \\ 0,52 \end{bmatrix}\right)^T \begin{bmatrix} 273,4375 & -78,125 \\ -78,125 & 52,0833 \end{bmatrix} \left(x - \begin{bmatrix} 0,24 \\ 0,52 \end{bmatrix}\right)\right)$$

$$p(y_1, y_2 | B) \sim N(\mu_B, \Sigma_B)$$

$$\mu_B = \begin{bmatrix} 0,5925 \\ 0,3275 \end{bmatrix} \quad \Sigma_B = \begin{bmatrix} 0,02289 & -0,009758 \\ -0,009758 & 0,03149 \end{bmatrix} \quad |\Sigma| = 0,00062567$$

$$\Sigma^{-1} = \begin{bmatrix} 50,3325 & 15,59657 \\ 15,59657 & 36,5873 \end{bmatrix}$$

$$p(x|\mu_B, \Sigma_B) = \frac{1}{2\pi \times 0,0250134} \exp\left(-\frac{1}{2}\left(x - \begin{bmatrix} 0,5925 \\ 0,3275 \end{bmatrix}\right)^T \begin{bmatrix} 50,3325 & 15,5966 \\ 15,5966 & 36,5873 \end{bmatrix} \left(x - \begin{bmatrix} 0,5925 \\ 0,3275 \end{bmatrix}\right)\right)$$

$\{y_3, y_4\}$

$$P(y_3=1, y_4=1 | A) = \frac{2}{3}$$

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

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

$$P(y_3=0, y_4=0 | B) = \frac{2}{4} = \frac{1}{2}$$

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

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

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

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

$\{y_5\}$

$$P(y_5=0 | A) = \frac{2}{3}$$

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

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

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

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

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

b. Under a MAP assumption, classify each testing observation showing all your calculus.

MAP assumption:  $h_{\text{MAP}} = \underset{h}{\operatorname{argmax}} P(D|h) P(h)$   
 $\hookrightarrow P(D|h) = \prod P(y_i|h)$  pois as conjugates & b independent

$$P(x_8|A) = p(y_3=0, y_4=1|A) p(y_5=0|A) p(y_1=0,38, y_2=0,52|A) =$$

$$= \frac{1}{3} \times \frac{1}{3} \times \frac{1}{2\pi \times 0,011085} \exp\left(-\frac{1}{2} \begin{bmatrix} 0,14 & 0 \end{bmatrix} \begin{bmatrix} 273,4375 & -78,125 \\ -78,125 & 52,0833 \end{bmatrix} \begin{bmatrix} 0,14 \\ 0 \end{bmatrix}\right) =$$

$$= \frac{1}{9} \times \frac{1}{2\pi \times 0,011085} e^{-2,6796895} \approx 0,10941$$

$$P(x_8|A) p(A) = 0,10941 \times \frac{3}{7} = 0,04689$$

$$P(x_8|B) = p(y_3=0, y_4=1|B) p(y_5=0|B) p(y_1=0,38, y_2=0,52|A) =$$

$$= \frac{1}{4} \times \frac{1}{4} \times 1,96237 \approx 0,12265$$

$$P(x_8|B) p(B) = 0,12265 \times \frac{4}{7} \approx 0,07008$$

$$p(A|x_8) = \frac{p(x_8|A) p(A)}{p(x_8|A) p(A) + p(x_8|B) p(B)} = \frac{0,04689}{0,04689 + 0,07008} \approx 0,40087$$

$$P(B|x_8) = 1 - p(A|x_8) = 0,59913$$

$$h_{\text{MAP}} = B$$

$$p(x_9|A) = p(y_3=0, y_4=1|A) p(y_5=1|A) p(y_1=0,42, y_2=0,59|A) =$$

$$= \frac{1}{3} \times \frac{1}{3} \times 0,40307 \approx 0,04478$$

$$P(x_9|B) = p(y_3=0, y_4=1|B) p(y_5=1|B) p(y_1=0,42, y_2=0,59|B) =$$

$$= \frac{1}{4} \times \frac{1}{2} \times 1,72857 = 0,21607$$

$$P(X_9|A)P(A) = 0,04478 \times \frac{3}{7} \approx 0,01919$$

$$P(X_9|B)P(B) = 0,21607 \times \frac{4}{7} \approx 0,12347$$

$$P(A|X_9) = \frac{P(X_9|A)P(A)}{P(X_9|A)P(A) + P(X_9|B)P(B)} = \frac{0,01919}{0,01919 + 0,12347} \approx 0,13451$$

$$P(B|X_9) = 1 - 0,13451 = 0,86549 \quad h_{MAP} = B$$

c. Consider that the default decision threshold of  $\theta = 0.5$  can be adjusted according to

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

Under a maximum likelihood assumption, what thresholds optimize testing accuracy?

ML assumption:  $h_{ML} = \underset{h}{\operatorname{argmax}} P(D|h) \Rightarrow$  assume que  $P(A) = P(B)$   
logo  $P(D|h) \propto P(h|D)$

$$P(X_8|A) = 0,10941 \quad P(X_9|A) = 0,04478$$

$$P(X_8|B) = 0,12265 \quad P(X_9|B) = 0,21607$$

$$P(A|X_8) = \frac{P(X_8|A)}{P(X_8|A) + P(X_8|B)} = \frac{0,10941}{0,10941 + 0,12265} \approx 0,47147$$

$$P(B|X_8) = 1 - P(A|X_8) = 1 - 0,47147 = 0,52853 \quad h_{ML} = B$$

$$P(A|X_9) = \frac{0,04478}{0,04478 + 0,21607} \approx 0,17167$$

$$h_{ML} = B$$

$$P(B|X_9) = 1 - 0,17167 = 0,82833$$

$$f(X_8|\theta=0,5) = B \quad f(X_9|\theta=0,5) = B \quad P(\text{certo}) = 50\%$$

$$f(X_8|\theta \in [0,17167, 0,47147]) = A$$

$$P(\text{certo}) = 100\%$$

$$f(X_9|\theta \in [0,17167, 0,47147]) = B$$

Logo  $\theta$  no intervalo  $0,17167 \leq \theta < 0,47147$  otimiza a testing accuracy para as observações teste.

2. Let  $y_1$  be the target numeric variable,  $y_2$ - $y_6$  be the input variables where  $y_2$  is binarized under an equal-width (equal-range) discretization. For the evaluation of regressors, consider a 3-fold cross-validation over the full dataset ( $x_1$ - $x_9$ ) without shuffling the observations.

a. Identify the observations and features per data fold after the binarization procedure.

a. Binarização

$D_{y_2} = [0, 1] \quad \frac{1-0}{2} = 0,5$

$\text{bino} = \{0,11; 0,28; 0,36; 0,39; 0,48\}$   
 $[0,0,5]$

$\text{bino}_2 = \{0,52; 0,53; 0,59; 0,72\}$   
 $[0,5; 0,5]$

new  $y_2 = [0, 0, 1, 0, 0, 0, 1, 1, 1]$

3-fold cross-validation

Iteração 1

fold 1:  $x_1 - x_3$   
 fold 2:  $x_4 - x_6$   
 fold 3:  $x_7 - x_9$

D	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_1$	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
$\bar{x}_4$	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
$\bar{x}_7$	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

Iteração 2

D	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_1$	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
$\bar{x}_4$	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
$\bar{x}_7$	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

Iteração 3

D	$y_1$	$y_2$	$y_3$	$y_4$	$y_5$	$y_6$
$x_1$	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
$\bar{x}_4$	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
$\bar{x}_7$	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

- b. Consider a distance-weighted  $kNN$  with  $k = 3$ , Hamming distance ( $d$ ), and  $1/d$  weighting. Compute the MAE of this  $kNN$  regressor for the 1<sup>st</sup> iteration of the cross-validation (i.e. train observations have the lower indices).

b.  $x_1 - x_6$  - treino  
 $x_7 - x_9$  - teste

Hamming distance = número de atributos diferentes entre dois valores

$$\begin{aligned} d(x_7, x_1) &= 4 \\ d(x_7, x_2) &= 4 \\ \rightarrow d(x_7, x_3) &= 2 \\ \rightarrow d(x_7, x_4) &= 2 \\ \rightarrow d(x_7, x_5) &= 3 \\ d(x_7, x_6) &= 4 \end{aligned}$$

3 vizinhos mais próximos de  $x_7$ :

$$x_3, x_4, x_5$$

$$\begin{aligned} \rightarrow d(x_8, x_1) &= 2 \\ d(x_8, x_2) &= 4 \\ \rightarrow d(x_8, x_3) &= 4 \\ d(x_8, x_4) &= 4 \\ \rightarrow d(x_8, x_5) &= 3 \\ d(x_8, x_6) &= 5 \end{aligned}$$

3 vizinhos mais próximos de  $x_8$ :

$$x_1, x_3, x_5$$

$$\begin{aligned} d(x_9, x_1) &= 4 \\ d(x_9, x_2) &= 4 \\ \rightarrow d(x_9, x_3) &= 2 \\ \rightarrow d(x_9, x_4) &= 2 \\ \rightarrow d(x_9, x_5) &= 3 \\ d(x_9, x_6) &= 4 \end{aligned}$$

3 vizinhos mais próximos de  $x_9$ :

$$x_3, x_4, x_6$$

Distâncias ponderadas:

$$w_i = \frac{1}{d(x, x_i)}$$

$$\begin{aligned} x_7: \quad w_3 &= \frac{1}{2} = 0,5 \\ w_4 &= \frac{1}{2} = 0,5 \\ w_5 &= \frac{1}{3} \end{aligned}$$

$$x_8: \quad w_1 = \frac{1}{2} = 0,5$$

$$w_3 = \frac{1}{4} = 0,25$$

$$w_5 = \frac{1}{3}$$

$$x_9: \quad w_3 = \frac{1}{2} = 0,5$$

$$w_4 = \frac{1}{2} = 0,5$$

$$w_6 = \frac{1}{3}$$

$$\hat{z}_i = \frac{\sum w_i \times \text{output}(x_i)}{\sum w_i}$$

$$\hat{z}_7 = \frac{0,5 \times 0,32 + 0,5 \times 0,54 + \frac{1}{3} \times 0,66}{0,5 + 0,5 + \frac{1}{3}} = 0,4875$$

$$\hat{z}_8 = \frac{0,5 \times 0,24 + 1 \times 0,32 + \frac{1}{3} \times 0,66}{0,5 + 1 + \frac{1}{3}} = 0,36$$

$$\hat{z}_9 = \frac{0,5 \times 0,32 + 0,5 \times 0,54 + \frac{1}{3} \times 0,66}{0,5 + 0,5 + \frac{1}{3}} = 0,4875$$

$$MAE = \frac{\sum |z_i - \hat{z}_i|}{n}$$

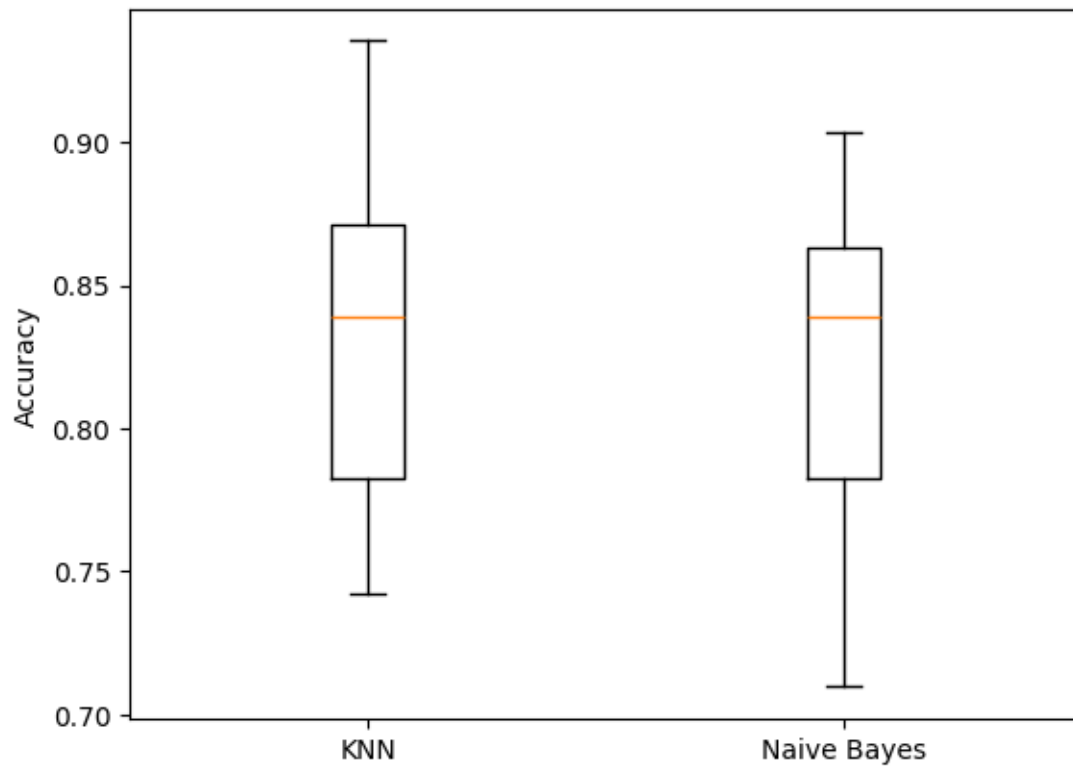
$$MAE = \frac{|0,41 - 0,4875| + |0,38 - 0,36| + |0,42 - 0,4875|}{3} = 0,055$$

## Part II: Programming

Considering the `column_diagnosis.arff` dataset available at the course webpage's homework tab. Using `sklearn`, apply a 10-fold stratified cross-validation with shuffling (`random_state = 0`) for the assessment of predictive models along this section.

1. Compare the performance of  $kNN$  with  $k = 5$  and naïve Bayes with Gaussian assumption (consider all remaining parameters for each classifier as `sklearn`'s default):
  - a. Plot two boxplots with the fold accuracies for each classifier.

```
1 from sklearn.model_selection import StratifiedKFold
2 from sklearn.neighbors import KNeighborsClassifier
3 from sklearn.naive_bayes import GaussianNB
4 import pandas as pd
5 from scipy.io.arff import loadarff
6 import matplotlib.pyplot as plt
7 import numpy as np
8
9 data = loadarff('column_diagnosis.arff')
10 df = pd.DataFrame(data[0])
11 y = df['class'].astype(str)
12 x = df.drop('class', axis=1).astype(np.float64)
13
14 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=0)
15
16 knn = KNeighborsClassifier(n_neighbors=5)
17 naive_bayes = GaussianNB()
18
19 knn_scores = []
20 naive_bayes_scores = []
21 for train_index, test_index in skf.split(x, y):
22     x_train, x_test = df.iloc[train_index, :-1], df.iloc[test_index, :-1]
23     y_train, y_test = df.iloc[train_index, -1], df.iloc[test_index, -1]
24     x_train = x_train.astype(np.float64)
25     x_test = x_test.astype(np.float64)
26     y_train = y_train.astype(str)
27     y_test = y_test.astype(str)
28     knn.fit(x_train, y_train)
29     knn_scores.append(knn.score(x_test, y_test))
30     naive_bayes.fit(x_train, y_train)
31     naive_bayes_scores.append(naive_bayes.score(x_test, y_test))
32
33 plt.boxplot([knn_scores, naive_bayes_scores])
34 plt.xticks([1, 2], ['KNN', 'Naive Bayes'])
35 plt.ylabel('Accuracy')
36 plt.show()
```



- b. Using `scipy`, test the hypothesis “*kNN* is statistically superior to naïve Bayes regarding accuracy”, asserting whether is true.

```
1 from scipy.stats import ttest_rel
2
3 stat, p = ttest_rel(knn_scores, naive_bayes_scores)
4 print('Statistics=%.3f, p=%.3f' % (stat, p))
```

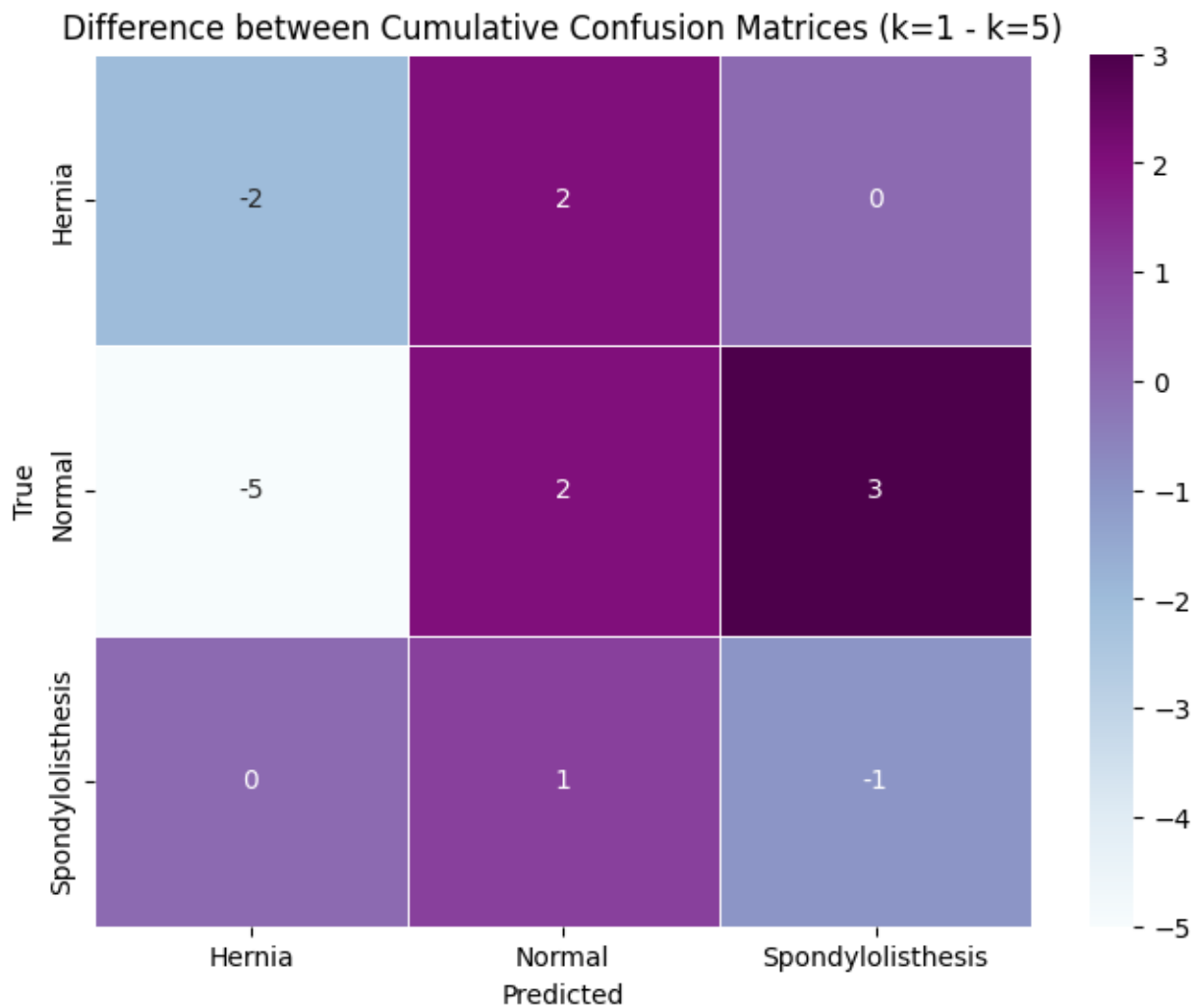
```
Statistics=0.921, p=0.381
```

O *p-value* obtido é bastante superior a 0.1 e a estatística bastante próxima de 1. Assim sendo, não podemos rejeitar a possibilidade de as amostras serem estatisticamente equivalentes, ou seja, a afirmação não é verdadeira.



2. Consider two  $kNN$  predictors with  $k = 1$  and  $k = 5$  (uniform weights, Euclidean distance, all remaining parameters as default). Plot the differences between the two cumulative confusion matrices of the predictors. Comment.

```
1 from sklearn.metrics import confusion_matrix
2 import seaborn as sns
3
4 knn1 = KNeighborsClassifier(n_neighbors=1, weights='uniform', metric='euclidean')
5 knn5 = KNeighborsClassifier(n_neighbors=5, weights='uniform', metric='euclidean')
6
7 skf = StratifiedKFold(n_splits=10, shuffle=True, random_state=0)
8
9 cumulative_cm1 = np.zeros((3, 3))
10 cumulative_cm5 = np.zeros((3, 3))
11
12
13 for train_index, test_index in skf.split(x, y):
14
15     x_train, x_test = df.iloc[train_index, :-1], df.iloc[test_index, :-1]
16     y_train, y_test = df.iloc[train_index, -1], df.iloc[test_index, -1]
17     x_train = x_train.astype(np.float64)
18     x_test = x_test.astype(np.float64)
19     y_train = y_train.astype(str)
20     y_test = y_test.astype(str)
21
22     knn1.fit(x_train, y_train)
23     knn5.fit(x_train, y_train)
24
25     y_pred1 = knn1.predict(x_test)
26     y_pred5 = knn5.predict(x_test)
27
28     cm1 = confusion_matrix(y_test, y_pred1)
29     cm5 = confusion_matrix(y_test, y_pred5)
30
31     cumulative_cm1 += cm1
32     cumulative_cm5 += cm5
33
34
35 diff_cm = cumulative_cm1 - cumulative_cm5
36
37 diff_cm = pd.DataFrame(diff_cm, index=['Hernia', 'Normal', 'Spondylolisthesis'],
38     ↪ columns=['Hernia', 'Normal', 'Spondylolisthesis'])
39 plt.figure(figsize=(8, 6))
40 sns.heatmap(diff_cm, annot=True, cmap="BuPu", linewidths=0.5)
41 plt.title("Difference between Cumulative Confusion Matrices (k=1 - k=5)")
42 plt.ylabel('True')
43 plt.xlabel('Predicted')
44 plt.show()
```



Ao observar a matriz resultante da diferença entre as matrizes de confusão dos dois classificadores, é possível concluir que não há uma diferença significativa de um classificador para o outro, logo, o classificador com  $k = 1$  é tão bom quanto o classificador com  $k = 5$ .

3. Considering the unique properties of `column_diagnosis`, identify three possible difficulties of naïve Bayes when learning from the given dataset.
  - 1- Devido à suposição de independência condicional feita pelo Naïve Bayes entre as variáveis, pode ocorrer uma redução na exatidão do modelo nos casos em que, de fato, existe uma dependência. Isso é evidenciado pelo fato de que as características biomédicas no dataset fornecido estão inter-relacionadas de alguma forma, o que viola essa suposição.
  - 2- O dataset apresenta uma reduzida dimensionalidade, podendo levar a uma redução da exatidão.
  - 3- Observa-se um desequilíbrio no número de observações para cada classe, nomeadamente, existem 150 observações classificadas como "spondylolisthesis" e apenas 60 como "hernia". Isso pode dificultar a aprendizagem e classificação correta das classes minoritárias, resultando em um viés em direção às classes majoritárias.