

I. Pen-and-paper

1)

1 Dist.	Weight	Class		2 Dist.	Weight	Class
2	2.5	0.4 +		1	2.5	0.4 +
3	1.5	0.6667 +		3	1.5	0.6667 +
4	0.5	2 +		4	2.5	0.4 +
5	1.5	0.6667 -		5	1.5	0.6667 -
6	1.5	0.6667 -		6	1.5	0.6667 -
7	1.5	0.6667 -		7	1.5	0.6667 -
8	2.5	0.4 -		8	0.5	2 -
predicted: + TRUE +				predicted: - TRUE +		
3 Dist.	Weight	Class		4 Dist.	Weight	Class
1	1.5	0.6667 +		1	0.5	2 +
2	1.5	0.6667 +		2	2.5	0.4 +
4	1.5	0.6667 +		3	1.5	0.6667 +
5	2.5	0.4 -		5	1.5	0.6667 -
6	2.5	0.4 -		6	1.5	0.6667 -
7	0.5	2 -		7	1.5	0.6667 -
8	1.5	0.6667 -		8	2.5	0.4 -
predicted: - TRUE +				predicted: + TRUE +		
5 Dist.	Weight	Class		6 Dist.	Weight	Class
1	1.5	0.6667 +		1	1.5	0.6667 +
2	1.5	0.6667 +		2	1.5	0.6667 +
3	2.5	0.4 +		3	2.5	0.4 +
4	1.5	0.6667 +		4	1.5	0.6667 +
6	0.5	2 -		5	0.5	2 -
7	2.5	0.4 -		7	2.5	0.4 -
8	1.5	0.6667 -		8	1.5	0.6667 -
predicted: - TRUE -				predicted: - TRUE -		
7 Dist.	Weight	Class		8 Dist.	Weight	Class
1	1.5	0.6667 +		1	2.5	0.4 +
2	1.5	0.6667 +		2	0.5	2 +
3	0.5	2 +		3	1.5	0.6667 +
4	1.5	0.6667 +		4	2.5	0.4 +
5	2.5	0.4 -		5	1.5	0.6667 -
6	2.5	0.4 -		6	1.5	0.6667 -
8	1.5	0.6667 -		7	1.5	0.6667 -
predicted: + TRUE -				predicted: + TRUE -		

$$\text{RECALL} = \text{TP}/(\text{TP}+\text{FN}) = 2/(2+2) = 0.5$$

$$2) \quad P(P|y_1, y_2, y_3) = P(P) \cdot P(y_1, y_2, y_3|P) / P(y_1, y_2, y_3)$$

$$P(N|y_1, y_2, y_3) = P(N) \cdot P(y_1, y_2, y_3|N) / P(y_1, y_2, y_3)$$

Because it's a classifier and the denominator of the two formulas is the same, we don't need to calculate that, and because y_3 is independent from the others we can separate them:

$$P(y_1, y_2, y_3|P) = P(y_1, y_2|P) \cdot P(y_3|P)$$

$$P(y_1, y_2, y_3|N) = P(y_1, y_2|N) \cdot P(y_3|N)$$

Next, we calculate each member of the formula:

$$P(P) = 5/9 \quad | \quad P(N) = 4/9$$

$$P(y_1, y_2|P) = \{A0-2/5, A1-1/5, B0-1/5, B1-1/5\} \quad | \quad P(y_1, y_2|N) = \{A0-0, A1-1/4, B0-1/2, B1-1/4\}$$

$$P(y_3|P) = N(y_3|P_{\text{mean}}, P_{\text{stan.dev.}}) \quad | \quad P(y_3|N) = N(y_3|N_{\text{mean}}, N_{\text{stan.dev.}})$$

$$P_{\text{mean}} = 0.84 \quad | \quad N_{\text{mean}} = 0.975$$

$$P_{\text{stan.dev.}} = 0.251 \quad | \quad N_{\text{stan.dev.}} = 0.17078$$

3) Totalmean = 0.9 | Totalstan.dev. = 0.21794

$$P(P|y_1, y_2, y_3) = 5/9 * P(y_1, y_2|P) * N(y_3|P_{mean}, P_{stan.dev.}) / P(y_1, y_2, y_3)$$

$$P(y_1, y_2, y_3) = P(y_1, y_2) * N(y_3|Totalmean, Totalstan.dev.)$$

$$P(A, 1, 0.8) = 2/9 * N(0.8|0.9, 0.21794) = 0.3661367$$

$$P(P|A, 1, 0.8) = 5/9 * 1/5 * N(0.8|0.84, 0.251) / 0.3661367 = 0.47625084$$

$$P(B, 1, 1) = 2/9 * N(1|0.9, 0.21794) = 0.3661367$$

$$P(P|B, 1, 1) = 5/9 * 1/5 * N(1|0.84, 0.251) / 0.3661367 = 0.39365367$$

$$P(B, 0, 0.9) = 3/9 * N(0.9|0.9, 0.21794) = 0.61017142$$

$$P(P|B, 0, 0.9) = 5/9 * 1/5 * N(0.9|0.84, 0.251) / 0.61017142 = 0.2812767$$

4) $P(P|X) > 0.5 \Rightarrow$ Positive

All the entries are predicted negative, which is an 33% accurate

$P(P|X) > 0.3 \Rightarrow$ Positive

The two first entries are predicted positive and the last one predicted negative, which is true, so the accuracy is 100%

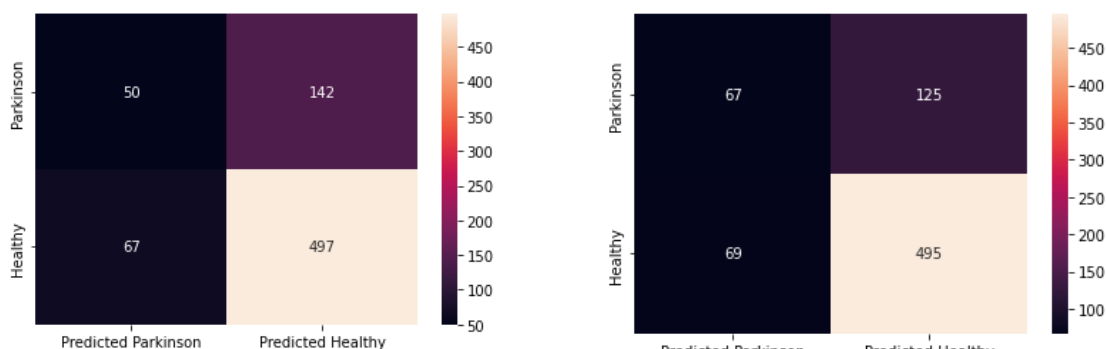
$P(P|X) > 0.7 \Rightarrow$ Positive

All the entries are predicted negative, which is 33% accurate

0.3 decision threshold optimizes testing accuracy

II. Programming and critical analysis

5)



6) $P_{val} = 0.9104476998751558$

Given this value, we can't assume that kNN is statistically superior (in terms of accuracy) to Naïve Bayes.

7) Three reasons for the fact that we can't conclude that kNN is statistically superior to Naïve Bayes can be:

1. The number of neighbors ($k=5$) might not be adjusted to the context of the problem;
2. The data might be too extensive;
3. The fact that there are too few feature dependencies, Naïve Bayes might be favored and its accuracy rises.

III. APPENDIX

5)

```
import pandas as pd
from sklearn.model_selection import StratifiedKFold, cross_val_score
from sklearn.neighbors import KNeighborsClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from scipy import stats

import numpy as np
import copy as cp
import matplotlib.pyplot as plt

import seaborn as sns
from typing import Tuple
from sklearn.metrics import confusion_matrix
from scipy.io.arff import loadarff

from sklearn.metrics import classification_report, confusion_matrix

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

X = df.drop(columns='class')
y = df["class"]

model_neig = KNeighborsClassifier(n_neighbors = 5)
model_gaus = GaussianNB()
kfold = StratifiedKFold(n_splits=10, random_state=0, shuffle=True)
```

Aprendizagem 2022/23
Homework II – Group 103

```
def confusio(model):  
    confusio_matrix = [[0,0],[0,0]]  
  
    for train_ndx, test_ndx in kfold.split(X, y):  
  
        train_X, test_X = X.iloc[train_ndx], X.iloc[test_ndx]  
        train_y, test_y = y.iloc[train_ndx], y.iloc[test_ndx]  
  
        model.fit(train_X, train_y)  
        predicted_y = model.predict(test_X)  
  
        matrix = confusion_matrix(test_y, predicted_y)  
  
        tn, fp, fn, tp = matrix.ravel()  
        confusio_matrix[0][0] += tn  
        confusio_matrix[0][1] += fp  
        confusio_matrix[1][0] += fn  
        confusio_matrix[1][1] += tp  
  
    cm = np.array(confusio_matrix)  
    confusion = pd.DataFrame(cm, index=['Parkinson', 'Healthy'], columns=['Predicted Parkinson', 'Predicted Health  
y'])  
  
    sns.heatmap(confusion, annot=True, fmt='g')  
  
#model_KNN  
confusio(model_neig)  
  
#model_gaussian  
confusio(model_gaus)
```

6)

```
#p_value  
acc_model = cross_val_score(model, X, y, cv=kfold, scoring='accuracy')  
acc_model_1 = cross_val_score(model_1, X, y, cv=kfold, scoring='accuracy')  
res = stats.ttest_rel(acc_model, acc_model_1, alternative='greater')  
print("p1>p2? pval=",res.pvalue)
```

END