

I Pen-and-paper

1.

	P	N
P	8	4
N	3	5

$TP = 5 + 3 = 8$
 $TN = 5$
 $FP = 2 + 2 = 4$
 $FN = 3$

2.

y ¹	y ²	O
B	> 2	N
B	> 2	N
B	> 2	N
B	> 2	N
B	> 2	N
B	> 2	P
B	> 2	P
B	> 2	P
B	≤ 2	P
B	≤ 2	P
B	≤ 2	P
B	≤ 2	N
B	≤ 2	N
A	?	P
A	?	P
A	?	P
A	?	P
A	?	P
A	?	N
A	?	N

y¹

A

P(5/7)

B

N(7/13)

$TP = 5$
 $TN = 7$
 $FP = 13 - 7 = 6$
 $FN = 7 - 5 = 2$

$precision = \frac{TP}{TP + FP} = \frac{5}{5 + 6} = \frac{5}{11}$

$recall = \frac{TP}{TP + FN} = \frac{5}{5 + 2} = \frac{5}{7}$

$F1 = 2 \times \frac{precision \times recall}{precision + recall} = 2 \times \frac{\frac{5}{11} \times \frac{5}{7}}{\frac{5}{11} + \frac{5}{7}} = \frac{5}{9}$

3.

6 último nó do caminho esquerdo poderia conter apenas elemento da mesma classe, terminando a ~~expansão~~ expansão do caminho, ou poderia não haver mais atributos por testar.

4.

$$\begin{aligned}
 I(\text{table}) &= & p(N) &= \frac{9}{20} \\
 &= -\frac{9}{20} \log_2\left(\frac{9}{20}\right) - \frac{11}{20} \log_2\left(\frac{11}{20}\right) & p(P) &= \frac{11}{20} \\
 &= 0,993 \text{ bits}
 \end{aligned}$$

 y_1

$$C_A = \{P, P, P, P, P, N, N\}$$

$$C_B = \{P, P, P, P, P, P, N, N, N, N, N, N, N\}$$

$$I(C_A) = -\frac{5}{7} \log_2\left(\frac{5}{7}\right) - \frac{2}{7} \log_2\left(\frac{2}{7}\right) = 0,8631 \text{ bits}$$

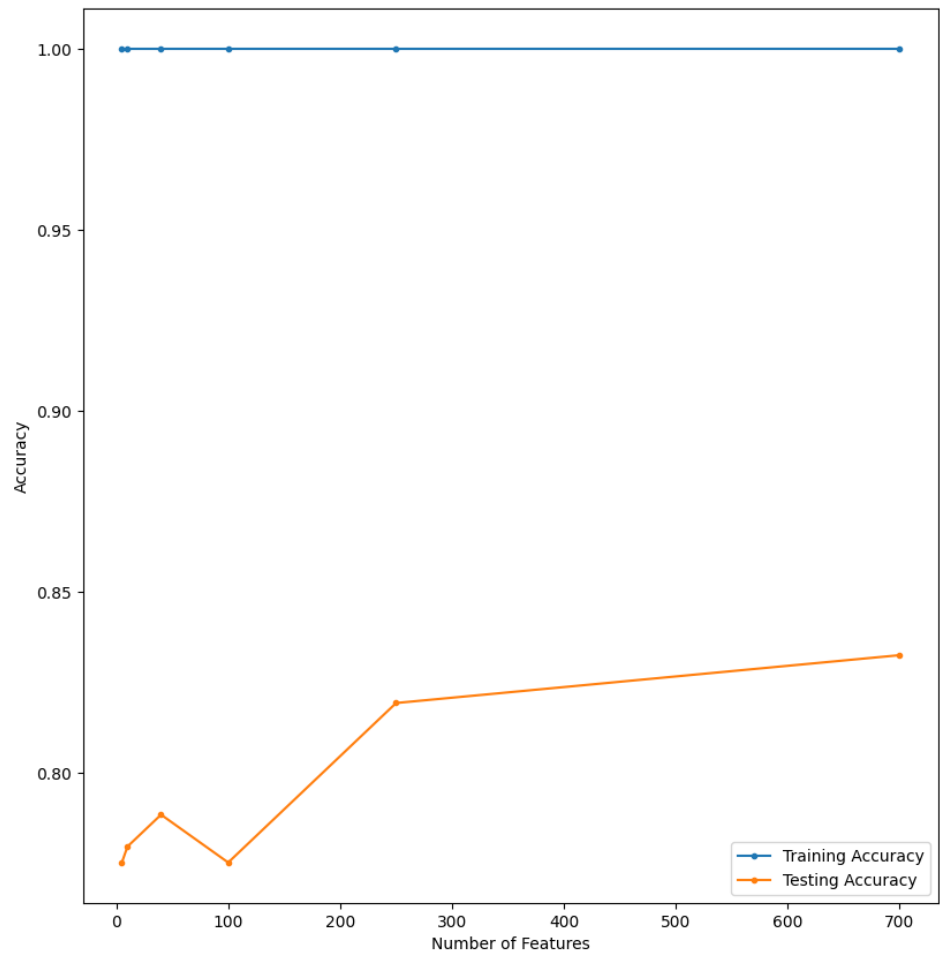
$$I(C_B) = -\frac{6}{13} \log_2\left(\frac{6}{13}\right) - \frac{7}{13} \log_2\left(\frac{7}{13}\right) = 0,9957 \text{ bits}$$

$$\begin{aligned}
 E(y_1) &= \frac{7}{20} \times 0,8631 + \frac{13}{20} \times 0,9957 = \\
 &= 0,94929 \text{ bits}
 \end{aligned}$$

$$g_{\text{ain}}(y_1) = I(\text{table}) - E(y_1) = 0,04371 \text{ bits}$$

II Programming

1.



2. From the gathered results, we can conclude that there is a large difference between the training accuracy and the testing accuracy. The results show a training accuracy of 1, therefore, the model is overfitting to the training set and doesn't generalize so nicely to the testing set.

III Appendix

```
1 from scipy.io.arff import loadarff
2 from sklearn import metrics, tree
3 from sklearn.feature_selection import SelectKBest, mutual_info_classif
4 from sklearn.model_selection import train_test_split
5 import matplotlib.pyplot as plt
6 import pandas as pd
7
8 if __name__ == "__main__":
9     data, meta = loadarff("pd_speech.arff")
10    df = pd.DataFrame(data)
11    df["class"] = df["class"].str.decode("utf-8")
12    df.dropna(inplace=True)
13
14    X, y = df.iloc[:, 0:-1], df.iloc[:, -1]
15    X_train, X_test, y_train, y_test = train_test_split(X,
16                                                         y,
17                                                         train_size=0.7,
18                                                         stratify=y,
19                                                         random_state=1)
20
21    predictor = tree.DecisionTreeClassifier()
22    num_features = (5, 10, 40, 100, 250, 700)
23    training_accuracy = []
24    testing_accuracy = []
25
26    for features in num_features:
27        selector = SelectKBest(mutual_info_classif, k=features)
28        selector.fit(X_train, y_train)
29        X_train_new = selector.transform(X_train)
30        X_test_new = selector.transform(X_test)
31
32        predictor.fit(X_train_new, y_train)
33        y_pred_train = predictor.predict(X_train_new)
34        y_pred_test = predictor.predict(X_test_new)
35
36        training_accuracy += [metrics.accuracy_score(y_train, y_pred_train)]
37        testing_accuracy += [metrics.accuracy_score(y_test, y_pred_test)]
38
39    plt.plot(num_features, training_accuracy, marker=".")
40    plt.plot(num_features, testing_accuracy, marker=".")
41    plt.xlabel("Number of Features")
42    plt.ylabel("Accuracy")
43    plt.legend(["Training Accuracy", "Testing Accuracy"])
44    plt.show()
```