

Aprendizagem 2022
Homework I – Group 019
Diogo Gaspar 99207, Rafael Oliveira 99311

Part I: Pen and paper

1. Draw the training confusion matrix.

After reviewing the decision tree described in the question's statement, we can assert that its confusion matrix is the following:

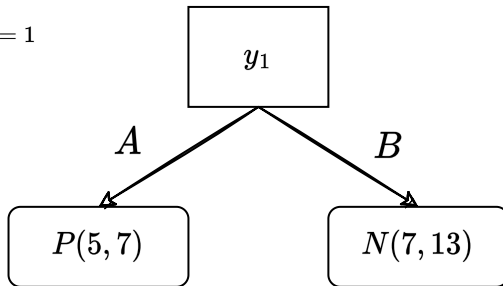
		Real	
		P	N
Projected	P	8	4
	N	3	5

Figure 1: Confusion Matrix

2. Identify the training F1 after a post-pruning of the given tree under a maximum depth of 1.

Post-Pruning

$d = 1$



New confusion matrix:

		Real	
		P	N
Projected	P	5	2
	N	6	7

Precision

Sensitivity
(= Recall)

Figure 2: Post-pruning confusion matrix

We know that:

$$d = 1 : F_1 = \frac{2 * \text{Precision} * \text{Sensitivity}}{\text{Precision} + \text{Sensitivity}}$$

Therefore, we must calculate the associated precision and sensitivity as follows:

$$\text{Sensitivity} = \frac{TP}{TP + FN} = \frac{5}{5 + 6} = \frac{5}{11}, \quad \text{Precision} = \frac{TP}{TP + FP} = \frac{5}{5 + 2} = \frac{5}{7}$$

We can, now, calculate the F1 score:

$$d = 1 : \quad F_1 = \frac{2 * \frac{5}{7} * \frac{5}{11}}{\frac{5}{7} + \frac{5}{11}} = \frac{5}{9}$$

3. Identify two different reasons as to why the left tree path was not further decomposed.

We can find multiple reasons for the left-most tree path not being further decomposed.

The first, perhaps most obvious reason is that we want to **avoid overfitting** the tree to the training data: given a threshold of $\frac{5}{7}$, it's already possible to make a very fair, general assessment about all instances that fall under the $y_1 = A$ condition (asserting they fall onto the P label). It's obviously not guaranteed that not further decomposing the tree will lead to a better labeling come testing time, but generally speaking it ends up being a good practice to avoid overfitting.

Related to the concept of overfitting (specifically, its prevention), there also aren't many instances (only 7) in the left path, which makes it so that further decomposing that path doesn't change the classification error by much (if at all), while also adding unnecessary complexity to the tree - therefore, it's **not likely** that there will be a significant **entropy reduction** by further decomposing the left-most tree path.

4. Compute the information gain of variable y_1 .

We know that:

$$\text{IG}(y_{out} \mid y_1) = H(y_{out}) - H(y_{out} \mid y_1) \quad (1)$$

$$H(Y) = - \sum_{y \in Y} p(y) \log_2 p(y) \quad (2)$$

By looking at the decision tree's confusion matrix (computed in the first question), we can infer that $P(y_{out} = P) = \frac{11}{20}$ and $P(y_{out} = N) = \frac{9}{20}$. After gathering the probabilities above, and considering y_{out} as our *output variable*, with possible values P or N , we'll be able to note that:

$$H(y_{out}) = -\frac{11}{20} \log_2 \frac{11}{20} - \frac{9}{20} \log_2 \frac{9}{20} \quad (3)$$

To finish being able to calculate the previously mentioned information gain, we still need to mention a few missing probabilities associated with the given decision tree:

- $P(y_1 = A) = \frac{7}{20}$, $P(y_1 = B) = \frac{13}{20}$, which can be asserted by looking at the decision tree's confusion matrix and adding the values in each column (the **real** values).
- $P(y_{out} = P \mid y_1 = A) = \frac{5}{7}$, $P(y_{out} = N \mid y_1 = A) = \frac{2}{7}$, which can be gathered by looking at the left-most tree path and checking the true and false P values.
- $P(y_{out} = P \mid y_1 = B) = \frac{6}{13}$, $P(y_{out} = N \mid y_1 = B) = \frac{7}{13}$, which can be gathered by looking at the right-most tree path's leaves, checking each value's true or false values and adding them up.

$$H(y_{out} \mid y_1) = \left[\frac{7}{20} \left(-\frac{5}{7} \log_2 \frac{5}{7} - \frac{2}{7} \log_2 \frac{2}{7} \right) + \frac{13}{20} \left(-\frac{6}{13} \log_2 \frac{6}{13} - \frac{7}{13} \log_2 \frac{7}{13} \right) \right] \quad (4)$$

We therefore have:

$$IG(y_{out} \mid y_1) = 0.04345941113$$

Part II: Programming

5. Using **sklearn**, apply a stratified 70-30 training-testing split with a fixed seed (`random_state=1`), and assess in a single plot the training and testing accuracies of a decision tree with no depth limits (and remaining default behavior) for a varying number of selected features in $\{5,10,40,100,250,700\}$. Feature selection should be performed before decision tree learning considering the discriminative power of the input variables according to mutual information criterion (`mutual_info_classif`).

The code utilized to generate the plot is present in the report's appendix.

6. Why is training accuracy persistently 1? Critically analyze the gathered results.

Since the decision tree has no depth-limit associated, the worst-case scenario sees the tree holding exactly one leaf per training sample - leading to there always being a "correct path" in the tree for a given observation - and thus the training accuracy is always 1, for any train-test split of the dataset. It's the "limit case" for an overfitting situation, where the tree is completely tailored to the training data.

The testing accuracy isn't necessarily 1, of course, just like it can be seen in the plot shown above: in the vast majority of cases, the testing set will hold samples which

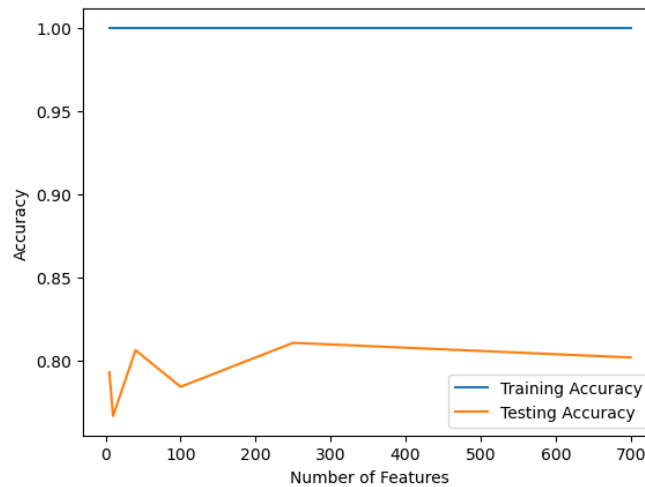


Figure 3: Training and Testing Accuracy vs. Number of Features

differ from all seen in the training set, effectively evicting the "guarantee" of a correct guess from the decision tree. It can, of course, still take a correct guess (with that being our goal), but it's not guaranteed.

Appendix

The code utilized in the first question of the programming section is shown below:

```

1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from scipy.io.arff import loadarff
6 from sklearn.model_selection import train_test_split
7 from sklearn.feature_selection import mutual_info_classif
8 from sklearn.tree import DecisionTreeClassifier
9 from sklearn.metrics import accuracy_score
10 from sklearn.feature_selection import SelectKBest
11
12 FEATURE_AMOUNT = [5, 10, 40, 100, 250, 700]
13 SEED = 1
14 SPLIT = 0.3
15
16 # Loading data
17
18 df = loadarff('data/pd_speech.arff')
19 df = pd.DataFrame(df[0])
20 df['class'] = df['class'].str.decode('utf-8')
21
22 # The first step should be splitting the dataset into stratified training and
23 # testing sets (70-30 split) with a fixed seed (random_state=1).
24
25 X, y = df.drop('class', axis=1), df['class']
26 X_train, X_test, y_train, y_test = train_test_split(
27     X, y,
28     test_size=SPLIT,

```

```

29 random_state=SEED,
30 stratify=y
31 )
32
33 # Now, time to finally assess (in a single plot) both the training and testing
    accuracies
34 # of a decision tree with no depth limits (and remaining default behavior) for a
    varying number of selected features.
35
36 CLF = DecisionTreeClassifier(random_state=SEED)
37 def train_and_test(X_train, X_test, y_train, y_test):
38     CLF.fit(X_train, y_train)
39     y_pred = CLF.predict(X_test)
40     return accuracy_score(y_test, y_pred)
41 def mutual_info_classif_custom(X, y):
42     return mutual_info_classif(X, y, random_state=SEED)
43
44 train_accuracies = []
45 test_accuracies = []
46
47 for n in FEATURE_AMOUNT:
48     selector = SelectKBest(score_func=mutual_info_classif_custom, k=n)
49     selector.fit(X_train, y_train)
50     X_train_new = selector.transform(X_train)
51     X_test_new = selector.transform(X_test)
52     train_accuracies.append(train_and_test(X_train_new, X_train_new, y_train,
        y_train))
53     test_accuracies.append(train_and_test(X_train_new, X_test_new, y_train, y_test)
        )
54
55 plt.plot(FEATURE_AMOUNT, train_accuracies, label='Training Accuracy')
56 plt.plot(FEATURE_AMOUNT, test_accuracies, label='Testing Accuracy')
57 plt.xlabel('Number of features')
58 plt.ylabel('Accuracy')
59 plt.legend()
60 plt.show()

```