

## I. Pen-and-paper

1)

1)  $\#P = 5 + 3 + (8 - 5) = 11$   
 $\#N = 5 + (7 - 5) + (5 - 3) = 9$

	real			P	N	
	P	N				
P	TP	FP	→	P	8	4
N	FN	TN		N	3	5

2)

2) Subárvore que corresponde à decomposição de  $y_2$

$\#P = 3 + (8 - 5) = 6$   
 $\#N = 5 + (5 - 3) = 7$

Como podemos verificar  $N$  é a variável predominante na subárvore, pelo que após o pruning temos:

$y_1$

```

graph TD
    y1["y1"] --> P["P(5/7)"]
    y1 --> N["N(7/13)"]
        
```

A nova matriz de confusão  $\hat{z}$ :

	P	N
P	5	2
N	6	7

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

Recall =  $\frac{TP}{TP+FN} = \frac{5}{5+6} = \frac{5}{11}$

$F_1 = \frac{1}{\frac{1}{2}(\frac{1}{\text{Precision}}) + \frac{1}{2}(\frac{1}{\text{Recall}})} = \frac{5}{9}$

3)

- i. Podemos ter que, para qualquer outra variável, o seu information gain quando  $y_1=A$  seja nulo, pelo que não se justifica realizar a decomposição.
- ii. Ao olharmos para a árvore conseguimos verificar que quando  $y_1=A$  temos uma precisão de 5/7, podendo, desta forma, pressupor que a expansão deste nó levaria a um cenário de overfitting, pelo que o seu crescimento foi interrompido.

4)

4 Pelas linhas anteriores temos que:  
 $\#P = 11$  e  $\#N = 9$  utilizamos class para denotar o output

$$IG(class | y_1) = I(class) - I(class | y_1)$$

$$I(class) = - \sum_{v_i \in class} p(v_i) \log(p(v_i)) =$$

$$= - \left( p(class=P) \log(p(class=P)) + p(class=N) \log(p(class=N)) \right) =$$

$$= - \left( \frac{11}{20} \log\left(\frac{11}{20}\right) + \frac{9}{20} \log\left(\frac{9}{20}\right) \right) = 0,9928$$

$$I(class | y_1) = \sum \frac{|y_i|}{|y_1|} I(class, y_1 = y_i) =$$

$$= \frac{|y_1=A|}{|y_1|} I(class, y_1=A) + \frac{|y_1=B|}{|y_1|} I(class, y_1=B)$$

$$= \frac{7}{20} \left( - \left( \frac{5}{7} \log\left(\frac{5}{7}\right) + \frac{2}{7} \log\left(\frac{2}{7}\right) \right) \right) +$$

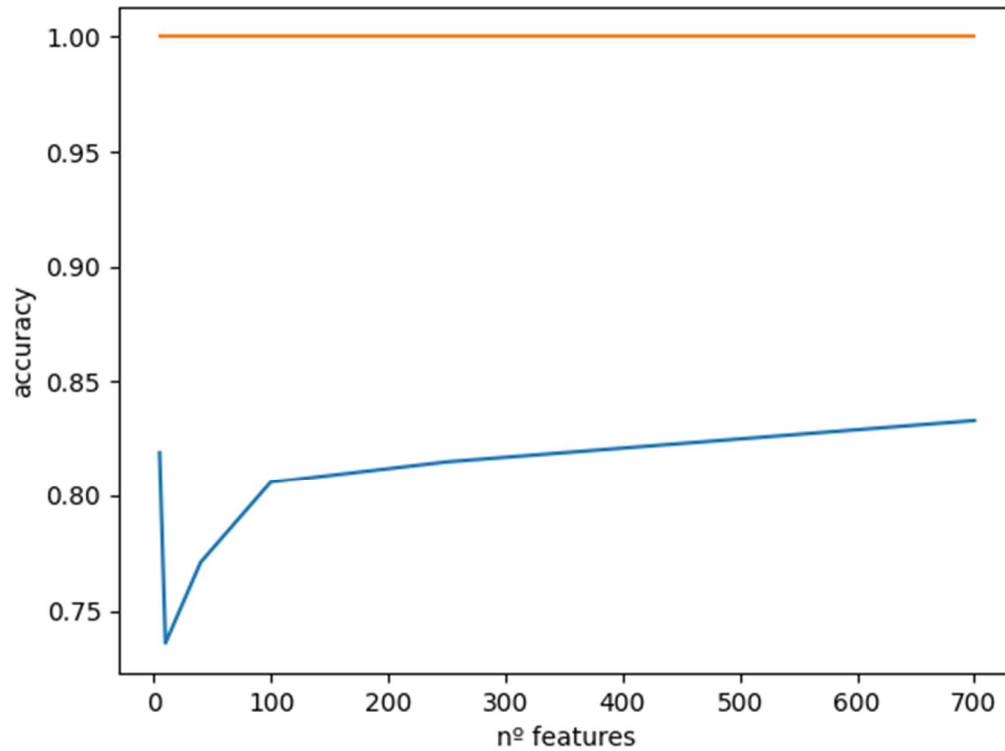
$$\frac{13}{20} \left( - \left( \frac{6}{13} \log\left(\frac{6}{13}\right) + \frac{7}{13} \log\left(\frac{7}{13}\right) \right) \right) =$$

$$= 0,9493$$

$$IG(class | y_1) = 0,9928 - 0,9493 = 0,0435$$

## II. Programming and critical analysis

5)



6)

Como não estabelecemos nenhum limite para a profundidade da árvore de decisão, esta irá ramificar-se indefinidamente, adaptando-se perfeitamente ao training set, fazendo com que a training accuracy seja 1. Ao observarmos o gráfico verificamos que a precisão no test set é significativamente inferior à do training set ( $<1$ , como seria de esperar), pelo que podemos concluir que estamos perante um cenário de overfitting.

### III. APPENDIX

```
##### Importing required libraries #####
import pandas as pd
from scipy.io.arff import loadarff
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import mutual_info_classif
from sklearn import tree, metrics
import matplotlib.pyplot as plt

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

##### Creating the training-testing split #####
X, y = df.drop('class', axis=1), df['class']
X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=0.7, stratify=y,
random_state=1)

##### Feature ranking based on discriminative power #####
dp = mutual_info_classif(X, y, random_state=1)
dict1, count = {}, 0

for feature in X_train.columns.values:
    dict1[feature] = dp[count]
    count += 1

dict2 = dict(sorted(dict1.items(), key=lambda item: item[1], reverse=True))
features = list(dict2.keys())

##### Running classifier and attesting results #####
train, test = [], []
predictor = tree.DecisionTreeClassifier(random_state=1)

for n in [5, 10, 40, 100, 250, 700]:
    predictor.fit(X_train[features[0:n]], y_train)
    y_pred1 = predictor.predict(X_test[features[0:n]])
    y_pred2 = predictor.predict(X_train[features[0:n]])
    test.append(round(metrics.accuracy_score(y_test, y_pred1), 3))
    train.append(round(metrics.accuracy_score(y_train, y_pred2), 3))

##### Creating plot #####
plt.plot([5, 10, 40, 100, 250, 700], test)
plt.plot([5, 10, 40, 100, 250, 700], train)
plt.xlabel('n° features')
plt.ylabel('accuracy')
plt.show()
```

END