

I. Pen-and-paper

I

1.)

	P	N
P	8	3
N	4	5

2.) $F_1 = \frac{2}{\frac{1}{8} + \frac{1}{11}} = \frac{2}{\frac{1}{5.6}} = \frac{5.6}{1} = 5.6$

$P_P = \frac{5.6}{7} \quad P_N = \frac{5.6}{5+6}$

3.) gain = 0, for example all $(y_1 | y_2 = A)$ had no doubt on the decision ($P(y_2 | y_1 = A) = 1$)
gain not high enough to justify possible future overfitting problems

4.) $gain(y_1) = E(C) - E(C | y_1)$

$E(C = P) = \frac{5+3+3}{20} = \frac{11}{20}$

$E(C = N) = \frac{2+5+2}{20} = \frac{9}{20}$

$E(C) = -\frac{11}{20} \log_2 \frac{11}{20} - \frac{9}{20} \log_2 \frac{9}{20}$

≈ 0.992774

$E(C | y_1 = A) = -\frac{5}{11} \log_2 \frac{5}{11} - \frac{3}{11} \log_2 \frac{3}{11} = 0.863121$

$E(C | y_1 = B) = -\frac{3}{11} \log_2 \frac{3}{11} - \frac{6}{11} \log_2 \frac{6}{11} = 0.997777$

$E(C | y_1) = \frac{3}{20} \times 0.863121 + \frac{17}{20} \times 0.997777$

$gain(y_1) = E(C) - E(C | y_1) = 0.99 - 0.94 = 0.05759$

	y_1	y_2	
5x	A	1	P
2x	A	3	N
3x	B	3	P
5x	B	3	N
3x	B	1	P
2x	B	1	N

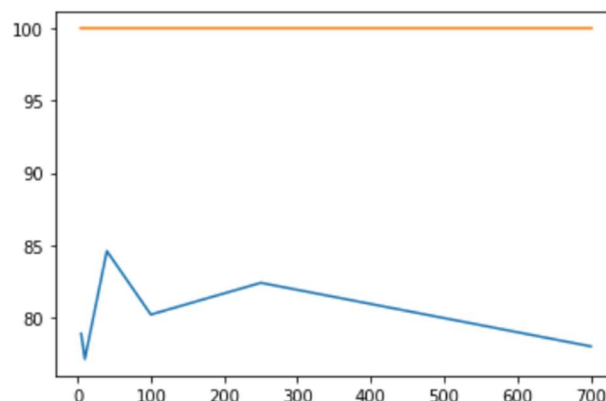
Decision Tree:

```

graph TD
    Y1((y1)) -- A --> NodeA
    Y1 -- B --> NodeB
    NodeA -- P --> LeafA1[5]
    NodeA -- N --> LeafA2[3]
    NodeB -- P --> LeafB1[3]
    NodeB -- N --> LeafB2[6]
  
```

II. Programming and critical analysis

1)



- 2) While when testing, the algorithm is guessing the answer, when evaluating on already trained instances it is certain about its class. So, unless it undervalues any train line, it will guess its class correctly. Which could also show some signs to overfitting.

III. APPENDIX

```
import pandas as pd, matplotlib.pyplot as plt
from scipy.io.arff import loadarff
from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import train_test_split
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import mutual_info_classif

data = loadarff('drive/MyDrive/ML/pd_speech.arff')
df = pd.DataFrame(data[0])
df['class'] = df['class'].str.decode('utf-8')
X = df.drop('class', axis=1)
y = df['class']

Kvalues = [5,10,40,100,250,700]
test_scores = []
train_scores = []

for i in Kvalues:
    X_new = SelectKBest(mutual_info_classif, k=i).fit_transform(X, y)
    X_train, X_test, y_train, y_test = train_test_split(X_new, y, train_size = 0.7, random_state = 1)

    predictor = DecisionTreeClassifier(random_state=1)
    predictor.fit(X_train, y_train)

    train_scores.append(predictor.score(X_train, y_train)*100)
    test_scores.append(predictor.score(X_test, y_test)*100)

plt.plot(Kvalues, test_scores)
plt.plot(Kvalues, train_scores)
plt.show()
```

END