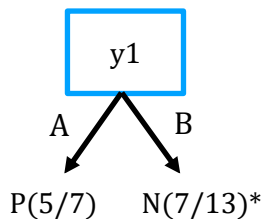


I. Pen-and-paper

1) $TP = 5 + 3 = 8$; $FP = 7 - 5 + 5 - 3 = 4$; $TN = 5$; $FN = 8 - 5 = 3$

Predicted/Actual	P_{actual}	N_{actual}
$P_{\text{predicted}}$	8_{TP}	4_{FP}
$N_{\text{predicted}}$	3_{FN}	5_{TN}

2) Pruning:



$$\begin{aligned}\#(y1=B) &= 8 + 5 = 13 \\ \#P &= 3 + (8 - 5) = 6 \\ \#N &= 13 - 6 = 7\end{aligned}$$

$$\begin{aligned}\text{Precision} &= \frac{TP_{\text{prun}}}{TP_{\text{prun}} + FP_{\text{prun}}} = \\ &= \frac{5}{5 + (7 - 5)} = \frac{5}{7}\end{aligned}$$

$$\begin{aligned}\text{Recall} &= \frac{TP_{\text{prun}}}{TP_{\text{prun}} + FN_{\text{prun}}} = \\ &= \frac{5}{5 + (13 - 7)} = \frac{5}{11}\end{aligned}$$

* Classified as N because when considering only $y1=B$, $\#N=7 > 6 = \#P$

$$\text{Training } F1 = \frac{2 \times P \times R}{P + R} = \frac{5}{9}$$

3) We found two possibilities:

- So as to avoid overfitting because the larger the tree the greater the chances of the tree capturing the noise of the training data and failing to generalize to the test data. So some chosen criterion might have stopped the further expansion of the tree to avoid overfitting.
- The considered features lacked discriminative power. So, for example, if information gain was chosen as the selection criteria to build the tree, then for all the considered features their information gain was either 0 or failed to pass some considered threshold.

4) $\#P_{\text{Actual}} = TP + FN = 8 + 3 = 11$ $\#N_{\text{actual}} = TN + FP = 5 + 4 = 9$ $\#TOTAL = TP + FP + TN + FN = 20$

$$I(\text{table}) = -\frac{11}{20} \log_2 \frac{11}{20} - \frac{9}{20} \log_2 \frac{9}{20} \cong 0.9928 \text{ bits}$$

$$I(C_A) = -\frac{5}{7} \log_2 \frac{5}{7} - \frac{7-5}{7} \log_2 \frac{7-5}{7} \cong 0.8631 \text{ bits}$$

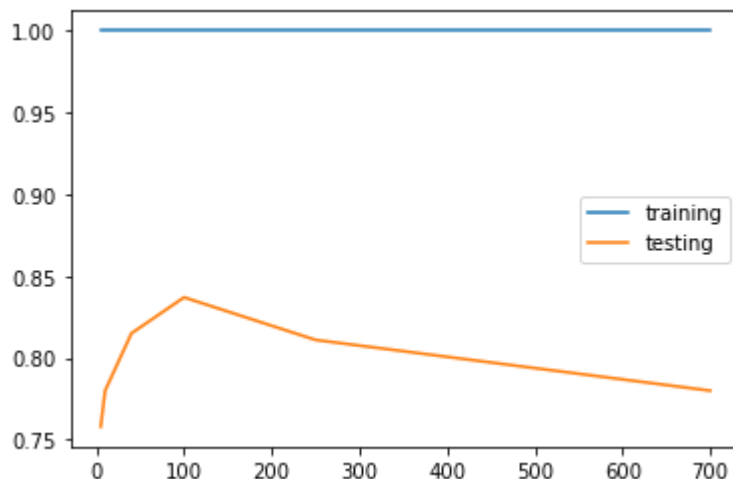
$$I(C_B) = -\frac{7}{13} \log_2 \frac{7}{13} - \frac{13-7}{13} \log_2 \frac{13-7}{13} \cong 0.9957 \text{ bits}$$

$$I(y_1) = \frac{\#(y1=A)}{\#TOTAL} \times I(C_A) + \frac{\#(y1=B)}{\#TOTAL} \times I(C_B) = \frac{20-13}{20} \times I(C_A) + \frac{13}{20} \times I(C_B) = 0.9493 \text{ bits}$$

$$\text{Gain}(y1) = I(\text{table}) - E(y1) = 0.0435 \text{ bits}$$

II. Programming and critical analysis

1) We plotted the training and testing accuracies of our decision tree.



2) Why the training accuracy is consistently 1:

- We chose the K best features for the training data. This means that even with one feature, chosen from a pool of 753, the training accuracy is already a whopping 0.998. Also, since no limit was imposed on the depth of the tree and no criteria was used to trim it, the model overfits the training data. Thus, for the considered number of features the training accuracy is invariably 1.

Critical analysis of the gathered results:

- For a small number of features the bias is larger. As the number of features increased along with the model complexity, the bias went down and the variance and testing accuracy went up. Eventually, when going from 100 to 250 features, the variance got too big and we started overfitting the training data and getting worse accuracy on the test data, the model got too complex and failed to generalize so well. From 200 to 700 features it somewhat stabilized. We still see a decline in testing accuracy, though. About 100 features seem to offer the best tradeoff between bias and variance and for the considered number of features it gave the best testing accuracy (about 0.85).

III. APPENDIX

```
from scipy.io.arff import loadarff
import pandas as pd
from sklearn.model_selection import train_test_split
import numpy as np
from sklearn import tree, metrics
from sklearn.feature_selection import SelectKBest, mutual_info_classif

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

#print(df.value_counts('class'))
#print(df.isna().sum())

number_of_features = [5, 10, 40, 100, 250, 700]
x, y = df.drop('class', axis=1), np.ravel(df['class'])
x_train, x_test, y_train, y_test = train_test_split(x, y, train_size = 0.7,
random_state=1)
data_accuracy_results = {'training':[], 'testing':[]}
for i in number_of_features:
    select = SelectKBest(score_func = mutual_info_classif, k=i)
    select.fit(x_train, y_train)
    x_train_new = select.transform(x_train)
    x_test_new = select.transform(x_test)
    predictor = tree.DecisionTreeClassifier()
    predictor.fit(x_train_new, y_train)
    y_pred_train = predictor.predict(x_train_new)
    y_pred_test = predictor.predict(x_test_new)
    data_accuracy_results['training'].append(metrics.accuracy_score(y_train,
y_pred_train))
    data_accuracy_results['testing'].append(metrics.accuracy_score(y_test,
y_pred_test))

df_accuracy_results = pd.DataFrame(data_accuracy_results, index = number_of_features)
df_accuracy_results.plot.line()
```

END