



LEIC-T 2023/2024
Aprendizagem - Machine Learning

Homework I - Group 007

Miguel Teixeira - 103449

Rodrigo Alves - 103299

I

a) $x_1 = (1, 3, 4, 6)$
 $x_2 = (-30, -10, 0, 20)$

Proposi:

$$\text{CORR}(x_1, x_2) = \frac{\text{COV}(x_1, x_2)}{G(x_1) G(x_2)} \quad \bar{x}_1 = 3,5 \quad \bar{x}_2 = -5$$
$$\text{COV}(x_1, x_2) = \frac{\sum_{k=1}^4 (x_{1k} - \bar{x}_1)(x_{2k} - \bar{x}_2)}{4} = \frac{(1-3,5)(-30-(-5)) + (3-3,5)(-10-(-5))}{4}$$
$$+ \frac{(4-3,5)(0-(-5)) + (6-3,5)(20-(-5))}{4} = \frac{62,5 + 2,5 + 2,5 + 62,5}{4} = 32,5$$
$$G(x_1) = \sqrt{\frac{1}{4} \sum_{k=1}^4 (x_{1k} - \bar{x}_1)^2} = \sqrt{\frac{\sum_{k=1}^4 (x_{1k} - \bar{x}_1)^2}{4}} = \sqrt{\frac{(1-3,5)^2 + (3-3,5)^2 + (4-3,5)^2 + (6-3,5)^2}{4}}$$
$$= \frac{\sqrt{6,25 + 0,25 + 0,25 + 6,25}}{2} = \frac{\sqrt{13}}{2}$$
$$G(x_2) = \sqrt{\frac{(-30-(-5))^2 + (-10-(-5))^2 + (0-(-5))^2 + (20-(-5))^2}{4}} = \frac{\sqrt{625 + 25 + 25 + 625}}{2}$$
$$= \frac{10\sqrt{13}}{2}$$
$$\text{CORR}(x_1, x_2) = \frac{32,5 \times 4}{10 \times 13} = \frac{130}{130} = 1$$

Spearman:

$$R(x_1) = (1, 2, 3, 4)$$

$$R(\bar{x}_1) = 2,5 \quad R(\bar{x}_2) = 2,5$$

$$R(x_2) = (1, 2, 3, 4)$$

$$\text{Spearman} = \frac{\text{Cov}(R(x_1), R(x_2))}{S(R(x_1)) S(R(x_2))}$$

$$\begin{aligned} \text{Cov}(R(x_1), R(x_2)) &= \frac{(1-2,5)(1-2,5) + (2-2,5)(2-2,5) + (3-2,5)(3-2,5) + (4-2,5)(4-2,5)}{4} \\ &= \frac{2,25 + 0,25 + 0,25 + 2,25}{4} = \frac{5}{4} \end{aligned}$$

$$S(R(x_1)) = \sqrt{\frac{(1-2,5)^2 + (2-2,5)^2 + (3-2,5)^2 + (4-2,5)^2}{4}} = \frac{\sqrt{5}}{2} = S(R(x_2))$$

In both cases, the correlation coefficients are 1, suggesting that the two variables, x_1 and x_2 , are perfectly positively correlated.

$$\text{Spearman} = \frac{\frac{5}{4}}{\frac{\sqrt{5}}{2} \times \frac{\sqrt{5}}{2}} = 1$$

R.: Spearman and the Pearson correlation have the same value, that is 1, so x_1 and x_2 are the same

6)

Pearson correlation:

$$\bar{x}_1 = 3,5 \quad \bar{x}_2 = 13,875$$

$$\text{Cov}(x_1, x_2) = \frac{\text{Cov}(x_1, x_2)}{S(x_1) S(x_2)}$$

$$\begin{aligned} \text{Cov}(x_1, x_2) &= \frac{(1-3,5)(-3-13,875) + (3-3,5)(-9,5-13,875) + (4-3,5)(29-13,875) + (6-3,5)(30-13,875)}{4} \\ &= \frac{42,1875 + 7,1875 + 7,5625 + 40,3125}{4} = 24,3125 \end{aligned}$$

$$S(x_1) = \sqrt{\frac{(1-3,5)^2 + (3-3,5)^2 + (4-3,5)^2 + (6-3,5)^2}{4}} = \frac{\sqrt{13}}{2}$$

$$S(x_2) = \sqrt{\frac{(-3-13,875)^2 + (-9,5-13,875)^2 + (29-13,875)^2 + (30-13,875)^2}{4}} = \frac{31,308}{2}$$

$$\begin{aligned} \text{Corr}(x_1, x_2) &= \frac{24,3125}{\frac{\sqrt{13}}{2} \times \frac{31,308}{2}} = 0,862 \end{aligned}$$

Spearman Correlation:

$$R(x_1) = (1, 2, 3, 4)$$

$$R(x_2) = (1, 2, 3, 4)$$

A partir do exercício anterior (1a), em que os ranks não os mesmos, concluímos que:

$$\text{cov}(R(x_1), R(x_2)) = \frac{5}{4}$$

$$R(\bar{x}_1) = 2,5$$

$$R(\bar{x}_2) = 2,5$$

$$\Delta(R(x_1)) = \Delta(R(x_2)) = \frac{\sqrt{5}}{2}$$

$$\text{corr}(R(x_1), R(x_2)) = \frac{\text{cov}(R(x_1), R(x_2))}{\Delta(R(x_1)) \Delta(R(x_2))} = 1$$

R.: Spearman and the Pearson Correlation have different values, so x_1 and x_2 are different

2
a)

F_1	F_2	F_3	F_4	Output
c	a	b	x	n
a	a	c	a	t
a	b	b	a	t
c	b	c	x	m
a	b	b	c	f

$$P(n) = P(m) = P(f) = \frac{1}{5}$$

$$P(t) = \frac{2}{5}$$

$$I(\text{table}) = -\frac{3}{5} \log \frac{1}{5} - \frac{2}{5} \log \frac{2}{5} = 1,922 \text{ bit}$$

$$F_1: C_c = \{n, m\}$$

$$C_a = \{t, f\}$$

$$P(n) = \frac{1}{2}$$

$$P(m) = \frac{1}{2}$$

$$P(t) = \frac{2}{3}$$

$$P(f) = \frac{1}{3}$$

$$I = -\frac{1}{2} \log \frac{1}{2} = 1$$

$$I = -\frac{2}{3} \log \frac{2}{3} - \frac{1}{3} \log \frac{1}{3} = 0,9183$$

$$I(F_1) = \frac{2}{5} \times 1 + \frac{3}{5} \times 0,9183 = 0,95098$$

$$\text{gain}(F_1) = 1,922 - 0,95098 = 0,97102$$

$$F_2: C_a = \{n, t\}$$

$$P(n) = \frac{1}{2}$$

$$P(t) = \frac{1}{2}$$

$$I = -\frac{2}{2} \log_2 \frac{1}{2} = 1$$

$$C_b = \{t, m, f\}$$

$$P(t) = \frac{1}{3}$$

$$P(m) = \frac{1}{3}$$

$$P(f) = \frac{1}{3}$$

$$I = -\frac{3}{3} \log_2 \frac{1}{3} = 1,585$$

$$I(F_2) = \frac{2}{5} \times 1 + \frac{3}{5} \times 1,585 = 1,351$$

$$\text{gain}(F_2) = 1,922 - 1,351 = 0,571$$

$$F_3: C_b = \{n, t, f\}$$

$$P(n) = \frac{1}{3}$$

$$P(t) = \frac{1}{3}$$

$$P(f) = \frac{1}{3}$$

$$I = -\log_2 \frac{1}{3} = 1,585$$

$$C_c = \{t, m\}$$

$$P(t) = \frac{1}{2}$$

$$P(m) = \frac{1}{2}$$

$$I = -\frac{2}{2} \log_2 \frac{1}{2} = 1$$

$$I(F_3) = \frac{2}{5} \times 1 + \frac{3}{5} \times 1,585 = 1,351$$

$$\text{gain}(F_3) = 1,922 - 1,351 = 0,571$$

$$F_4: C_a = \{n, m\}$$

$$P(n) = \frac{1}{2}$$

$$P(m) = \frac{1}{2}$$

$$I = -\frac{2}{2} \log_2 \frac{1}{2} = 1$$

$$C_b = \{t, f\}$$

$$P(t) = 1$$

$$I = -\log_2 1 = 0$$

$$C_c = \{f\}$$

$$P(f) = 1$$

$$I = -\log_2 1 = 0$$

$$I(F_4) = \frac{2}{5} \times 1 + \frac{2}{5} \times 0 + \frac{1}{5} \times 0 = \frac{2}{5}$$

$$\text{gain}(F_4) = 1,922 - 0,4 = 1,522$$

F_4 is the root because it provides the best gain

(b) F_4 is the first node in the tree

$F_4 = x$

F_1	F_2	F_3	Output
c	a	b	m
c	b	c	m

$F_4 = a$

F_1	F_2	F_3	Output
a	a	c	t
a	b	b	t

Done

$F_4 = c$

F_1	F_2	F_3	Output
a	b	b	f

Done

Still has uncertainty, we have to decide between testing F_1 , F_2 or F_3



$$\begin{aligned}
 & C_c \{m, m\} \\
 & E(F_1(C_c)) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) \\
 & = \boxed{1} \\
 & E(F_1) = \boxed{1}
 \end{aligned}$$

$C_a \{m\}$

$C_b \{m\}$

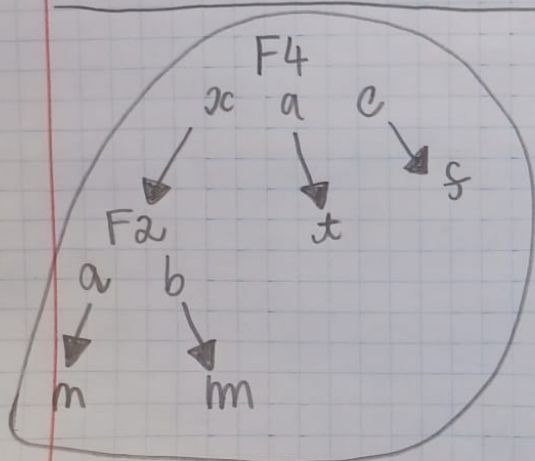
$$E(F_2) = E_{F_1}(C_a) + E_{F_1}(C_b)$$

$$= \boxed{0}$$

$$\begin{aligned}
 & C_b \{m\} \quad E_{F_2}(C_b) = 0 \\
 & C_c \{m\} \quad E_{F_2}(C_c) = 0 \\
 & E(F_2) = \boxed{0}
 \end{aligned}$$

$$E(F_2) = E(F_3) = 0,$$

So we can choose either one, we chose F_2



$F_4 = x$			$F_2 = b$		
F_4	F_3	Output	F_1	F_3	Output
c	b	m	c	c	m

(c)

	m	m	t	s
m	1	0	0	0
m	0	1	0	0
t	0	0	2	0
s	0	0	0	1

There are no predictions or false data

III Software Experiments

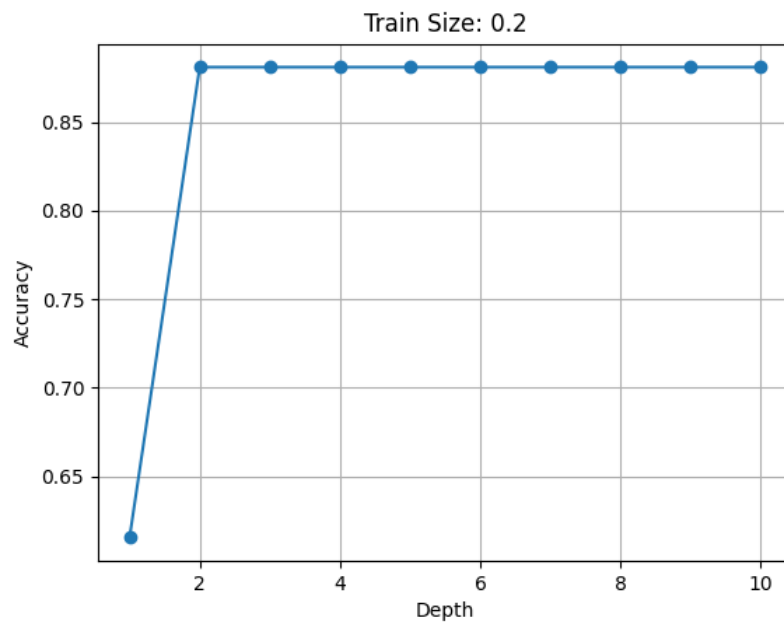
```

1  import matplotlib.pyplot as plt
2  from sklearn import metrics, datasets, tree
3  from sklearn.model_selection import train_test_split
4
5  # Load the wine dataset
6  wine = datasets.load_wine()
7  X, y = wine.data, wine.target
8
9  # Define the different train sizes
10 train_sizes = [0.2, 0.5, 0.7]
11 group_number = 7
12
13 # Loop through different train sizes
14 for train_size in train_sizes:
15     # Split the data
16     X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=train_size, stratify=y, random_state=int(group_number))
17
18     # Create a decision tree classifier and evaluate its accuracy for different depths
19     depths = range(1, 11) # Test depths from 1 to 10
20     accuracies = []
21
22     for depth in depths:
23         # Create a decision tree classifier with the current depth
24         clf = tree.DecisionTreeClassifier(max_depth=depth, random_state=int(group_number))
25         clf.fit(X_train, y_train)
26
27         # Predict on the test set
28         y_pred = clf.predict(X_test)
29
30         # Calculate accuracy
31         accuracy = metrics.accuracy_score(y_test, y_pred)
32         accuracies.append(accuracy)
33
34         # Print depth and accuracy
35         print(f"Train Size: {train_size}, Depth: {depth}, Accuracy: {accuracy:.4f}")
36
37     # Plot accuracy vs. depth for the current train size
38     plt.figure()
39     plt.plot(depths, accuracies, marker='o', linestyle='-')
40     plt.title(f"Train Size: {train_size}")
41     plt.xlabel("Depth")
42     plt.ylabel("Accuracy")
43     plt.grid(True)
44     plt.show()
45

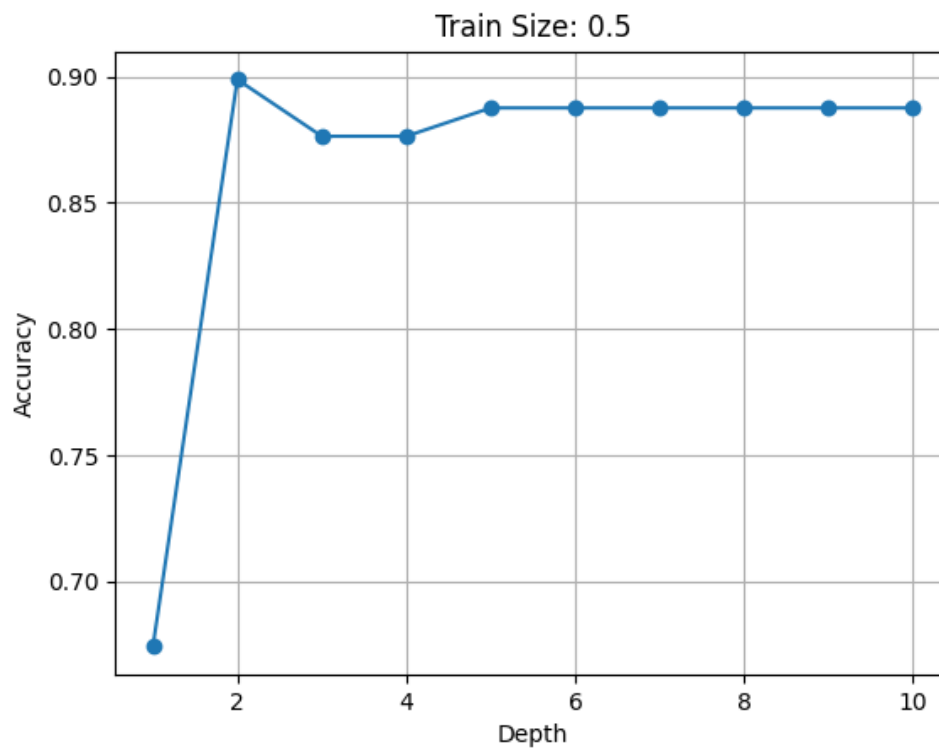
```

1. (a)

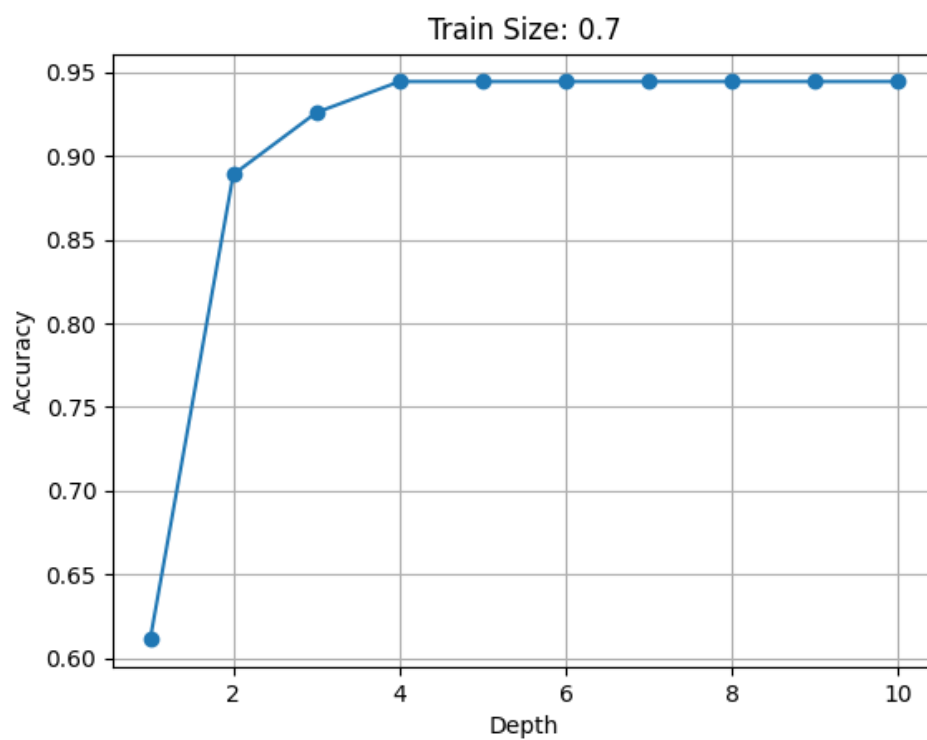
Train Size: 0.2, Depth: 1, Accuracy: 0.6154
Train Size: 0.2, Depth: 2, Accuracy: 0.8811
Train Size: 0.2, Depth: 3, Accuracy: 0.8811
Train Size: 0.2, Depth: 4, Accuracy: 0.8811
Train Size: 0.2, Depth: 5, Accuracy: 0.8811
Train Size: 0.2, Depth: 6, Accuracy: 0.8811
Train Size: 0.2, Depth: 7, Accuracy: 0.8811
Train Size: 0.2, Depth: 8, Accuracy: 0.8811
Train Size: 0.2, Depth: 9, Accuracy: 0.8811
Train Size: 0.2, Depth: 10, Accuracy: 0.8811



Train Size: 0.5, Depth: 1, Accuracy: 0.6742
Train Size: 0.5, Depth: 2, Accuracy: 0.8989
Train Size: 0.5, Depth: 3, Accuracy: 0.8764
Train Size: 0.5, Depth: 4, Accuracy: 0.8764
Train Size: 0.5, Depth: 5, Accuracy: 0.8876
Train Size: 0.5, Depth: 6, Accuracy: 0.8876
Train Size: 0.5, Depth: 7, Accuracy: 0.8876
Train Size: 0.5, Depth: 8, Accuracy: 0.8876
Train Size: 0.5, Depth: 9, Accuracy: 0.8876
Train Size: 0.5, Depth: 10, Accuracy: 0.8876

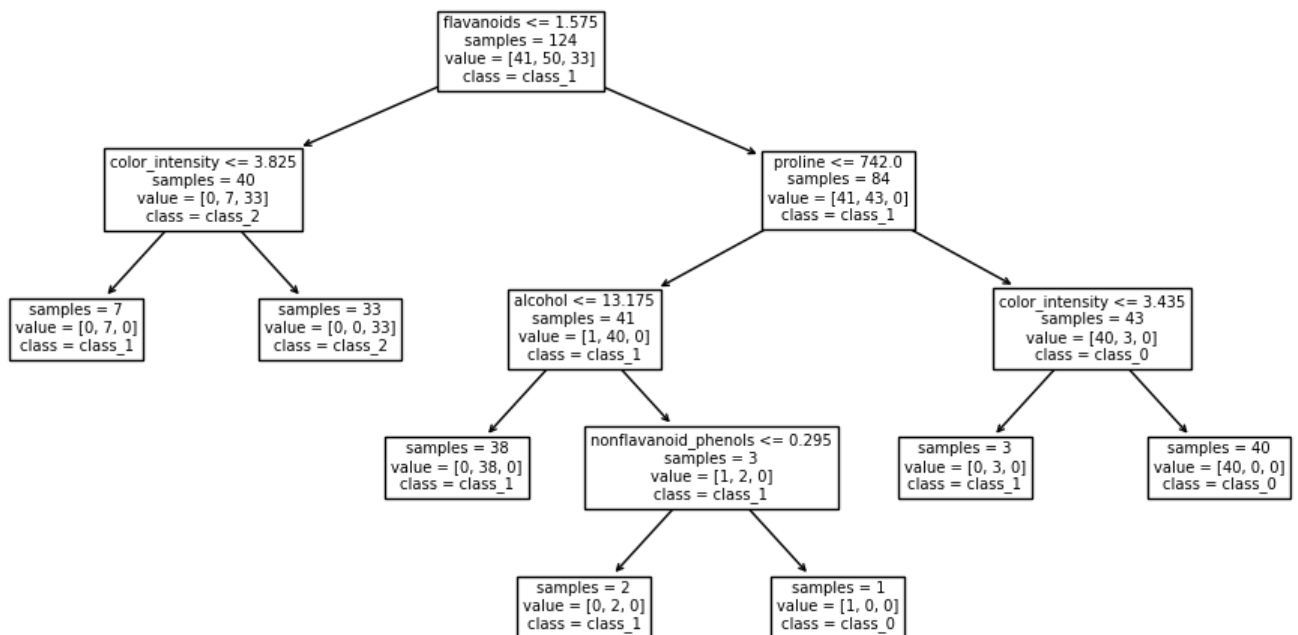


Train Size: 0.7, Depth: 1, Accuracy: 0.6111
Train Size: 0.7, Depth: 2, Accuracy: 0.8889
Train Size: 0.7, Depth: 3, Accuracy: 0.9259
Train Size: 0.7, Depth: 4, Accuracy: 0.9444
Train Size: 0.7, Depth: 5, Accuracy: 0.9444
Train Size: 0.7, Depth: 6, Accuracy: 0.9444
Train Size: 0.7, Depth: 7, Accuracy: 0.9444
Train Size: 0.7, Depth: 8, Accuracy: 0.9444
Train Size: 0.7, Depth: 9, Accuracy: 0.9444
Train Size: 0.7, Depth: 10, Accuracy: 0.9444



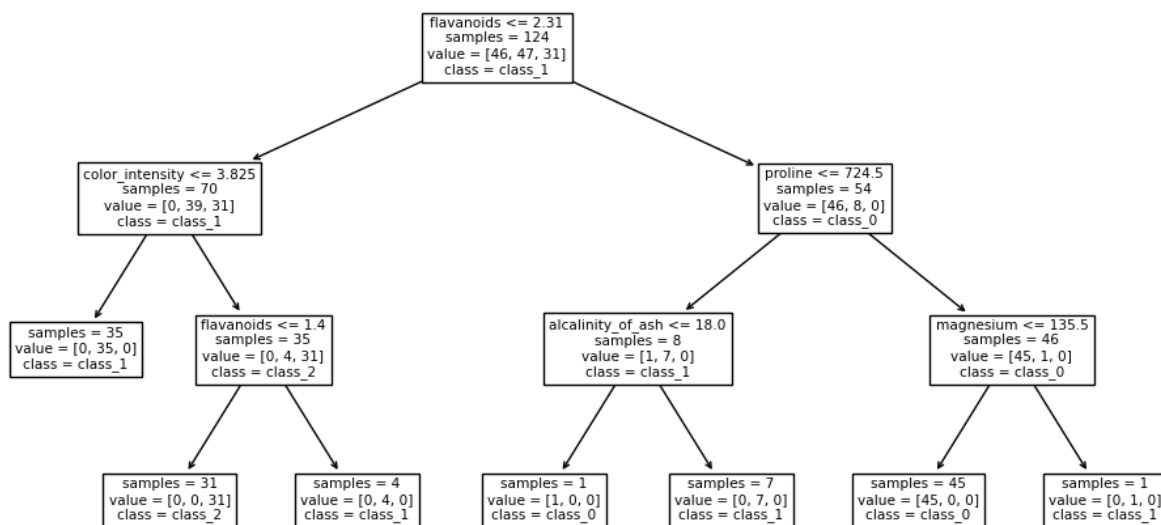
(b)

```
1 import matplotlib.pyplot as plt
2 from sklearn import datasets, tree
3 from sklearn.model_selection import train_test_split
4
5 # Load the Wine dataset
6 wine = datasets.load_wine()
7 X, y = wine.data, wine.target
8
9 # Define the classifier (Decision Tree with entropy)
10 predictor = tree.DecisionTreeClassifier(criterion='entropy')
11
12 # Split the data
13 train_size = 0.7
14 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=train_size, stratify=y, random_state=7)
15
16 # Train the classifier on the training data
17 predictor.fit(X_train, y_train)
18
19 # Plot the decision tree
20 figure = plt.figure(figsize=(12, 6))
21 tree.plot_tree(predictor, feature_names=wine.feature_names, class_names=wine.target_names, impurity=False)
22
23 # Save the decision tree as an image (e.g., PNG)
24 plt.savefig('decision_tree.png')
25
26 # Display the plot (optional)
27 plt.show()
28
```



(c)

```
1 import matplotlib.pyplot as plt
2 from sklearn import datasets, tree
3 from sklearn.model_selection import train_test_split
4
5 # Load the Wine dataset
6 wine = datasets.load_wine()
7 X, y = wine.data, wine.target
8
9 # Define the classifier (Decision Tree with entropy)
10 predictor = tree.DecisionTreeClassifier(criterion='entropy')
11
12 # Split the data
13 train_size = 0.7
14
15 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size=train_size, random_state=7)
16
17 # Train the classifier on the training data
18 predictor.fit(X_train, y_train)
19
20 # Plot the decision tree
21 figure = plt.figure(figsize=(12, 6))
22 tree.plot_tree(predictor, feature_names=wine.feature_names, class_names=wine.target_names, impurity=False)
23 plt.show()
24
```



The results may be different when not using `stratify=y` because, without stratification, the class distribution in the training and testing sets may not be representative of the overall dataset, potentially leading to biased model evaluations, while `stratify=y` ensures a stratified fashion, preserving the same proportion of target classes in both subsets.