

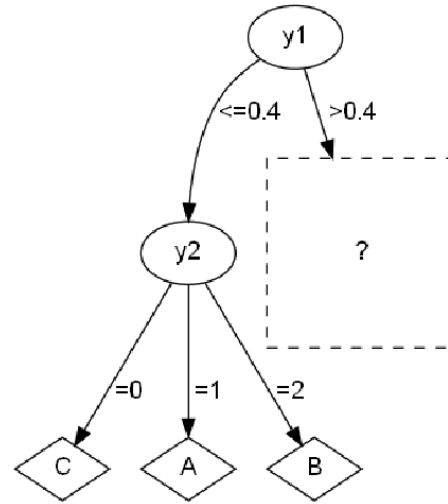
Aprendizagem 2023
Homework I – Group 28

Gonalo Barias (ist1103124) & Raquel Braunschweig (ist1102624)

Part I: Pen and Paper

Consider the partially learnt decision tree from the dataset D . D is described by four input variables – one numeric with values in $[0, 1]$ and 3 categorical – and a target variable with three classes.

| D | y_1 | y_2 | y_3 | y_4 | y_{out} |
|----------|-------|-------|-------|-------|-----------|
| x_1 | 0.24 | 1 | 1 | 0 | A |
| x_2 | 0.06 | 2 | 0 | 0 | B |
| x_3 | 0.04 | 0 | 0 | 0 | B |
| x_4 | 0.36 | 0 | 2 | 1 | C |
| x_5 | 0.32 | 0 | 0 | 2 | C |
| x_6 | 0.68 | 2 | 2 | 1 | A |
| x_7 | 0.9 | 0 | 1 | 2 | A |
| x_8 | 0.76 | 2 | 2 | 0 | A |
| x_9 | 0.46 | 1 | 1 | 1 | B |
| x_{10} | 0.62 | 0 | 0 | 1 | B |
| x_{11} | 0.44 | 1 | 2 | 2 | C |
| x_{12} | 0.52 | 0 | 2 | 0 | C |



1. Complete the given decision tree using **Information gain with Shannon entropy** (\log_2). Consider that: i) a minimum of 4 observations is required to split an internal node, and ii) decisions by ascending alphabetic order should be placed in case of ties.

The entropy of y_{out} is given by:

$$\begin{aligned}
 E(y_{out}|y_1 > 0.4) = & -P(y_{out} = A, y_1 > 0.4) \log_2 (P(y_{out} = A, y_1 > 0.4)) \\
 & -P(y_{out} = B, y_1 > 0.4) \log_2 (P(y_{out} = B, y_1 > 0.4)) \\
 & -P(y_{out} = C, y_1 > 0.4) \log_2 (P(y_{out} = C, y_1 > 0.4))
 \end{aligned} \tag{1}$$

We can calculate $E(y_{out})$:

$$E(y_{out}|y_1 > 0.4) = -\left(\frac{3}{7} \log_2 \left(\frac{3}{7}\right) + \frac{2}{7} \log_2 \left(\frac{2}{7}\right) + \frac{2}{7} \log_2 \left(\frac{2}{7}\right)\right) \approx 1.5567$$

The next step is calculating $E(y_{out}|y_1 > 0.4, y_x)$, in which x will take the values of 2, 3 or 4:

$$\begin{aligned}
 E(y_{out}|y_1 > 0.4, y_x) = & -P(y_x = 0)E(y_{out}|y_1 > 0.4, y_x = 0) \\
 & -P(y_x = 1)E(y_{out}|y_1 > 0.4, y_x = 1) \\
 & -P(y_x = 2)E(y_{out}|y_1 > 0.4, y_x = 2)
 \end{aligned} \tag{2}$$

And the information gain of variable y_x is given by

$$IG(y_x) = E(y_{out}|y_1 > 0.4) - E(y_{out}|y_1 > 0.4, y_x) \quad (3)$$

Let's start with $x = 2$:

$$P(y_2 = 0, y_1 > 0.4) = \frac{3}{7}$$

$$P(y_2 = 1, y_1 > 0.4) = \frac{2}{7}$$

$$P(y_2 = 2, y_1 > 0.4) = \frac{2}{7}$$

$$E(y_{out}|y_1 > 0.4, y_2 = 0) = -\left(\frac{1}{3} \log_2 \left(\frac{1}{3}\right) + \frac{1}{3} \log_2 \left(\frac{1}{3}\right) + \frac{1}{3} \log_2 \left(\frac{1}{3}\right)\right) \approx 1.58493$$

$$E(y_{out}|y_1 > 0.4, y_2 = 1) = -\left(\frac{0}{2} \log_2 \left(\frac{0}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) = 1$$

$$E(y_{out}|y_1 > 0.4, y_2 = 2) = -\left(\frac{2}{2} \log_2 \left(\frac{2}{2}\right) + \frac{0}{2} \log_2 \left(\frac{0}{2}\right) + \frac{0}{2} \log_2 \left(\frac{0}{2}\right)\right) = 0$$

Therefore, replacing these values on equation (2), gives us:

$$E(y_{out}|y_1 > 0.4, y_2) = \frac{3}{7} \times 1.58493 + \frac{2}{7} \times 1 + \frac{2}{7} \times 0 \approx 0.965.$$

Finally, we can calculate the information gain, as per (3),

$$IG(y_2) = 1.5567 - 0.965 = 0.5917$$

Now, let's calculate for $x = 3$:

$$P(y_3 = 0, y_1 > 0.4) = \frac{1}{7}$$

$$P(y_3 = 1, y_1 > 0.4) = \frac{2}{7}$$

$$P(y_3 = 2, y_1 > 0.4) = \frac{4}{7}$$

$$E(y_{out}|y_1 > 0.4, y_3 = 0) = -\left(\frac{0}{1} \log_2 \left(\frac{0}{1}\right) + \frac{1}{1} \log_2 \left(\frac{1}{1}\right) + \frac{0}{1} \log_2 \left(\frac{0}{1}\right)\right) = 0$$

$$E(y_{out}|y_1 > 0.4, y_3 = 1) = -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{0}{2} \log_2 \left(\frac{0}{2}\right)\right) = 1$$

$$E(y_{out}|y_1 > 0.4, y_3 = 2) = -\left(\frac{2}{4} \log_2 \left(\frac{2}{4}\right) + \frac{0}{4} \log_2 \left(\frac{0}{4}\right) + \frac{2}{4} \log_2 \left(\frac{2}{4}\right)\right) = 1$$

Therefore, replacing these values on equation (2), gives us:

$$E(y_{out}|y_1 > 0.4, y_3) = \frac{1}{7} \times 0 + \frac{2}{7} \times 1 + \frac{4}{7} \times 1 \approx 0.8571.$$

Finally, we can calculate the information gain, as per (3),

$$IG(y_3) = 1.5567 - 0.8571 = 0.6996$$

Finally, let's calculate for x = 4:

$$P(y_4 = 0, y_1 > 0.4) = \frac{2}{7}$$

$$P(y_4 = 1, y_1 > 0.4) = \frac{3}{7}$$

$$P(y_4 = 2, y_1 > 0.4) = \frac{2}{7}$$

$$E(y_{out}|y_1 > 0.4, y_4 = 0) = -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{0}{2} \log_2 \left(\frac{0}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{3}\right)\right) = 1$$

$$E(y_{out}|y_1 > 0.4, y_4 = 1) = -\left(\frac{1}{3} \log_2 \left(\frac{1}{3}\right) + \frac{2}{3} \log_2 \left(\frac{2}{3}\right) + \frac{0}{3} \log_2 \left(\frac{0}{3}\right)\right) \approx 0.9183$$

$$E(y_{out}|y_1 > 0.4, y_4 = 2) = -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{0}{2} \log_2 \left(\frac{0}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) = 1$$

Therefore, replacing these values on equation (2), gives us:

$$E(y_{out}|y_1 > 0.4, y_4) = \frac{2}{7} \times 1 + \frac{3}{7} \times 0.9183 + \frac{2}{7} \times 1 \approx 0.965.$$

Finally, we can calculate the information gain, as per (3),

$$IG(y_4) = 1.5567 - 0.965 = 0.5917$$

Upon computing the information gains for each attribute, it is evident that y_3 yields the highest value of 0.6996. Consequently, it is selected as the next node, resulting in the construction of the following decision tree:

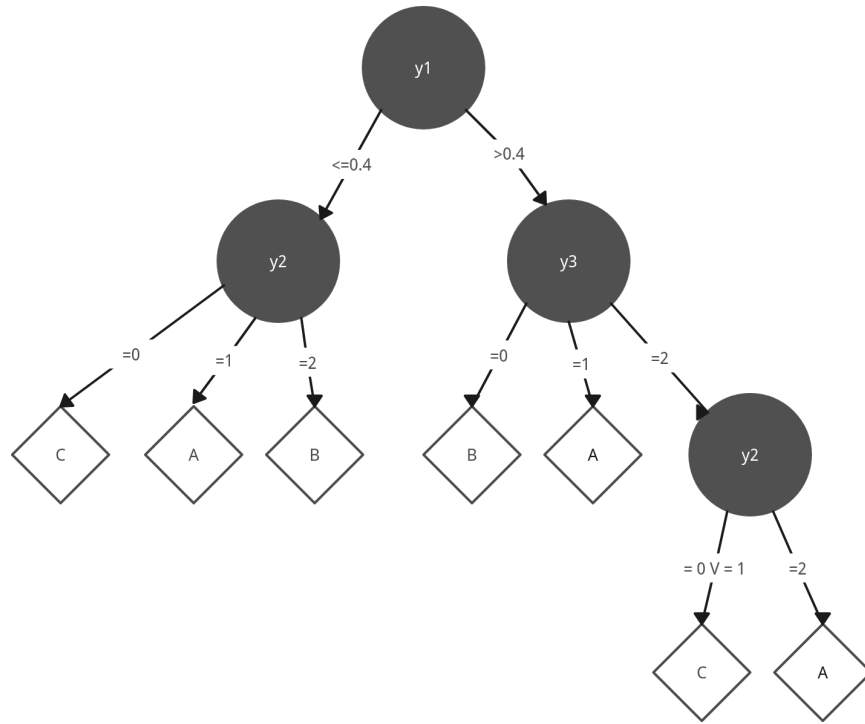


Figure 1: Decision Tree for exercise I.1

2. Draw the training confusion matrix for the learnt decision tree.

Following the learnt decision tree above, we can predict the values for each observation. For each observation, we look at the value for the first variable (y_1) and follow the branch that corresponds with its value. From the node we arrive at, we do the same thing for the next variable, and we keep doing this until we reach a leaf. The class present in this leaf will be the predicted value, while the real value is the value of y_{out} for that observation. Below we present the real values along with the predicted ones:

| | | | | | | | | | | | | | |
|-----------|---|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|----------|----------|
| | | x_1 | x_2 | x_3 | x_4 | x_5 | x_6 | x_7 | x_8 | x_9 | x_{10} | x_{11} | x_{12} |
| real | = | [A | B | B | C | C | A | A | A | B | B | C | C] |
| predicted | = | [A | B | C | C | C | A | A | A | A | B | C | C] |

Finally, we can show the count of each pair of real and predicted values in a confusion matrix (e.g. 4 pairs of AA from observations x_1, x_6, x_7 and x_8):

| | | Real | | | |
|-----------|---|------|---|---|----|
| | | A | B | C | |
| Predicted | A | 4 | 1 | 0 | 5 |
| | B | 0 | 2 | 0 | 2 |
| | C | 0 | 1 | 4 | 5 |
| | | 4 | 4 | 4 | 12 |

3. Identify which class has the lowest training F1 score.

$F1_{score}$ is given by the following equation:

$$F1_{score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

And precision and recall are given by:

$$\text{Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}} \quad (5)$$

$$\text{Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}} \quad (6)$$

Therefore, **let's start by calculating the precision for A, B and C** by replacing the values on (5):

$$\text{Precision}_A = \frac{4}{4+1} = \frac{4}{5}$$

$$\text{Precision}_B = \frac{2}{2+0} = 1$$

$$\text{Precision}_C = \frac{4}{4+1} = \frac{4}{5}$$

Now, it's time to calculate the recalls for A, B and C, using the equation on (6):

$$\text{Recall}_A = \frac{4}{4+0} = 1$$

$$\text{Recall}_B = \frac{2}{2+2} = \frac{1}{2}$$

$$\text{Recall}_C = \frac{4}{4+0} = 1$$

Finally, let's calculate the $F1_{score}$, using the equation (4):

$$F1_{score}A = 2 \cdot \frac{\frac{4}{5} \cdot 1}{\frac{4}{5} + 1} \approx 0.8889$$

$$F1_{score}B = 2 \cdot \frac{1 \cdot \frac{1}{2}}{1 + \frac{1}{2}} \approx 0.6667$$

$$F1_{score}C = 2 \cdot \frac{\frac{4}{5} \cdot 1}{\frac{4}{5} + 1} \approx 0.8889$$

The class with the lowest training F1 score is B, with a score of 0.6667.

4. **Considering y_2 to be ordinal, assess if y_1 and y_2 are correlated using the Spearman coefficient.**

To calculate the Spearman coefficient when there's rank, we have to use the following equation:

$$\text{Spearman}(y_x, y_y) = \text{Pearson}(y'_x, y'_y) = \frac{\text{cov}(y'_x, y'_y)}{\sigma_{y'_x} \sigma_{y'_y}} = \frac{\sum_{i=1}^n (y'_{xi} - \bar{y}'_x)(y'_{yi} - \bar{y}'_y)}{\sqrt{\sum_{i=1}^n (y'_{xi} - \bar{y}'_x)^2} \sqrt{\sum_{i=1}^n (y'_{yi} - \bar{y}'_y)^2}} \quad (7)$$

Firstly, **let's order y_1 and y_2** so we can calculate the ranks and y'_1 and y'_2 :

$$\text{ordered_}y_1 = [0.04, 0.06, 0.24, 0.32, 0.36, 0.44, 0.46, 0.52, 0.62, 0.68, 0.76, 0.9]$$

$$\text{ranks_}y_1 = [1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12]$$

$$y'_1 = [3, 2, 1, 5, 4, 10, 12, 11, 7, 9, 6, 8]$$

$$\text{ordered_}y_2 = [0, 0, 0, 0, 0, 0, 1, 1, 1, 2, 2, 2]$$

$$\text{ranks_}y_2 = [3.5, 3.5, 3.5, 3.5, 3.5, 3.5, 8, 8, 8, 11, 11, 11]$$

$$y'_2 = [8, 11, 3.5, 3.5, 3.5, 11, 3.5, 11, 8, 3.5, 8, 3.5]$$

Now, we have all we need to calculate **the Spearman coefficient** using the expression at (7). Here is the result:

$$\text{Spearman}(y_1, y_2) = \text{Pearson}(y'_1, y'_2) \approx 0.07966$$

Since the rank correlation (under Spearman coefficient) obtained a value relatively close to 0, we can conclude that the variables y_1 and y_2 are almost non correlated.

5. **Draw the class-conditional relative histograms of y_1 using 5 equally spaced bins in $[0, 1]$. Find the root split using the discriminant rules from these empirical distributions.**

The following histograms were created by observing the values of y_1 and fixing the class in y_{out} . Since there are three classes A , B and C , we get three histograms. In each histogram we have 5 equally spaced bins in $[0, 1]$, meaning that each bin has a length of 0.2. For each class, when a value of y_1 falls into a bin, we increase the value of that bin by one.

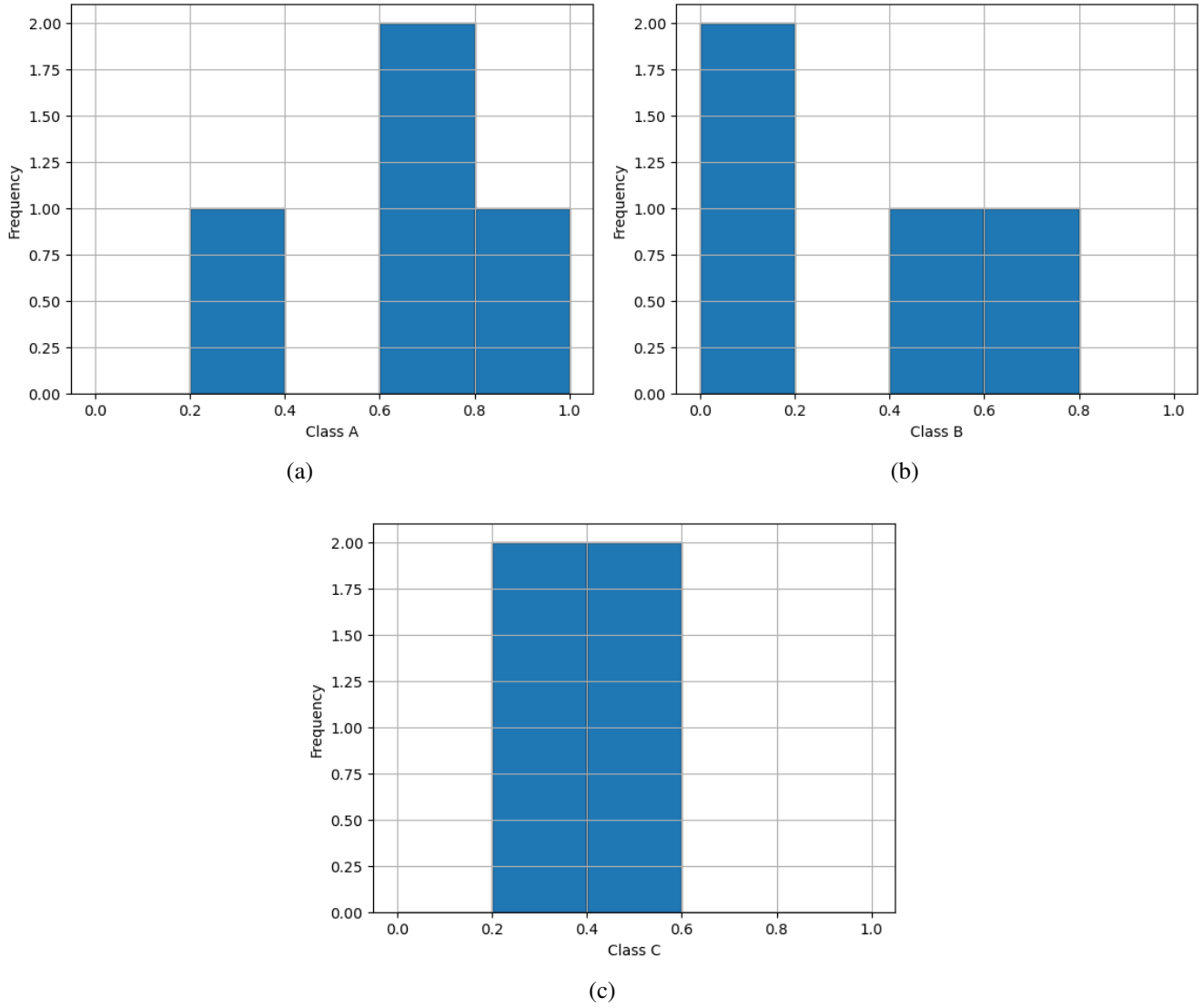


Figure 2: Class-conditional relative histograms of y_1 with 5 equally spaced bins in $[0, 1]$

In order to find the root split using the discriminant rules from the empirical distributions, we made a decision tree with only a root split on y_1 . Each branch will have the values of one bin, meaning there will be 5 branches (e.g. the first branch is for values of y_1 that belong to $[0, 0.2[$). The leaf for each branch corresponds to the class that has the highest value for that specific bin in its histogram.

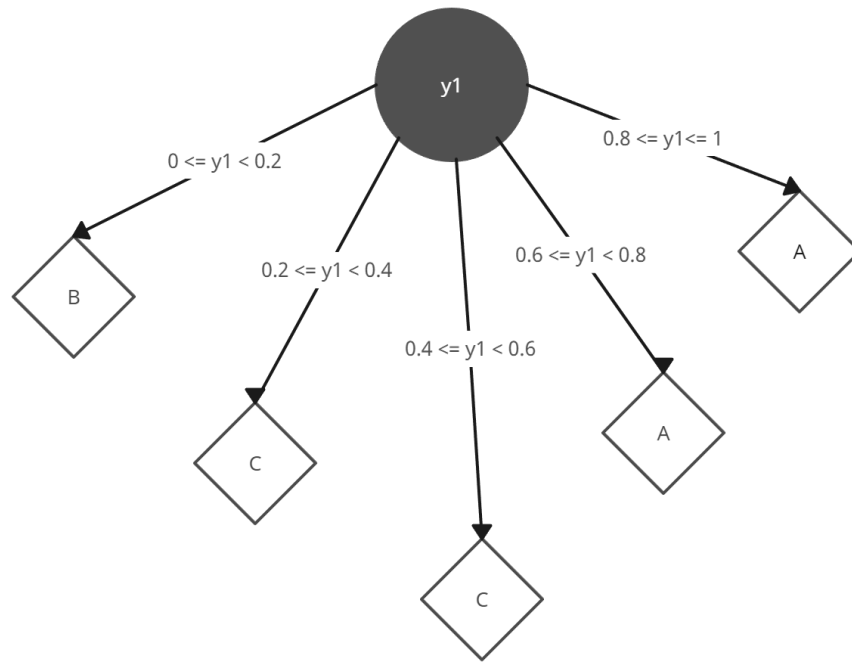


Figure 3: Decision Tree using the discriminant rules from the empirical distributions

Part II: Programming

Consider the `column_diagnosis.arff` data available at the homework tab, comprising 6 biomechanical features to classify 310 orthopaedic patients into 3 classes (normal, disk hernia, spondilolysthesis).

1. Apply `f_classif` from `sklearn` to assess the discriminative power of the input variables. Identify the input variable with the highest and lowest discriminative power. Plot the class-conditional probability density functions of these two input variables.

```

1 import numpy as np, pandas as pd, seaborn as sns, matplotlib.pyplot as plt
2 from scipy.io.arff import loadarff
3 from sklearn.feature_selection import f_classif
4
5 # Read the ARFF file and prepare data
6 data = loadarff("./data/column_diagnosis.arff")
7 df = pd.DataFrame(data[0])
8 df["class"] = df["class"].str.decode("utf-8")
9 X, y = df.drop("class", axis=1), df["class"]
10
11 # Apply f_classif
12 f_scores, _ = f_classif(X, y)
13
14 # Obtains the variables with the highest and lowest discriminative power.
15 h_disc_power_var = X.columns[np.argmax(f_scores)]
16 l_disc_power_var = X.columns[np.argmin(f_scores)]
17
18 plt.figure(figsize=(8, 6))
19

```



```

20 # Plot for the highest discriminative power variable
21 for class_label in np.unique(y):
22     class_data = X.loc[y == class_label, h_disc_power_var]
23     sns.kdeplot(
24         class_data,
25         label=f"Class {class_label} - {h_disc_power_var}",
26         linewidth=2,
27     )
28
29 # Plot for the lowest discriminative power variable
30 for class_label in np.unique(y):
31     class_data = X.loc[y == class_label, l_disc_power_var]
32     sns.kdeplot(
33         class_data,
34         label=f"Class {class_label} - {l_disc_power_var}",
35         linestyle="--",
36         linewidth=2,
37     )
38
39 plt.xlabel("Variables")
40 plt.ylabel("Density")
41
42 plt.legend()
43 plt.grid(True)
44 plt.show()

```

As you can see in the graph ahead, the highest discriminative power variable is *degree_spondilolisthesis* and the lowest discriminative power variable is *pelvic_radius*.

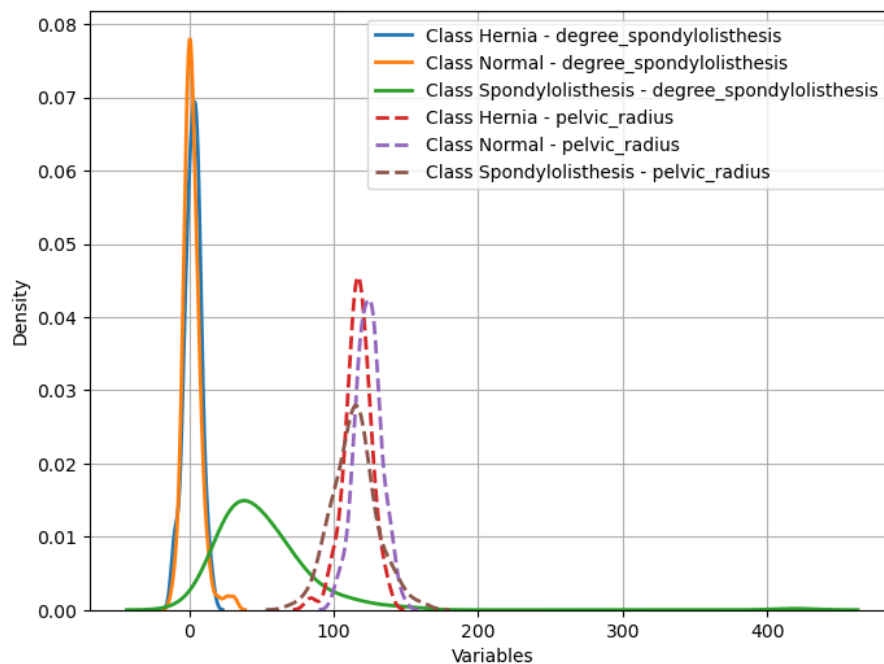


Figure 4: Class-conditional probability density functions of the highest and lowest discriminative power variables.

2. Using a stratified 70-30 training-testing split with a fixed seed (`random_state=0`), assess in

a single plot both the training and testing accuracies of a decision tree with depth limits in {1, 2, 3, 4, 5, 6, 8, 10} and the remaining parameters as default.

[Optional] Note that split thresholding of numeric variables in decision trees is non-deterministic in sklearn, hence you may opt to average the results using 10 runs per parameterization.

```
1 import pandas as pd, matplotlib.pyplot as plt
2 from scipy.io.arff import loadarff
3 from sklearn import metrics, tree
4 from sklearn.model_selection import train_test_split
5
6 # Read the ARFF file and prepare data
7 data = loadarff("./data/column_diagnosis.arff")
8 df = pd.DataFrame(data[0])
9 df["class"] = df["class"].str.decode("utf-8")
10 X, y = df.drop("class", axis=1), df["class"]
11
12 DEPTH_LIMIT = [1, 2, 3, 4, 5, 6, 8, 10]
13 training_accuracy, test_accuracy = [], []
14
15 # Split the dataset into a testing set (30%) and a training set (70%)
16 X_train, X_test, y_train, y_test = train_test_split(
17     X, y, test_size=0.3, stratify=y, random_state=0
18 )
19
20 for depth_limit in DEPTH_LIMIT:
21     # Create and fit the decision tree classifier
22     predictor = tree.DecisionTreeClassifier(
23         max_depth=depth_limit, random_state=0
24     )
25     predictor.fit(X_train, y_train)
26
27     # Use the decision tree to predict the outcome of the given observations
28     y_train_pred = predictor.predict(X_train)
29     y_test_pred = predictor.predict(X_test)
30
31     # Get the accuracy of each test
32     train_acc = metrics.accuracy_score(y_train, y_train_pred)
33     training_accuracy.append(train_acc)
34     test_acc = metrics.accuracy_score(y_test, y_test_pred)
35     test_accuracy.append(test_acc)
36
37 plt.plot(
38     DEPTH_LIMIT, training_accuracy,
39     label="Training Accuracy", marker="+", color="#f8766d"
40 )
41 plt.plot(
42     DEPTH_LIMIT, test_accuracy,
43     label="Test Accuracy", marker=".", color="#00bfc4"
44 )
45
46 plt.xlabel("Depth Limit")
47 plt.ylabel("Accuracy")
48
49 plt.legend()
50 plt.grid(True)
```

51 `plt.show()`

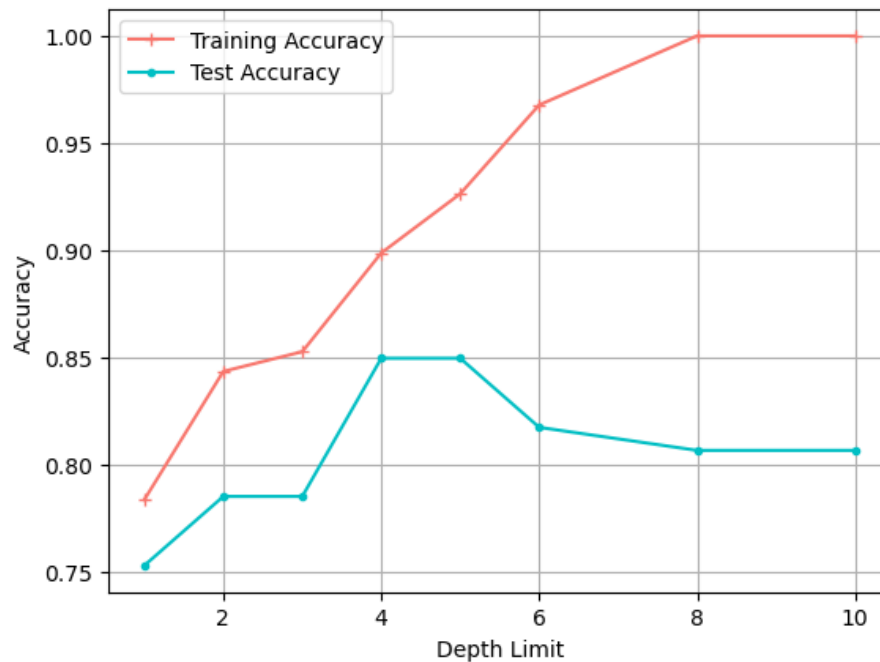


Figure 5: Accuracy of the trained decision tree, applied to both a test and training sets, for varying depth limits.

3. Comment on the results, including the generalization capacity across settings.

The graphic shows that as the depth limit of the decision tree increases, the training accuracy rises steadily until reaching 100%. This observation suggests that deeper trees can better fit the training data, capturing complex patterns and achieving higher accuracy when evaluated on the dataset.

The testing accuracy initially improves as the depth limit increases, indicating improved generalization. However, beyond a certain depth limit (around 4 or 5 in this case), the testing accuracy starts to decline. Therefore, we conclude there is a loss in generalization capacity, which causes overfitting, implying that, in this case, overly complex decision trees can fit noise in the training data, performing poorly on new, unseen data.

The optimal depth limit appears to be around 4 or 5, which seems to be the point where the accuracy is maximized for the testing data, striking a balance between model complexity and generalization to new data, this way we can avoid both underfitting (too simple) and overfitting (too complex).

4. To deploy the predictor, a healthcare team opted to learn a single decision tree (`random.state=0`) using *all* available data as training data, and further ensuring that each leaf has a minimum of 20 individuals in order to avoid overfitting risks.

(a) Plot the decision tree.

```
1 import matplotlib.pyplot as plt, pandas as pd, numpy as np
2 from scipy.io.arff import loadarff
3 from sklearn.tree import DecisionTreeClassifier, plot_tree
```

```

4
5 # Read the ARFF file and prepare data
6 data = loadarff("./data/column_diagnosis.arff")
7 df = pd.DataFrame(data[0])
8 df["class"] = df["class"].str.decode("utf-8")
9 X, y = df.drop("class", axis=1), df["class"]
10
11 # Create and train the decision tree classifier
12 clf = DecisionTreeClassifier(random_state=0, min_samples_leaf=20)
13 clf.fit(X, y)
14
15 # Set style and plot the decision tree
16 plt.figure(figsize=(15, 10))
17 plot_tree(clf, filled=True, feature_names=list(X.columns),
18           class_names=list(np.unique(y)), rounded=True, fontsize=12)
19 plt.show()

```

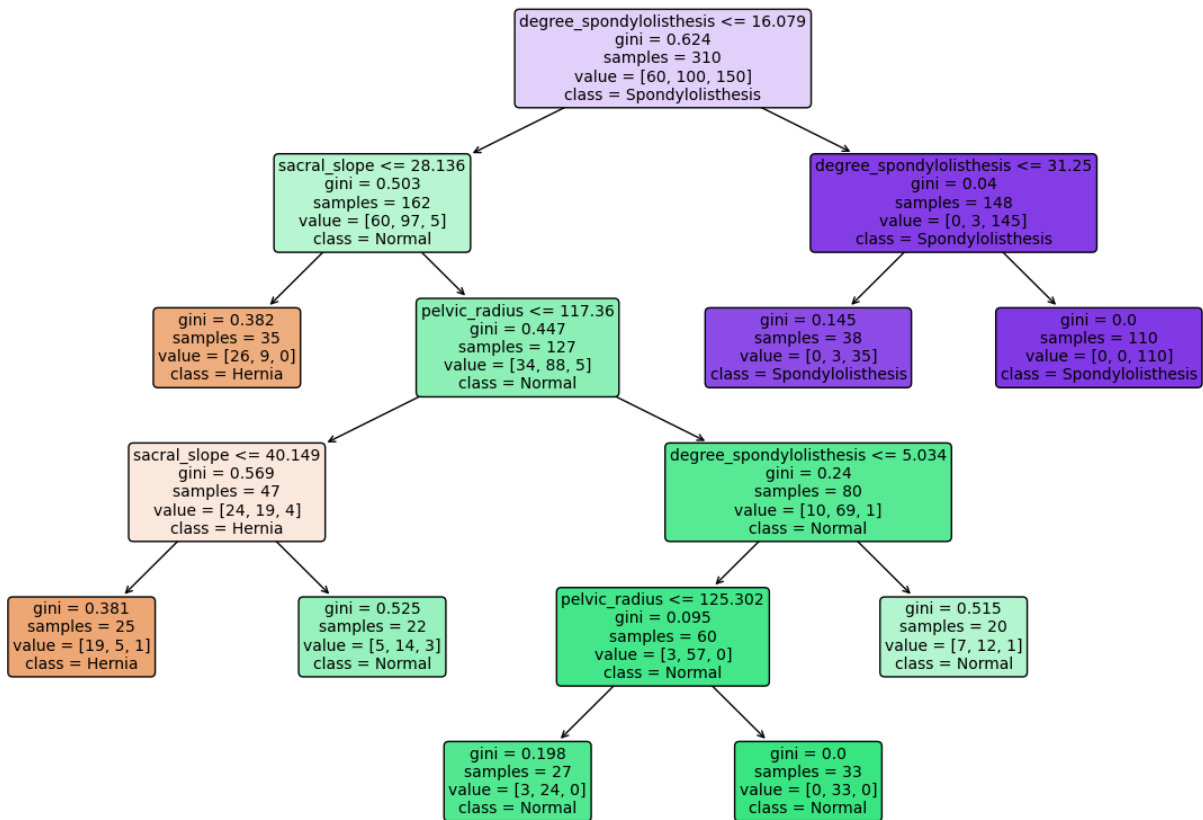


Figure 6: Decision Tree for exercise II.4a

(b) **Characterize a hernia condition by identifying the hernia-conditional associations.**

The hernia condition can be characterized by:

- i. $\text{degree_spondylolisthesis} \leq 16.079$ and $\text{sacral_slope} \leq 28.136$
- ii. $\text{degree_spondylolisthesis} \leq 16.079$, $\text{sacral_slope} > 28.136$, $\text{pelvic_radius} \leq 117.36$ and $\text{sacral_slope} \leq 40.149$

END