

CIDM 6355 Data Mining Methods HW1

(100 points in total; Due 11:59 PM Central Time, September 15, 2024)

Requirements: Follow the instruction, take the required screenshots with date and time (see the examples in RapidMiner Lab instruction), and answer all the questions. Sharing your queries, screenshots, or answers with other students is considered as cheating, which will be reported to the university authority. A screenshot without showing reliable date and time will receive a penalty of 50% of points. If identical screenshots are found from two or more students, such a misconduct will be reported to the university authority. Please type your name as below to indicate that you understand and comply with all the requirements in this homework.

Name: Sera Hill

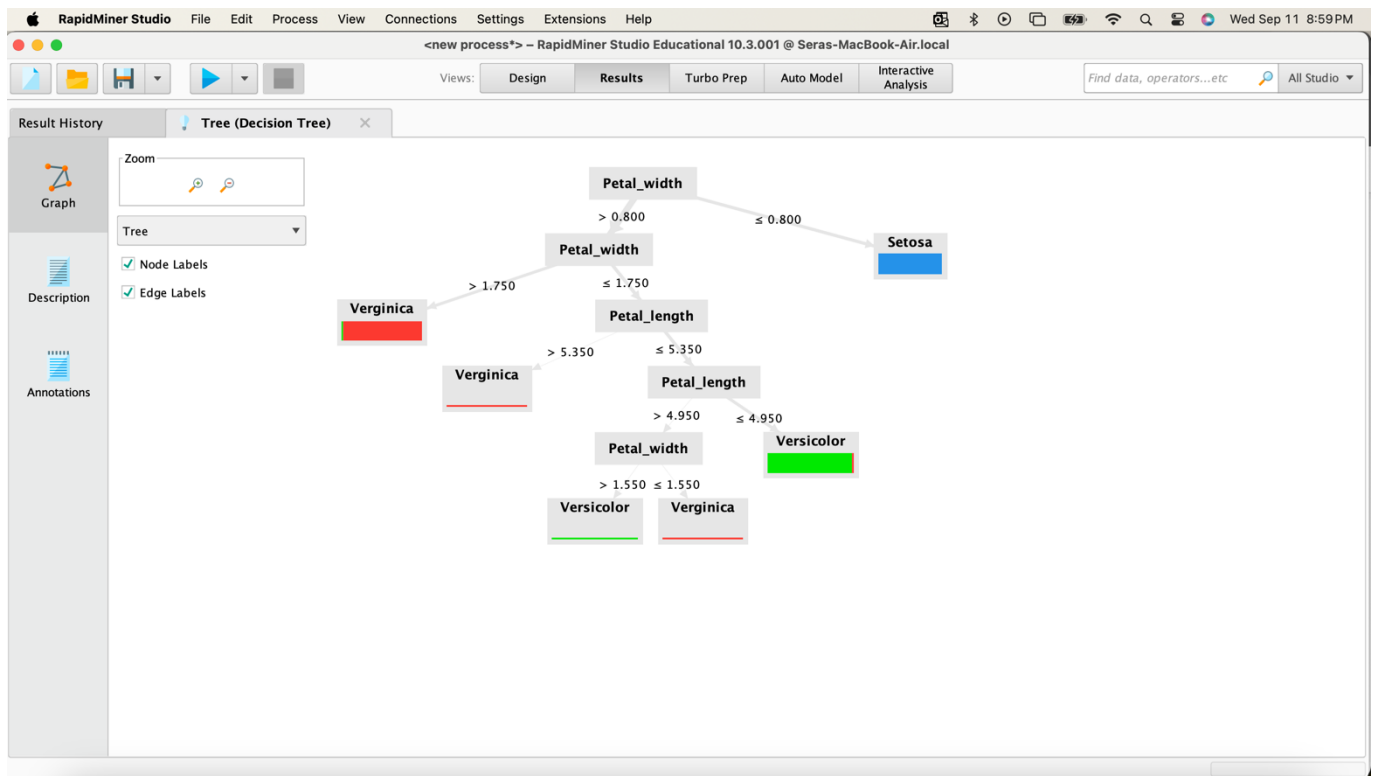
Part 1: Answer all the questions in Week 4 RapidMiner Lab (Step 1.3, 1.8.1-1.8.5, and 2.2.3 – 2.2.4) and an additional question via HW1-Part 1 Submission (40 points) on WTCLASS. You have two attempts and the higher one will be counted into your grade. Please DO NOT include them here; otherwise, they won't be graded here.

Part 2 Lab Screenshots and Deliverables (60 points)

Take the required screenshots with date and time and answer all the questions. **Windows and MacBook show the date and time differently, so your screenshot is acceptable as long as it displays the date and time, no matter how. MacBook displays the date and time on the top right corner.** If you do not know how to take a screenshot, please check this website <https://www.take-a-screenshot.org/> for more instructions. If you do not know how to show the date and time on your MAC Book, Google your question or try [this site](#). Sharing your queries, screenshots, or answers with other students is considered as cheating, which will be reported to the university authority.

1) Screenshots in RapidMiner Lab (10 points)

- Screenshot 1: A screenshot of the decision tree graph with date and time at Step 1.8 (5 points)



- Screenshot 2: A screenshot of prediction results for the 19 observations with date and time in Step 2.2 (5 points)

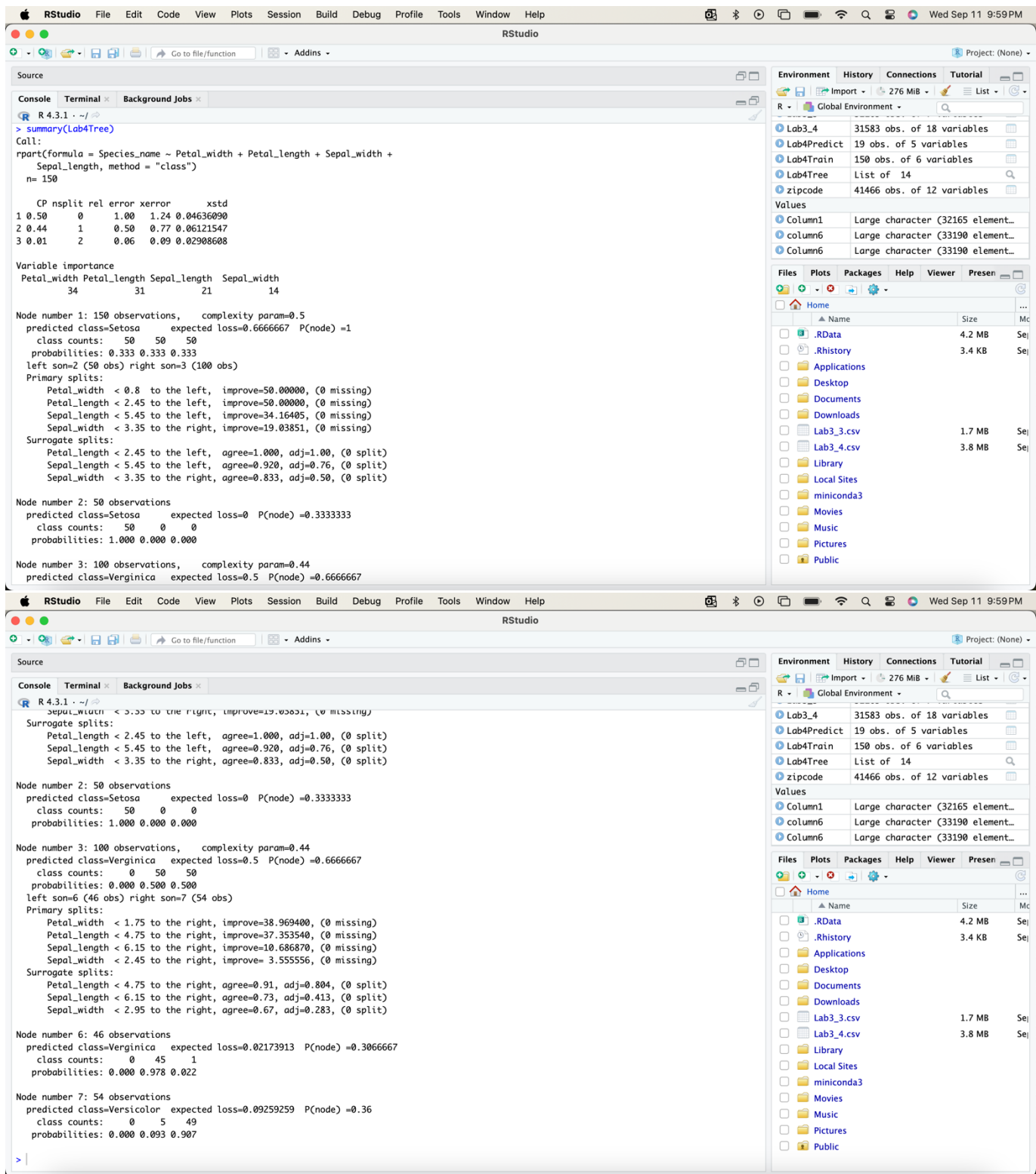
The screenshot displays the 'ExampleSet (Apply Model)' view in RapidMiner Studio, showing the prediction results for 19 observations. The table includes the following columns: Row No., prediction..., confidence..., confidence..., confidence..., Petal_width, Petal_length, Sepal_width, and Sepal_length.

Row No.	prediction...	confidence...	confidence...	confidence...	Petal_width	Petal_length	Sepal_width	Sepal_length
1	Verginica	0	0.022	0.978	2.700	3.200	3.600	6.300
2	Versicolor	0	0.979	0.021	1.400	4.800	4.400	5.900
3	Verginica	0	0.022	0.978	2.400	3.100	3.500	7.100
4	Versicolor	0	0.979	0.021	1.600	3.700	3.700	5.500
5	Verginica	0	0.022	0.978	1.900	1.500	3.600	6.300
6	Setosa	1	0	0	0.200	1.300	2.700	6.200
7	Setosa	1	0	0	0.300	4.500	4.200	6.700
8	Verginica	0	0.022	0.978	2.500	3.200	3.600	6.900
9	Verginica	0	0.022	0.978	2.900	7.700	4.400	5.500
10	Versicolor	0	0.979	0.021	1.200	3.500	4.500	4.500
11	Setosa	1	0	0	0.800	1.700	3.300	5.100
12	Verginica	0	0.022	0.978	2	1.200	4.200	5.900
13	Versicolor	0	0.979	0.021	1.700	1.500	3.100	6.300
14	Setosa	1	0	0	0.100	7	3.500	7
15	Setosa	1	0	0	0.600	5.100	3.900	6.300
16	Verginica	0	0.022	0.978	2.700	3.100	2.600	6.800
17	Verginica	0	0.022	0.978	2.600	2.900	4.200	6.500
18	Verginica	0	0.022	0.978	2	6.600	4.400	7.200

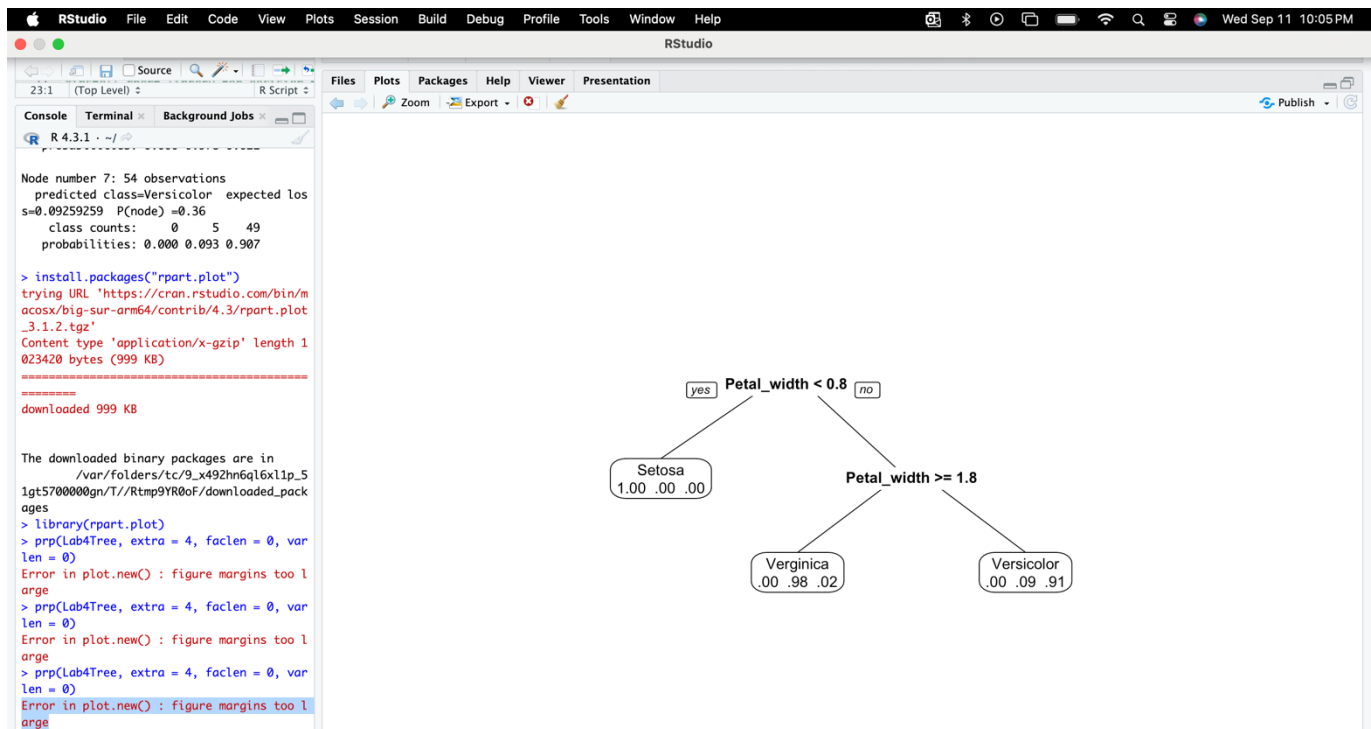
ExampleSet (19 examples, 4 special attributes, 4 regular attributes)

2) Deliverables in R Lab (50 points)

- **Deliverable R1:** take a screenshot of your decision tree model with date and time (5 points).

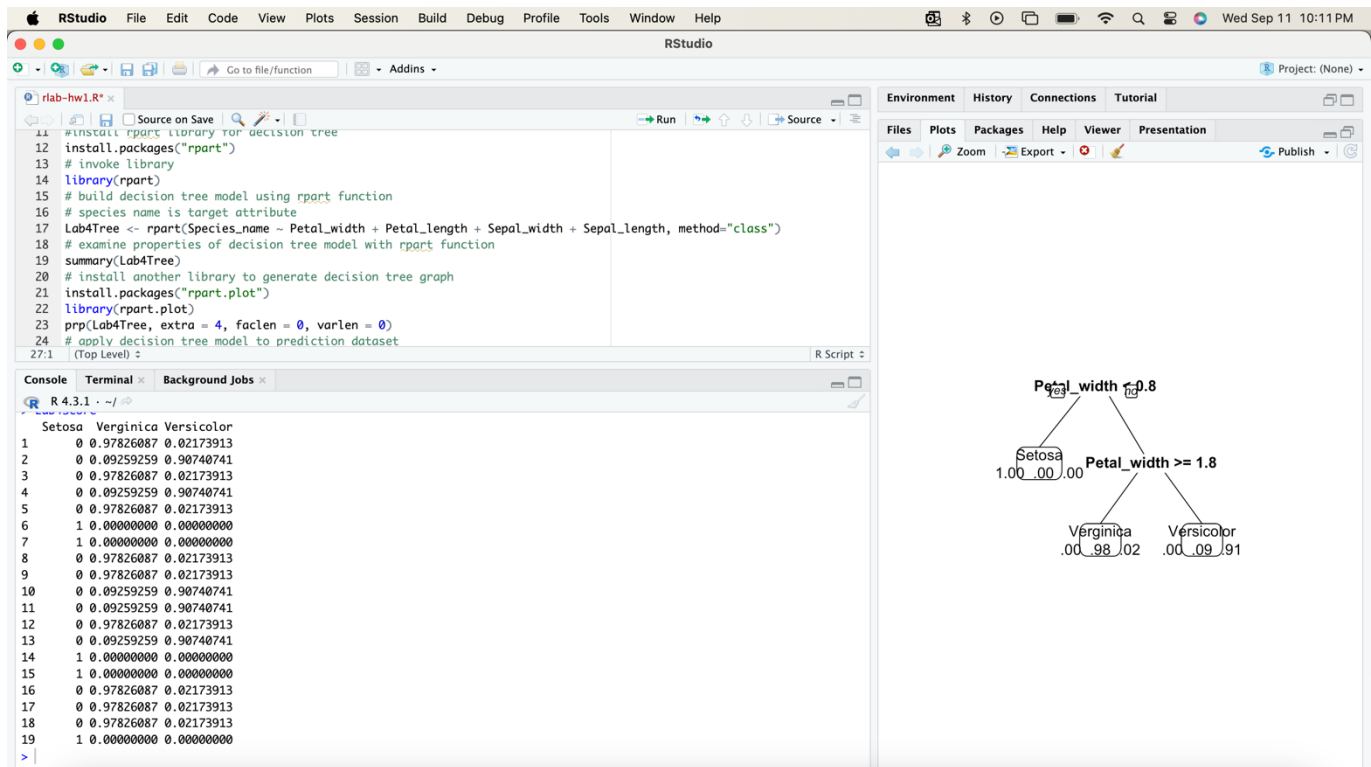


- Deliverable R2:** take a screenshot of your decision tree graph with date and time and briefly describe it. Your description must include the root node, split nodes, and leaf nodes. (10 points: 5 points for your screenshot and 5 points for your description).



The root node is the petal width < 0.8. There is only one split node: Petal_width >= 1.8. There are 3 leaf nodes that each have the species of flower: Setosa, Verginica, and Versicolor.

- **Deliverable R3:** after you apply the decision tree model to your prediction dataset, take a screenshot of the prediction result with date and time and briefly describe how the result help you determine the predicted class of each case. (10 points: 5 points for your screenshot and 5 points for your description).



These results help us predict what species the flowers will be. The closer the number is to 1, the greater the chance of that prediction being correct. These predictions follow what the decision tree is telling us, because there is no chance of the flower being anything besides Setosa if the petal width is less than .8. And the chances of the flower being Verginica if the petal width is greater or equal to 1.8 is much higher than Versicolor.

- **Deliverable R4:** take a screenshot of your decision tree model with date and time. Try to use the resources provided to understand its output (5 points).

```

library(party-hw1.R*)
6 dim(Lab4Train)
7 # view summary stats of training data
8 summary(Lab4Train)
9 # attach training dataset for ease of writing and maintaing code
10 attach(Lab4Train)
11 # install library party for decision tree
12 install.packages("party")
13 # invoke library
14 library("party")
15 #convert species name
16 Lab4Train[,6] <- as.factor(Lab4Train[,6])
17 Lab4Train$Species_name <- as.factor(Lab4Train$Species_name)
18 # build decision tree using ctree function
19 Lab4Tree2 <- ctree(Species_name ~ Petal_width + Sepal_width + Sepal_length, data = Lab4Train)
20 #examine properties of decision tree
21 Lab4Tree2
21:1 (Top Level)

```

Conditional inference tree with 4 terminal nodes

Response: Species_name
Inputs: Petal_width, Petal_length, Sepal_width, Sepal_length
Number of observations: 150

```

1) Petal_length <= 1.9; criterion = 1, statistic = 140.264
2)* weights = 50
1) Petal_length > 1.9
3) Petal_width <= 1.7; criterion = 1, statistic = 67.894
4) Petal_length <= 4.8; criterion = 0.999, statistic = 13.865
5)* weights = 46
4) Petal_length > 4.8
6)* weights = 8
3) Petal_width > 1.7
7)* weights = 46

```

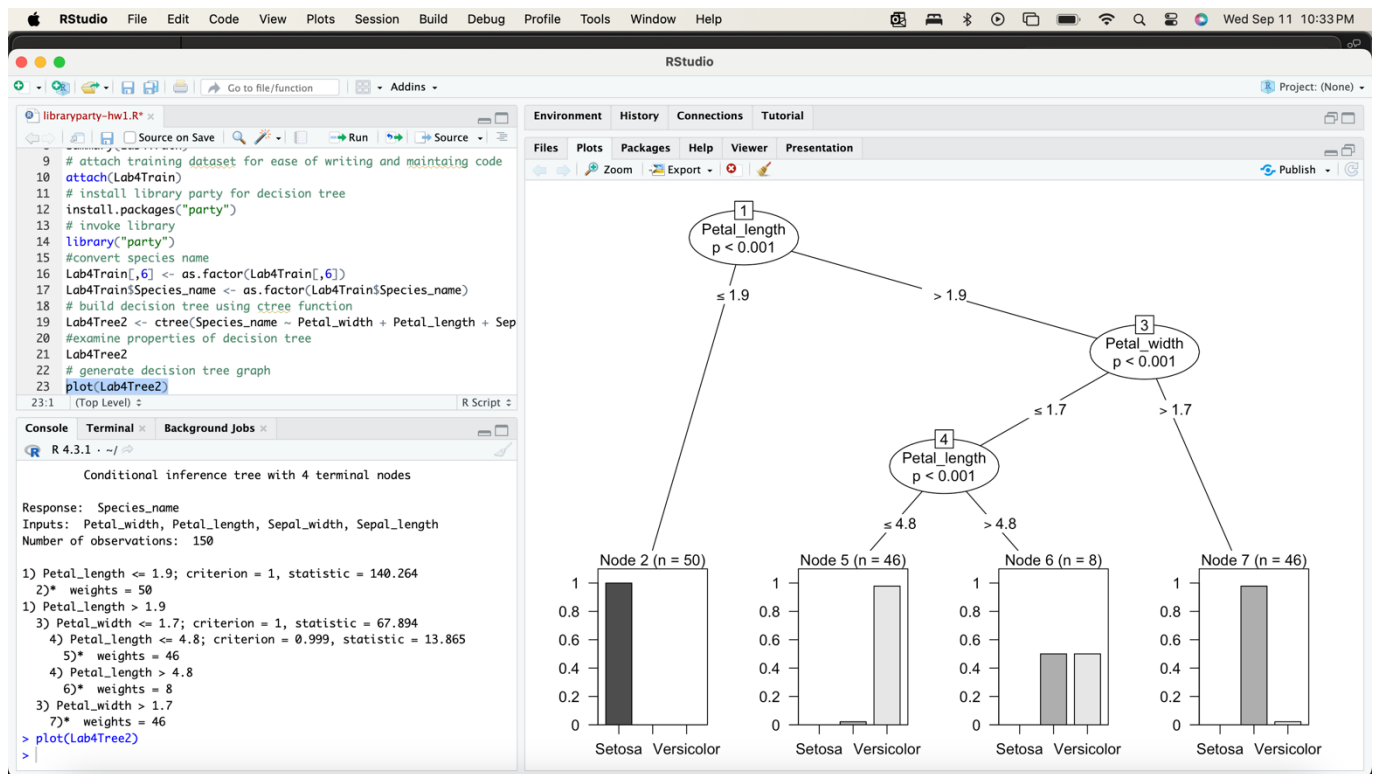
Environment

Object	Class	Size
Lab3_4	data.frame	31583 obs. of 18 variables
Lab4Predict	data.frame	19 obs. of 5 variables
Lab4Train	data.frame	150 obs. of 6 variables
Lab4Tree2	Formal class 'BinaryTree'	
zipcode	data.frame	41466 obs. of 12 variables
Column1	Large character	(32165 elements, 2.1 MB)
column6	Large character	(33190 elements, 1.3 MB)
Column6	Large character	(33190 elements, 1.3 MB)

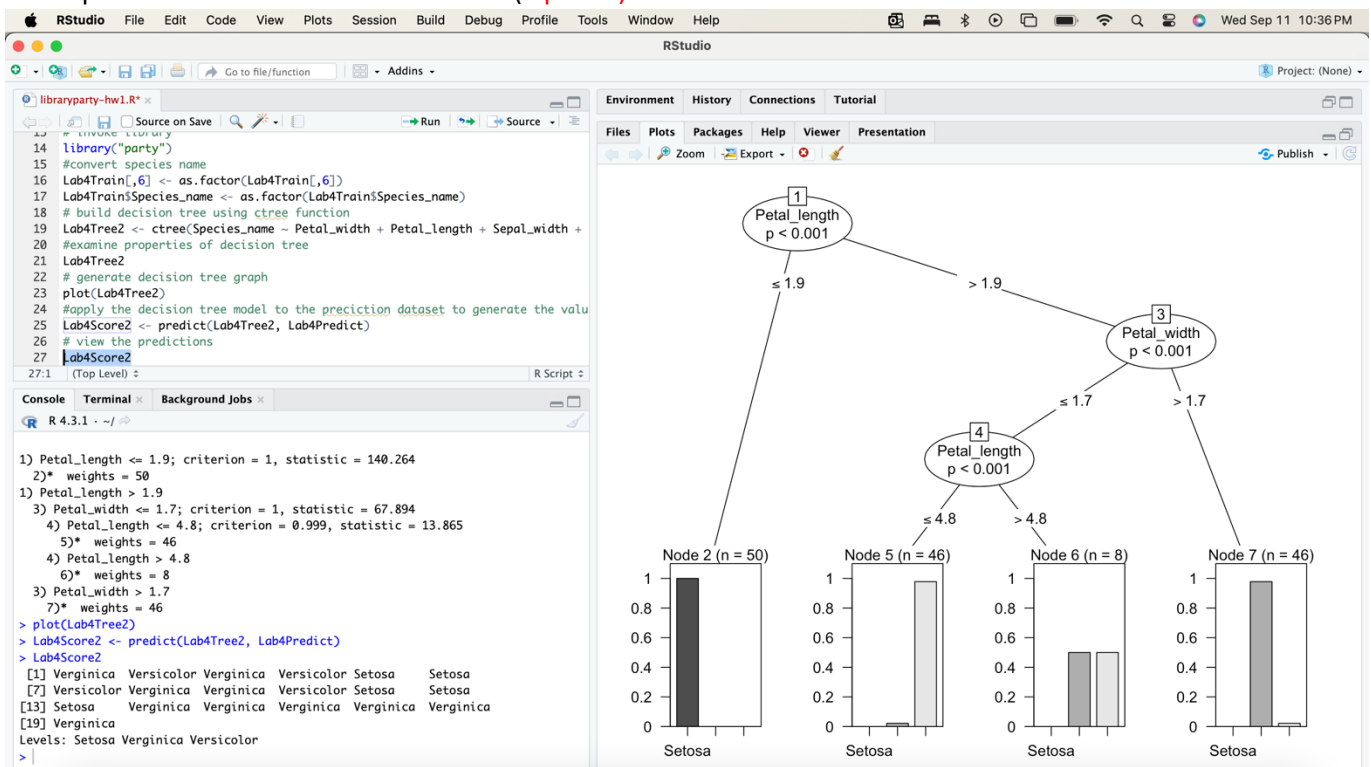
Files

Name	Size	Modified
.RData	4.2 MB	Sep 8, 2024, 9:42 AM
.Rhistory	4 KB	Sep 11, 2024, 10:20 PM
Applications		
Desktop		
Documents		
Downloads		
Lab3_3.csv	1.7 MB	Sep 8, 2024, 9:30 AM
Lab3_4.csv	3.8 MB	Sep 8, 2024, 9:42 AM
Library		
Local Sites		
miniconda3		
Movies		
Music		
Pictures		
Public		

- **Deliverable R5:** take a screenshot of your decision tree graph with date and time (5 points).



- **Deliverable R6:** after you apply the decision tree model to your prediction dataset, and take a screenshot of the prediction result with date and time (5 points).



- **Deliverable R7:** Choose one of the two decision tree models generated in R and compare it with the decision tree model generated in RapidMiner. Identify and discuss at least three differences between the two models. When discussing each difference, please include both R and RM. For example, "R does ..., but RM does not" (10 points).

RM ended up having more branches compared to the second R decision tree, the R decision tree model looks like it is simplified. R starts off with the petal length whereas RM starts off with petal width, and the value to determine which direction you go on the tree for the first petal length value is much smaller in R than the one for RM. The output in R combines the final leaf nodes into one graph that compares the different values, whereas in RM, the leaf nodes are almost all completely split out.