Explore the Iris Data Set Using R Studio

Didem Bulut Aykurt

MIS500-1 – Foundations of Data Analytics

Colorado State University-Global Campus

Dr. Steve Chung

July 24, 2022

When do I start to learn new knowledge that comes up the question of why I need that tool? Where will I use it? What kind of problem would I solve? Why do most companies work with the to R program? First, I would like to give some minor information about all my questions. I think the best reason to work R has over 12K packages and libraries available and handling big unstructured data, all free or lower than another program. Also, safe and security are the reasons all companies need them. Let's look at my work; R has an analysis tool to create graph&chart with big and heavy data. That has an excellent console to manage code, data, and results and display the graph and chart. Therefore, let's deep dive into the iris dataset with the R programing language.

Load and display the first six rows.

Module 1: Option #1: Critical Thinking

When displaying a statistical summary help to understand statistical information in a dataset with statistical terms; Iris data has five different variables, each of which has a close result mean and median number, which means most data point are relative to the mean, and most variable look a normal distribution thus we should check bar graph or boxplot to make sure each distribution to tell normal or not.



Sepal.Length and Petal. Lenght is the Iris dataset of a variable as sepal length mean bigger than petal length, which implies sepal length longer than petal length. Look at the petal length range more considerable than the sepal length; thus, the petal length data points are far, and the sepal length data points close to each other.

```
9
10  #store sepal length and petal length in a variable
11  sepal_length <- iris$Sepal.Length
12  petal_legth <- iris$Petal.Length
13  #calculate mean of sepal length and petal length variable
14  mean(sepal_length)
15  mean(petal_legth)
16  #calculate median of sepal length and petal length
17  median(sepal_length)
18  median(petal_legth)
19  #calculate range of sepal length and petal length
20  range(sepal_length)
21  range(petal_legth)
22  |
```

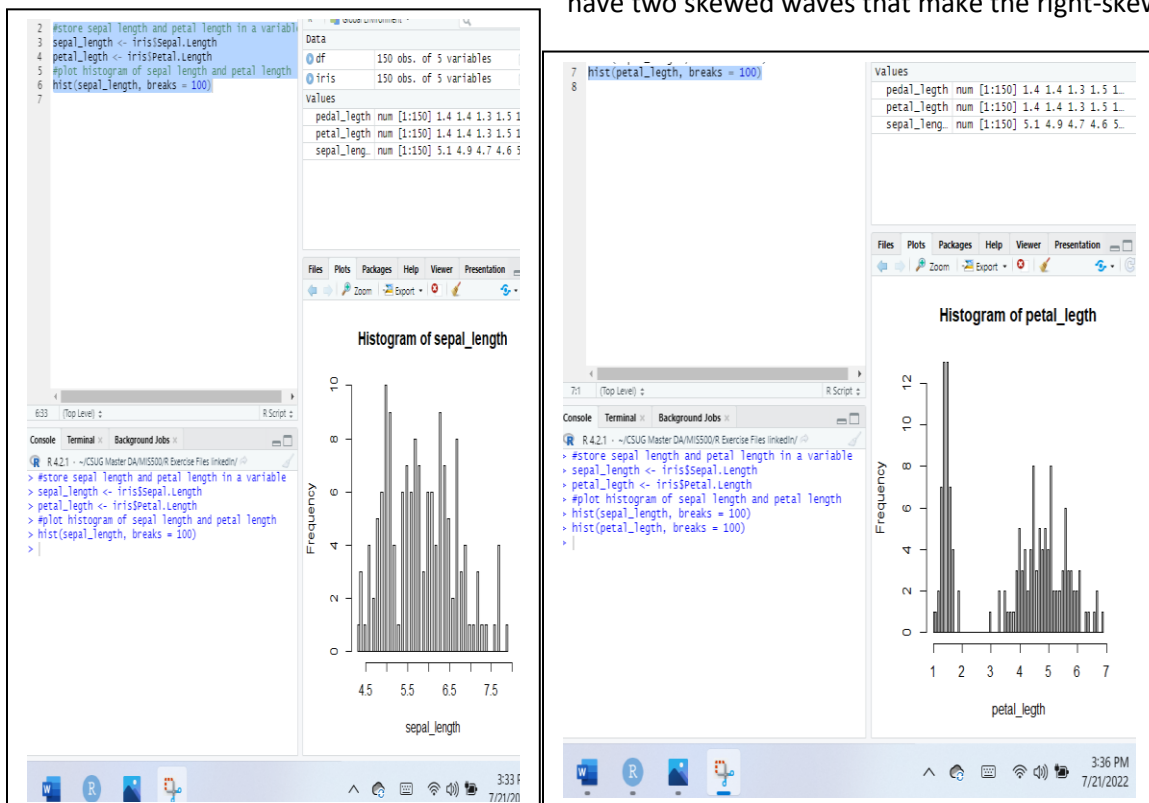22:1    ☰ LOAD AND PREPARE DATA ⬦                                    R Script ⬦

Console   Terminal ×   Background Jobs ×

ℝ  R 4.2.1 · ~/CSUG Master DA/MIS500/R Exercise Files linkedIn/

```
> #store sepal length and petal length in a variable
> sepal_length <- iris$Sepal.Length
> petal_legth <- iris$Petal.Length
> #calculate mean of sepal length and petal length variable
> mean(sepal_length)
[1] 5.843333
> mean(petal_legth)
[1] 3.758
> #calculate median of sepal length and petal length
> median(sepal_length)
[1] 5.8
> median(petal_legth)
[1] 4.35
> #calculate range of sepal length and petal length
> range(sepal_length)
[1] 4.3 7.9
> range(petal_legth)
[1] 1.0 6.9
> |
```
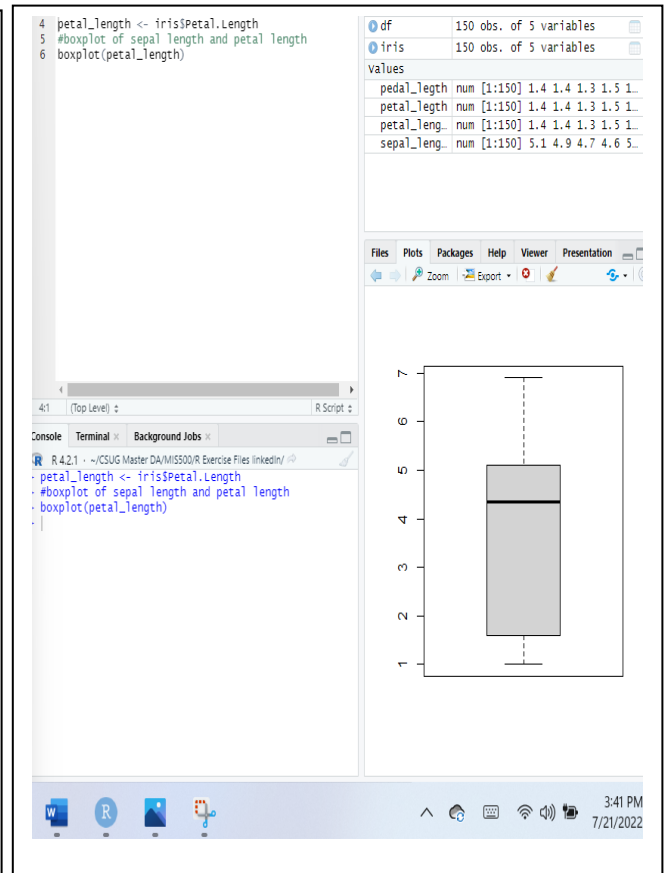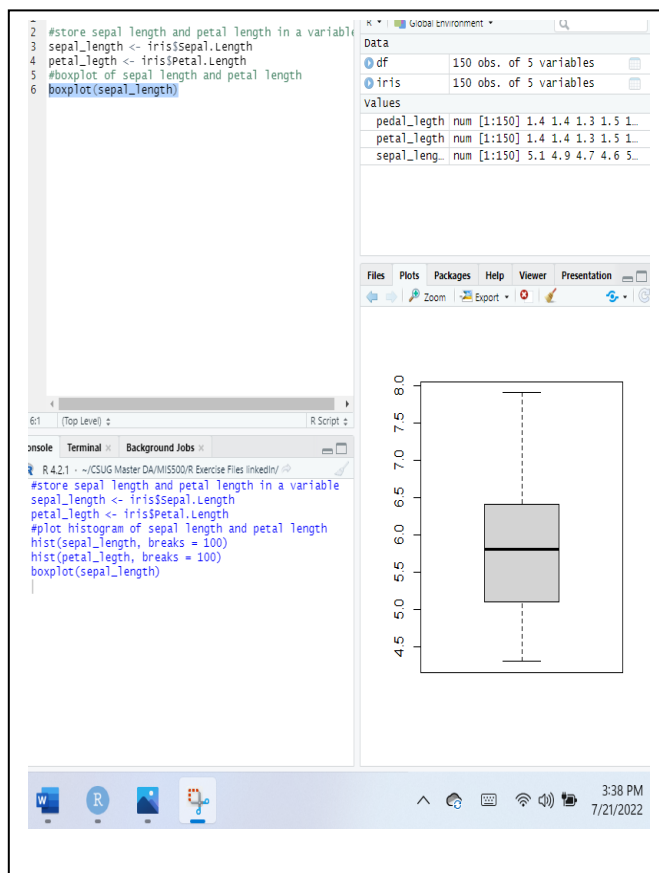
3:30 PM
7/21/2022

The sepal length bar graph has one pick, which means the normal distribution is symmetric from a peak of the curve and mainly observed data near the mean, but petal length variables far from the norm also have two skewed waves that make the right-skewed.



Histogram of sepal_length
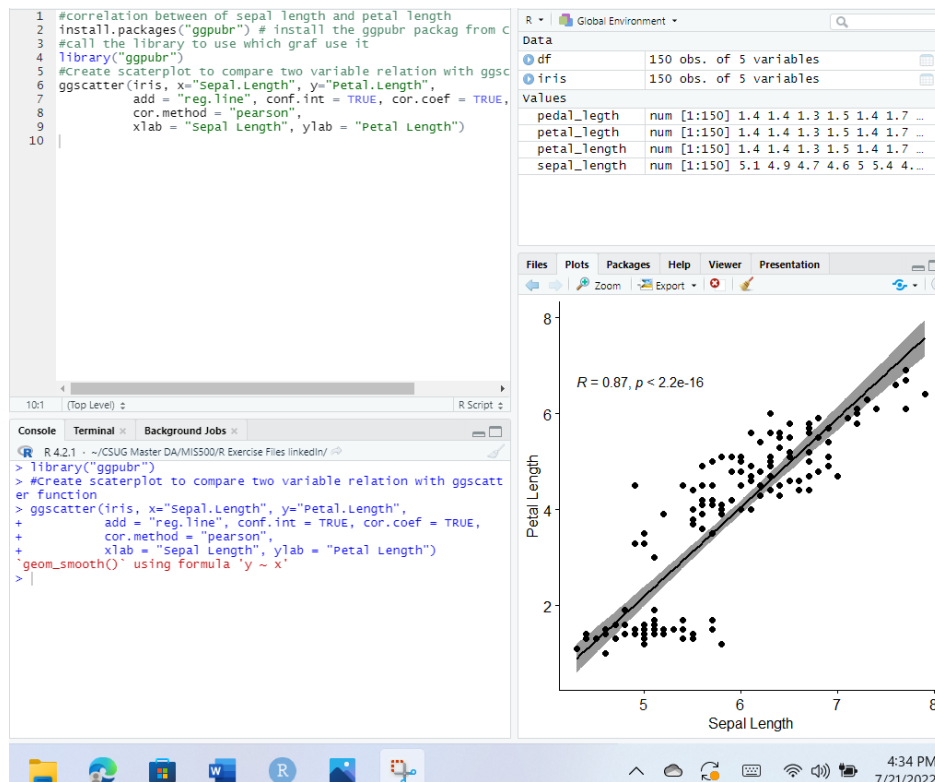


Histogram of petal_legth

Module 1: Option #1: Critical Thinking

The sepal length's box plot shows the symmetric distribution. Petal length's box plot tells positively skewed distribution—sepal length's box plot shorter the fewer spread data. Petal length box plot longer the more spread data.



Sepal length and petal length have positive relation, and R square is 87%, which means the regression model fits the data. Also, a p-value lower than .05 says to reject the null hypothesis. The null hypothesis implies that there is no relation between the two variables.

## Module 1: Option #1: Critical Thinking



## Conclusion

R is a new tool to add my skills. That is like python syntax. Where I got stuck starting the find R syntax. I am still looking cut sheet for syntax or any sources. Other statical terms make me so excited to find the meaning.

## Reference

[Correlation Test Between Two Variables in R - Easy Guides - Wiki - STHDA](#)

[Summary or Descriptive statistics in R - DataScience Made Simple](#)