

STATS 2107

Statistical Modelling and Inference II

Solutions

Workshop 6: The `rstatix` package

Matt Ryan

Semester 2 2022

Contents

The <code>rstatix</code> package	2
What is it	2
Installing the package	2
One more package	2
Some data to play with	2
What is the <code>iris</code> dataset	2
First taste of <code>rstatix</code>	3
1 sample t-test	3
What is wrong with base R	3
What is wrong with base R	3
Getting things from base R	3
<code>rstatix</code> brings dataframes!	4
Your turn	4
What to do	4
Two sample t-tests	7
Let's narrow our scope	7
The hypothesis to test	7
The how-to in <code>rstatix</code>	8
Conclusion?	8
Let's visualise this! (the code)	8
Let's visualise this! (the plot)	9

Your turn	9
What to do	9

The `rstatix` package

What is it

The `rstatix` package provides a pipe-friendly (`%>%` this thing) interface for performing hypothesis test, that fits nicely into the `tidyverse` framework and philosophy.

This package allows for a single framework to do many different hypothesis tests, a full list of which can be found in [here](#)

Installing the package

This is available on CRAN, so you can install with

```
install.packages("rstatix") # Install the package
library(rstatix) # Load the package
```

One more package

```
install.packages("ggpubr") # Install the package
library(ggpubr) # Load the package
```

Some data to play with

For todays workshop, we will use the classic `iris` dataset by Fisher. Load this dataset and convert it into a tibble with the following code

```
data(iris)
iris <- as_tibble(iris)
```

What is the `iris` dataset

The `iris` dataset contains the following observations on the 50 flowers from 3 different species

variable	units
Sepal.Length	centimeters
Sepal.Width	centimeters
Petal.Length	centimeters
Petal.Width	centimeters
Species	Setosa, Versicolor, Virginica

First taste of rstatix

We can use `rstatix` to easily generate summary statistics for our variables:

```
iris %>%  
  get_summary_stats() %>%  
  select(variable, min, q1,  
         median, q3, max) # Just so it all fits
```

```
## # A tibble: 4 x 6  
##   variable      min    q1 median    q3    max  
##   <chr>      <dbl> <dbl> <dbl> <dbl> <dbl>  
## 1 Petal.Length    1    1.6  4.35    5.1    6.9  
## 2 Petal.Width     0.1    0.3   1.3    1.8    2.5  
## 3 Sepal.Length    4.3    5.1   5.8    6.4    7.9  
## 4 Sepal.Width     2    2.8   3      3.3    4.4
```

1 sample t-test

What is wrong with base R

When performing hypothesis testing in normal R, it is functional and pretty to look at, but hard to use.

For example, consider the following hypothesis test. Let μ_{pl} be the true mean petal length of iris flowers in centimeters, and consider the hypothesis

$$H_0 : \mu_{pl} = 5 \quad \text{vs} \quad H_a : \mu_{pl} \neq 5.$$

What is wrong with base R

```
(petal_t_test <- t.test(iris$Petal.Length, mu = 5))
```

```
##  
## One Sample t-test  
##  
## data: iris$Petal.Length  
## t = -8.6169, df = 149, p-value = 9.235e-15  
## alternative hypothesis: true mean is not equal to 5  
## 95 percent confidence interval:  
## 3.473185 4.042815  
## sample estimates:  
## mean of x  
## 3.758
```

Getting things from base R

If we want to access the information about the t-test, we need to use `$` notation to get it, i.e.

```
petal_t_test$statistic
```

```
##           t  
## -8.616862
```

```
petal_t_test$p.value
```

```
## [1] 9.235041e-15
```

rstatix brings dataframes!

The beauty of `rstatix` is that it takes all of this information and bundles it into an easy to use data frame (good for things like plotting). We perform the t-test as follows:

```
petal_t_test_rstatix <- iris %>%  
  t_test(Petal.Length ~ 1, mu = 5)  
petal_t_test_rstatix
```

```
## # A tibble: 1 x 7  
##   .y.      group1 group2      n statistic    df      p  
## * <chr>      <chr> <chr>   <int>    <dbl> <dbl>  <dbl>  
## 1 Petal.Length 1      null model   150    -8.62   149 9.24e-15
```

Your turn

What to do

1. Obtain summary statistics for each Species of iris in the `iris` dataset. Which species has the largest average petal length? Which species has the smallest average petal width?

Solutions:

This is done simply using “`group_by`”

```
iris %>%  
  group_by(Species) %>%  
  get_summary_stats() %>%  
  select(Species, variable, min, q1, median, q3, max)
```

```
## # A tibble: 12 x 7  
##   Species    variable    min    q1 median    q3    max  
##   <fct>      <chr>      <dbl> <dbl> <dbl> <dbl> <dbl>  
## 1 setosa    Petal.Length    1    1.4    1.5    1.58    1.9  
## 2 setosa    Petal.Width     0.1  0.2    0.2    0.3     0.6  
## 3 setosa    Sepal.Length    4.3  4.8    5      5.2     5.8  
## 4 setosa    Sepal.Width     2.3  3.2    3.4    3.68    4.4  
## 5 versicolor Petal.Length    3    4      4.35   4.6     5.1
```

```
## 6 versicolor Petal.Width      1      1.2      1.3      1.5      1.8
## 7 versicolor Sepal.Length    4.9      5.6      5.9      6.3      7
## 8 versicolor Sepal.Width      2      2.52      2.8      3      3.4
## 9 virginica  Petal.Length    4.5      5.1      5.55     5.88     6.9
## 10 virginica Petal.Width      1.4      1.8      2      2.3      2.5
## 11 virginica Sepal.Length    4.9      6.22      6.5      6.9      7.9
## 12 virginica Sepal.Width      2.2      2.8      3      3.18     3.8
```

Solutions:

Virginica has the largest average petal length as measured by median. Setosa has the smallest average petal width as measured by median.

2. Let μ_{pw} be the true mean petal width of iris flowers in centimeters, and consider the hypothesis

$$H_0 : \mu_{pw} = 1.3 \quad \text{vs} \quad H_a : \mu_{pw} \neq 1.3.$$

Test this hypothesis at the $\alpha = 0.05$ level of significance using `rstatix`. Conclude your hypothesis in context.

Solutions:

```
petal_width_t_test_rstatix <- iris %>%
  t_test(Petal.Width ~ 1, mu = 1.3)
petal_width_t_test_rstatix
```

```
## # A tibble: 1 x 7
##   .y.      group1 group2      n statistic    df      p
## * <chr>    <chr> <chr>   <int>    <dbl> <dbl> <dbl>
## 1 Petal.Width 1      null model   150    -1.62   149 0.108
```

Solutions:

Since the p-value is $0.108 > 0.05$, there is insufficient evidence to suggest the true mean petal width of iris flowers is different from 1.3 cm.

3. Perform the hypothesis test from Q2 for each species of iris flower in the `iris` dataset. What do you notice?
-

Solutions:

```
petal_width_t_test_rstatix_species <- iris %>%
  group_by(Species) %>%
  t_test(Petal.Width ~ 1, mu = 1.3)
petal_width_t_test_rstatix_species
```

```
## # A tibble: 3 x 8
##   Species   .y.      group1 group2      n statistic    df      p
## * <fct>   <chr>   <chr>   <chr>   <int>   <dbl> <dbl>   <dbl>
## 1 setosa   Petal.Width 1      null model    50   -70.7    49 5.43e-51
## 2 versicolor Petal.Width 1      null model    50    0.930    49 3.57e- 1
## 3 virginica Petal.Width 1      null model    50    18.7    49 6.09e-24
```

Solutions:

At the 5% level, both setosa and virginica have petal width different to 1.3 cm, i.e. we would reject the hypothesis. This shows the power of having more information. If we lump them all together, we get a different picture than if we look at each species individually. This is shown nicely in the following two boxplots.

```
iris %>%
  ggplot(aes(y = Petal.Width)) +
  geom_boxplot() +
  theme_classic()
```

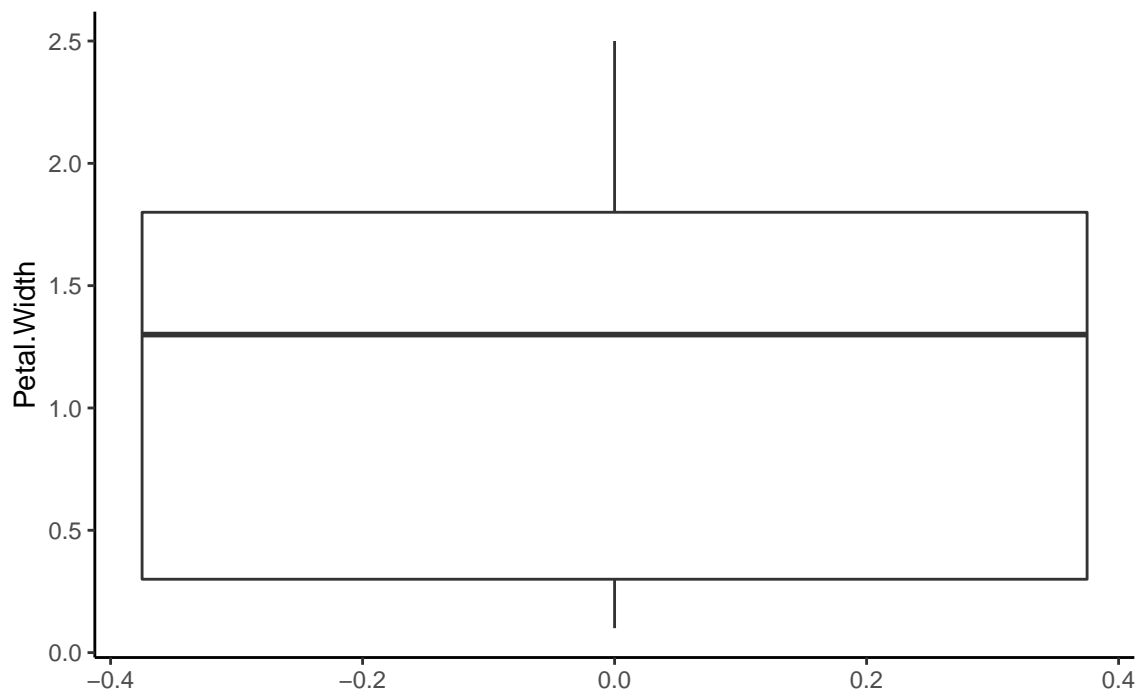


Figure 1: Boxplot of petal width for all iris flowers together from the iris dataset.

```
iris %>%
  ggplot(aes(x = Species, y = Petal.Width, fill = Species)) +
  geom_boxplot() +
  scale_fill_manual(values = cols) + # This line of code won't run for you
  theme_classic()
```

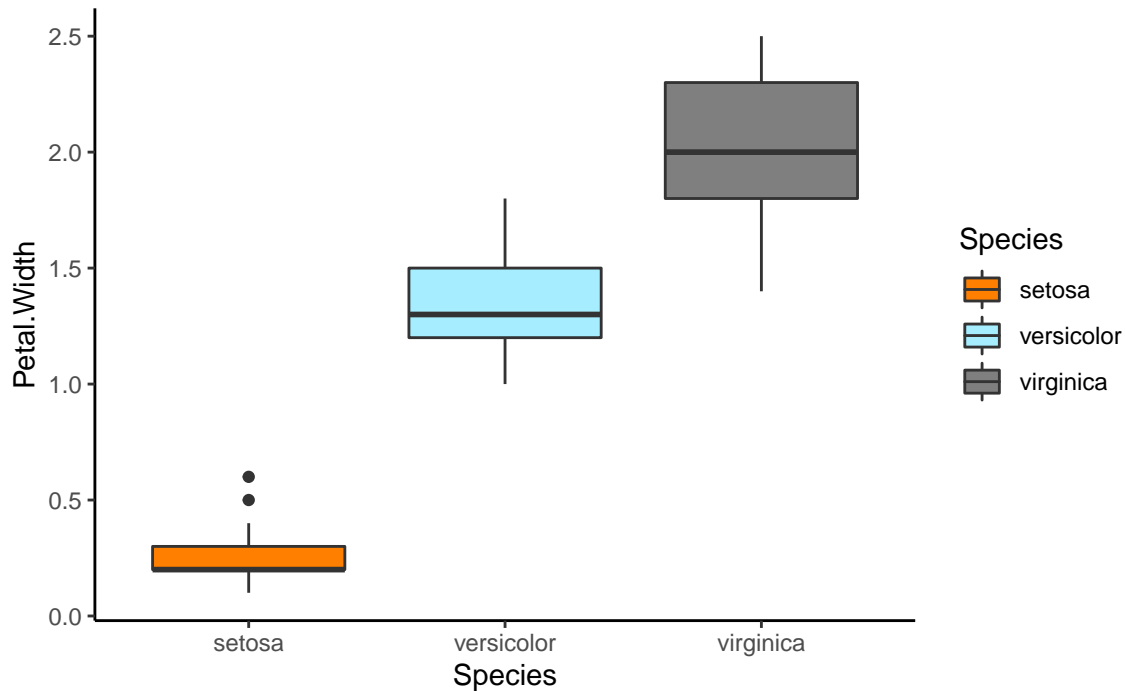


Figure 2: Side-by-side boxplots of petal width for each species of iris flower from the iris dataset.

Two sample t-tests

Let's narrow our scope

We will consider the two groups of flowers from the versicolor species and the virginica species. We can subset our data as

```
iris_sub <- iris %>%
  filter(Species != "setosa") %>%
  mutate(Species = fct_drop(Species))
head(iris_sub)
```

```
## # A tibble: 6 x 5
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
##   <dbl>         <dbl>         <dbl>         <dbl> <fct>
## 1         7         3.2         4.7         1.4 versicolor
## 2         6.4        3.2         4.5         1.5 versicolor
## 3         6.9        3.1         4.9         1.5 versicolor
## 4         5.5        2.3         4          1.3 versicolor
## 5         6.5        2.8         4.6         1.5 versicolor
## 6         5.7        2.8         4.5         1.3 versicolor
```

The hypothesis to test

Let $\mu_{versi-pl}$ and $\mu_{virgin-pl}$ be the true mean petal length of versicolor and virginica iris flowers in centimeters, respectively. We will test the hypothesis

$$H_0 : \mu_{\text{versi-pl}} = \mu_{\text{virgin-pl}} \quad \text{vs} \quad H_a : \mu_{\text{versi-pl}} \neq \mu_{\text{virgin-pl}}.$$

at the $\alpha = 0.05$ level.

The how-to in rstatix

This is simple in `rstatix` in the following way:

```
two_sample_petal_length <- iris_sub %>%
  t_test(Petal.Length ~ Species)
two_sample_petal_length
```

```
## # A tibble: 1 x 8
##   .y.      group1      group2      n1      n2 statistic      df      p
## * <chr>      <chr>      <chr>    <int> <int>    <dbl> <dbl>  <dbl>
## 1 Petal.Length versicolor virginica     50     50    -12.6   95.6 4.9e-22
```

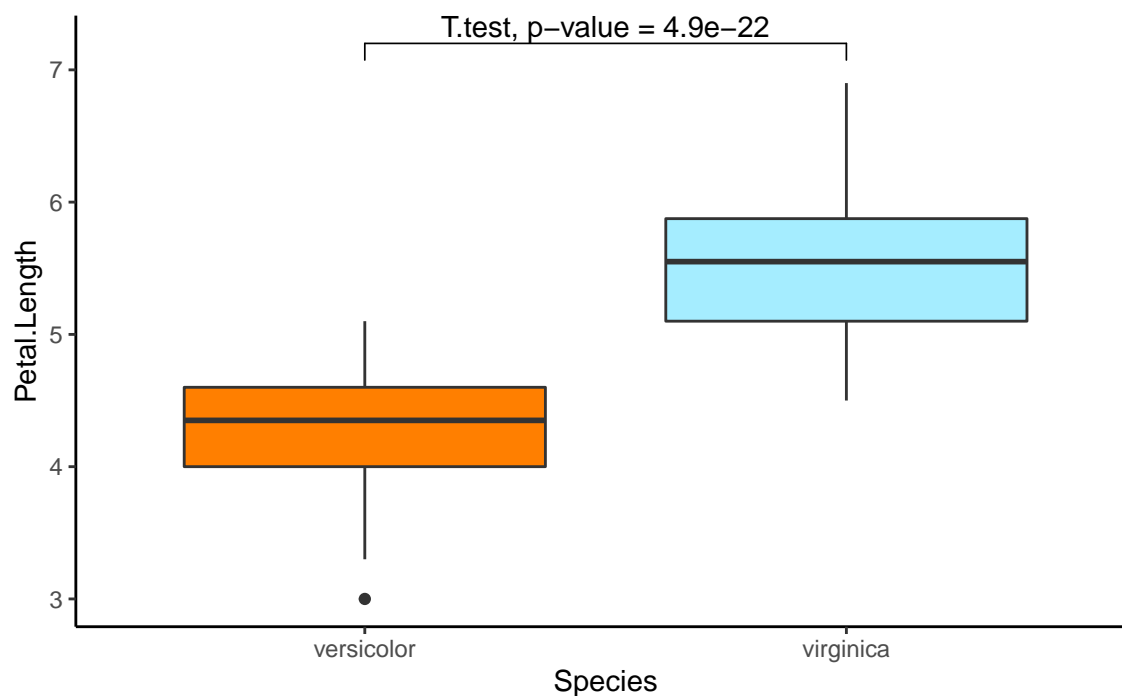
Conclusion?

Since the p-value is MUCH smaller than 0.05, we reject the hypothesis and conclude that there is sufficient evidence to suggest the true mean petal length between versicolor and virginica is different.

Let's visualise this! (the code)

```
p <- iris_sub %>%
  ggplot(aes(x = Species, y = Petal.Length)) +
  geom_boxplot(aes(fill = Species), show.legend = FALSE) +
  scale_fill_manual(values = cols) + # This line of code won't run for you
  theme_classic() +
  stat_pvalue_manual(two_sample_petal_length, # This is where ggpubr comes in!
                    label = "T.test, p-value = {p}",
                    y.position = 7.2)
```


Let's visualise this! (the plot)



Your turn

What to do

1. Let $\mu_{versi-pw}$ and $\mu_{setosa-pw}$ be the true mean petal width of versicolor and setosa iris flowers in centimeters, respectively. Test the hypothesis

$$H_0 : \mu_{versi-pw} = \mu_{setosa-pw} \quad \text{vs} \quad H_a : \mu_{versi-pw} \neq \mu_{setosa-pw} .$$

at the $\alpha = 0.05$ level. Conclude your hypothesis test in context.

Solutions:

First, we create the new data

```
iris_pw_data <- iris %>%  
  filter(Species != "virginica") %>%  
  mutate(Species = fct_drop(Species))
```

Solutions:

Now, we conduct the hypothesis test

```
two_sample_petal_width <- iris_pw_data %>%
  t_test(Petal.Width~ Species)
two_sample_petal_width
```

```
## # A tibble: 1 x 8
##   .y.      group1 group2      n1      n2 statistic    df      p
## * <chr>    <chr>  <chr>    <int> <int>    <dbl> <dbl>  <dbl>
## 1 Petal.Width setosa versicolor    50    50    -34.1  74.8 2.72e-47
```

Solutions:

Since the p-value is much less than 0.05, there is sufficient evidence to suggest that the true mean petal width of setosa and versicolor iris flowers is different.

2. Visualise your results from Q1 on a boxplot.

Solutions:

```
iris_pw_data %>%
  ggplot(aes(x = Species, y = Petal.Width)) +
  geom_boxplot(aes(fill = Species), show.legend = FALSE) +
  scale_fill_manual(values = cols) + # This line of code won't run for you
  theme_classic() +
  stat_pvalue_manual(two_sample_petal_width, # This is where ggpubr comes in!
                    label = "T.test, p-value = {p}",
                    y.position = 2)
```

3. **Challenge:** Perform pairwise t-tests on Sepal Length for each species of iris flower, and visualise this on a boxplot. Brownie points for using (and remembering) adjusted p-values. \begin{solutions}

\end{solutions}

```
sepal_length_test <- iris %>%
  t_test(Sepal.Length ~ Species)
sepal_length_test
```

```
## # A tibble: 3 x 10
##   .y.      group1 group2      n1      n2 stati~1    df      p    p.adj p.adj~2
## * <chr>    <chr>  <chr>    <int> <int>    <dbl> <dbl>  <dbl>  <dbl> <chr>
## 1 Sepal.Length setosa versi~    50    50   -10.5   86.5 3.75e-17 7.5 e-17 ****
## 2 Sepal.Length setosa virgi~    50    50   -15.4   76.5 3.97e-25 1.19e-24 ****
## 3 Sepal.Length versi~ virgi~    50    50    -5.63  94.0 1.87e- 7 1.87e- 7 ****
## # ... with abbreviated variable names 1: statistic, 2: p.adj.signif
```

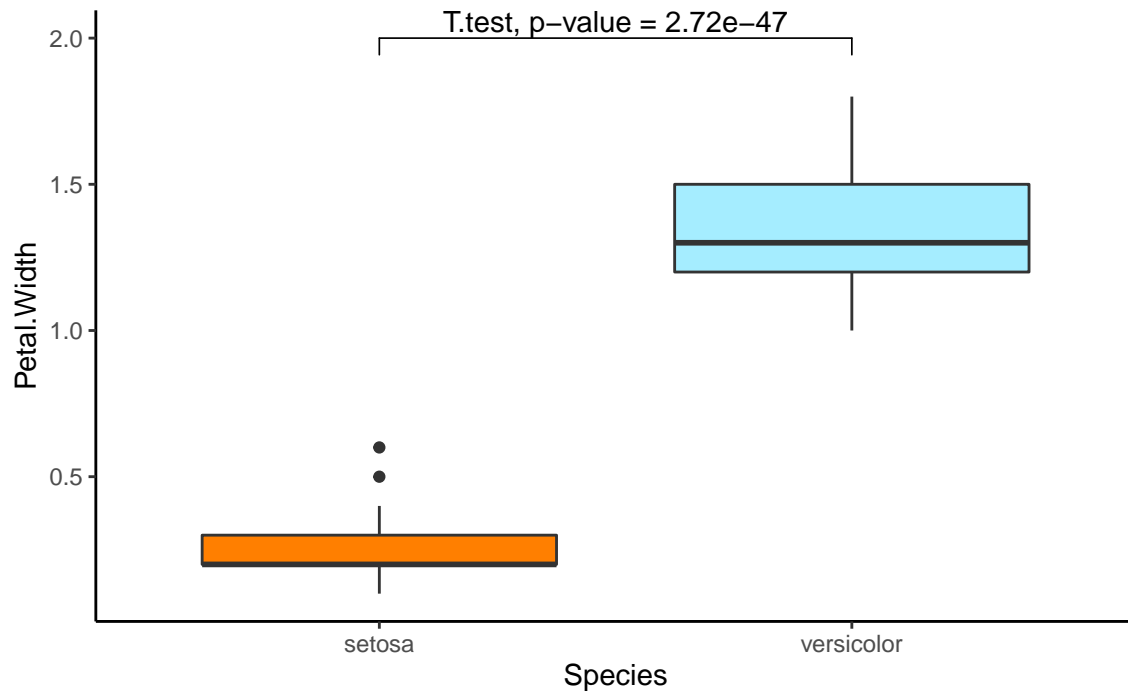


Figure 3: Side-by-side boxplots of petal width for setosa and versicolor iris flowers from the iris dataset. We indicate the p-value from a two-sample t-test on the boxplot, showing an extremely significant difference in mean petal width.

```
iris %>%
  ggplot(aes(x = Species, y = Sepal.Length)) +
  geom_boxplot(aes(fill = Species), show.legend = FALSE) +
  scale_fill_manual(values = cols) + # This line of code won't run for you
  theme_classic() +
  stat_pvalue_manual(sepal_length_test, # This is where ggpubr comes in!
    label = "T.test, p-value = {p.adj}",
    y.position = c(7, 8, 9))
```

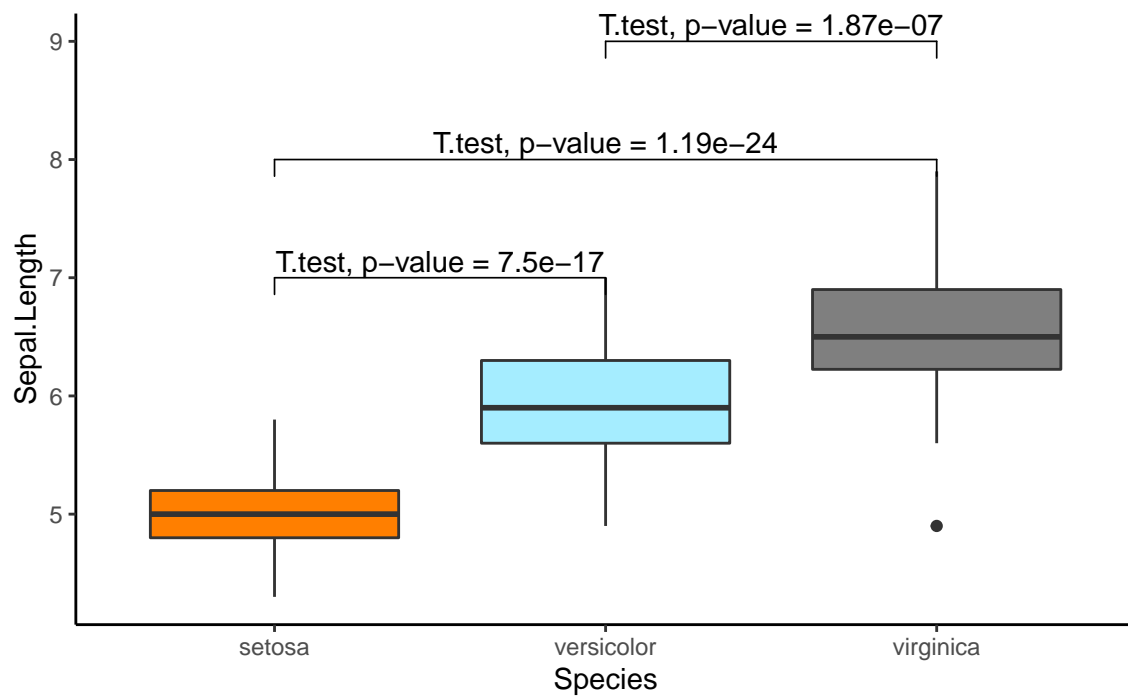


Figure 4: Side-by-side boxplots of Sepal length for each species of iris in the iris dataset. We have performed and visualised pairwise t-tests with p-values adjusted for multiple comparisons.