# R Programming: Data Analysis & Visualization

## Stacey Borrego (edited by Kristina Riemer)

## Helpful Links

- Data Carpentry: Data Analysis and Visualization in R for Ecologists
- Tidyverse: R packages for Data Science
- R for Data Science 2nd Edition by Hadley Wickham and Garret Grolemund
- Advanced R by Hadley Wickham
- ggplot2: Elegant graphics for data analysis by Hadley Wickham
- Posit Cheatsheets by Posit
- The R Gallery by Kyle W. Brown
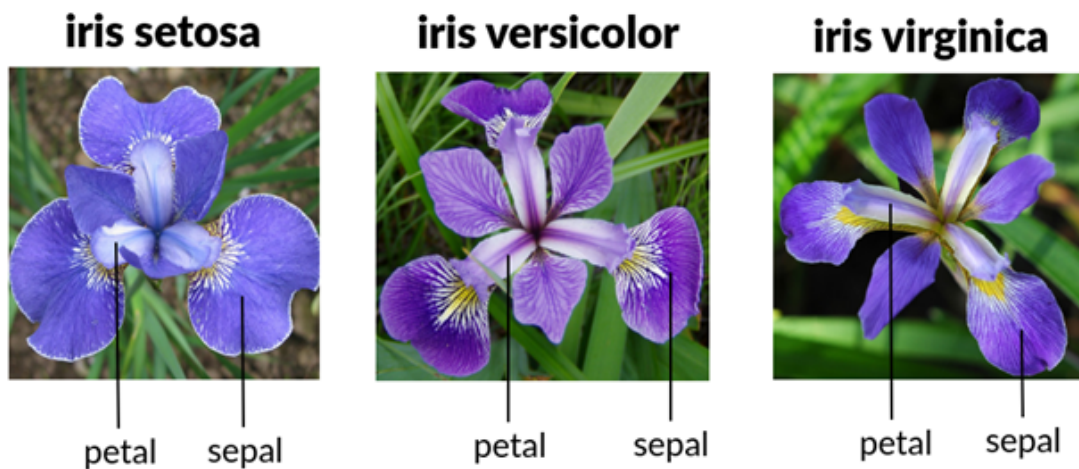- Introduction to R by Douglas, Roos, Mancini, Cuoto, & Lusseau



Figure 1: Understanding the Iris Dataset

## Starting with Data

This is what the data will look like:

| Column | Description |
| --- | --- |
| Sepal.Length | Sepal length in centimeters |
| Sepal.Width | Sepal width in centimeters |
| Petal.Length | Petal length in centimeters |
| Petal.Width | Petal width in centimeters |
| Species | Name of species name |

```
# To see data sets available
data()

# Load the data set of choice
# use the data set name as the argument to `data()`
data(iris)
```

Open the dataset in RStudio's Data Viewer

```
# Check the data type of the data just loaded
class(iris)

# View the whole data set
View(iris)
```

Inspect the data

```
# See a few rows of the data
head(iris)
tail(iris)

# See a specific number of rows of the data
head(iris, n = 10)
tail(iris, n = 10)
```

- Size:
  - `dim(iris)` - returns a vector with the number of rows in the first element, and the number of columns as the second element (the dimensions of the object)
  - `nrow(iris)` - returns the number of rows
  - `ncol(iris)` - returns the number of columns
- Content:
  - `head(iris)` - shows the first 6 rows
  - `tail(iris)` - shows the last 6 rows
- Names:

- `names(iris)` - returns the column names (synonym of colnames() for data.frame objects)
- `rownames(iris)` - returns the row names

- Summary:

  - `str(iris)` - structure of the object and information about the class, length and content of each column
  - `summary(iris)` - summary statistics for each column

**Indexing and subsetting data frames**

A data frame has rows and columns in 2 dimensions. To extract specific data, specify the "coordinates" in `[ ]` indicating row numbers first followed by column numbers.

```
# dim() gives the output of the rows and columns of the object
dim(iris)

# Extract specific values by specifying row and column indices
# in the format:
# data_frame[row_index, column_index]

# Extract the first row and first column from iris
iris[1, 1]

# Extract the first row and sixth column
iris[1, 6]

# To select all columns, leave the column index blank
# Extract the first row and all columns
iris[1, ]

# Extract all rows and the first column
iris[, 1]

# An even shorter way to select first column across all rows:
iris[1] # No comma!

# Select multiple rows or columns with vectors
# Extract the first three rows of the 1st and 2nd columns
iris[c(1, 2, 3), c(1, 2)]

# We can use the : operator to create a range and select all listed vectors
iris[1:3, 1:2]
```

```
# This is equivalent to head(iris)
iris[1:6, ]

# Subsetting with single square brackets ("[]") always returns a data frame.
# If you want a vector, use double square brackets ("[[]]")

# For instance, to get a column as a vector:
iris[[1]]
iris[["Species"]]

# To get the first value in our data frame:
iris[[1, 1]]
iris[1, "Species"]

# Data frames can be subset by calling their column names directly
# Use the $ operator with column names to return a vector
iris$Species
```

```
# Exceptions: subsetting the whole data frame, except the first
# and second columns
iris[, -c(1, 2)]

# Exceptions: subsetting all columns of the data frame
# except for the first six rows
iris[-(7:nrow(iris)), ]
```

## Plotting

Helpful links:

- [Installation](#)
- [ggplot2: Elegant Graphics for Data Analysis](#)

Plotting can be done in many different ways in R. Base R has several basic plotting functions which can be valuable for a quick peek at your data. However, a common and customizable option is to use the package `ggplot2`.

Every ggplot2 plot has three key components:

- **Data** (`data`)
- A set of **aesthetic mappings** between variables in the data and visual properties (`aes(x, y)`)
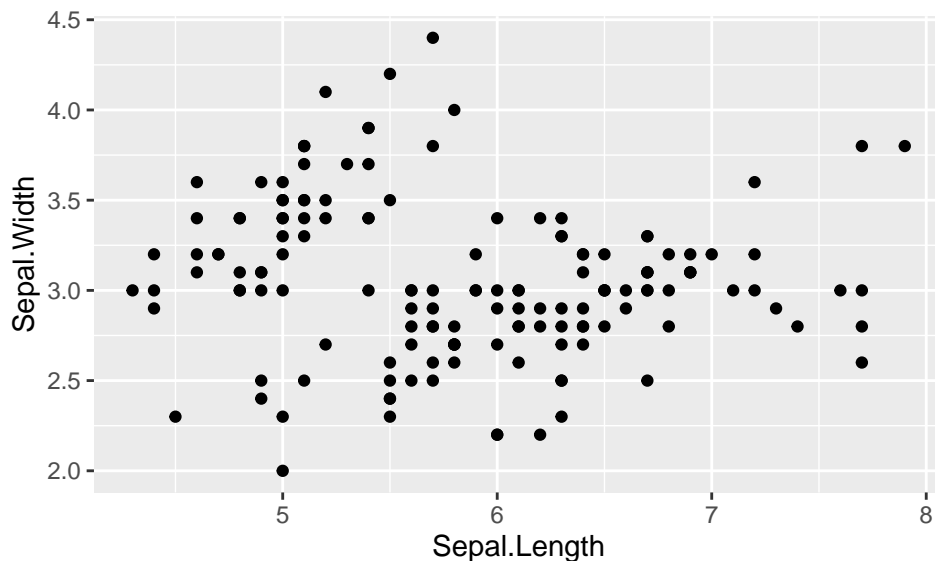
- At least one **layer** which describes how to render each observation. Layers are usually created with a geom function. It is important to note that layers are added with a **+**.

```r
# Install the following library
# install.packages("ggplot2")

# Load the ggplot2 package to access its unique functions
library(ggplot2)



# Step 1: Make the plot object by indicating the input data
# and plot aesthetics
# Example of creating the plot object for a scatter plot
plot_1 <- ggplot(data = iris,
                 aes(x = Sepal.Length,
                     y = Sepal.Width))

# Step 2: Add layers to the plot object to plot the data. Here we add
# the points of the scatter plot
# Example of a basic scatter plot adding a `geom_point()` layer
plot_1 +
  geom_point()
```
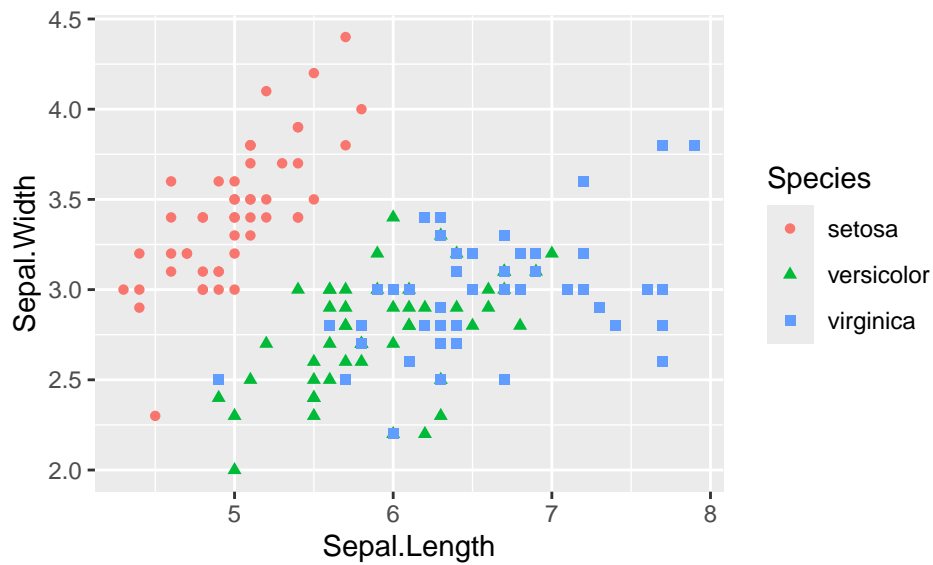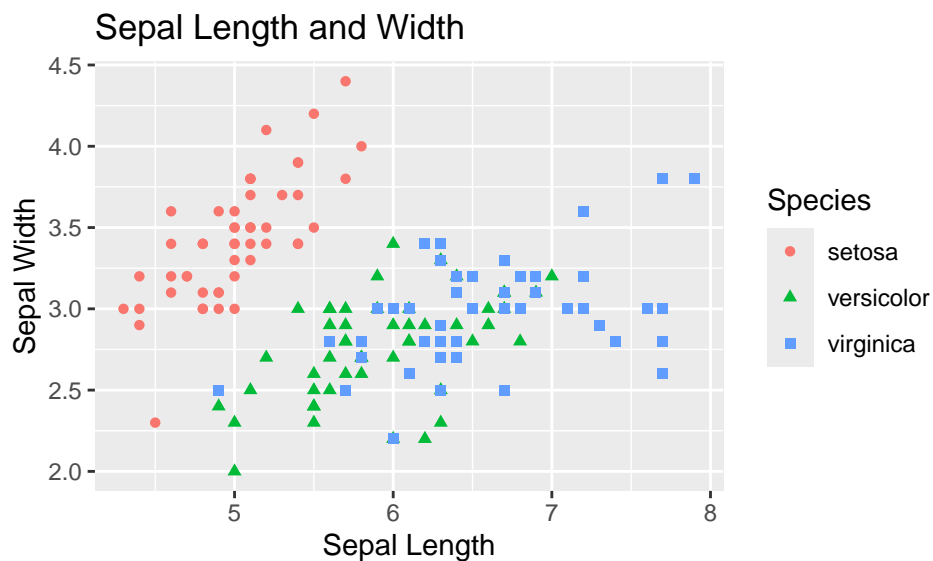


```r
# Find out the options for `geom_point()`
?geom_point

# Step 3: Modify each layer using specific arguments and aesthetics options
# Example of a colorful scatter plot adding aesthetic mappings to the
```

```
# `geom_point()` layer
plot_1 +
  geom_point(aes(color = Species,
                 shape = Species))
```


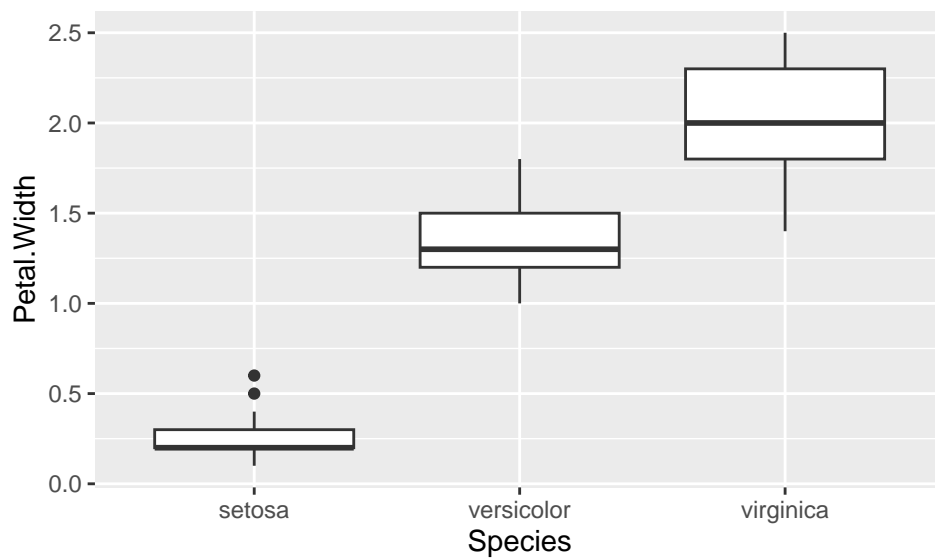
```
# Add labels using `labs()`
plot_1 +
  geom_point(aes(color = Species,
                 shape = Species)) +
  labs(x = "Sepal Length",
       y = "Sepal Width",
       title = "Sepal Length and Width")
```
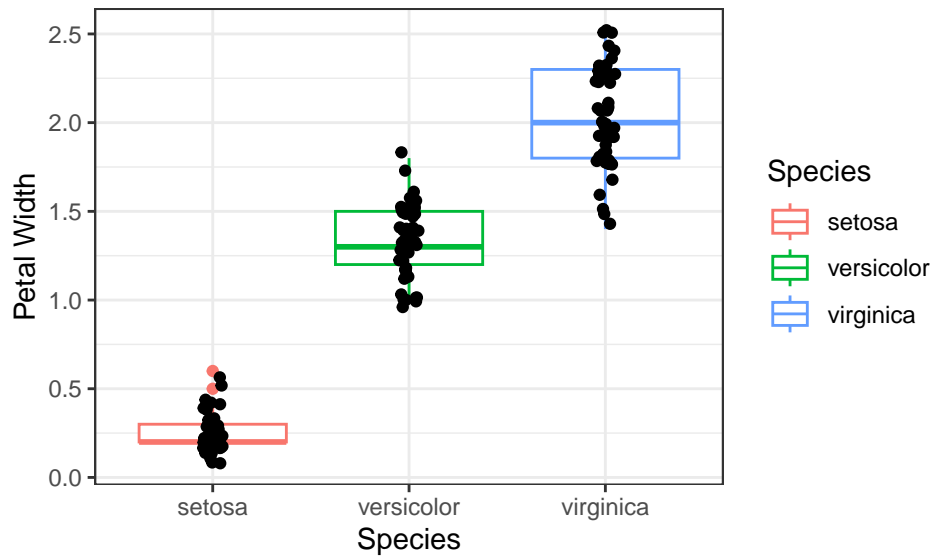
```
# Boxplot
plot_2 <- ggplot(data = iris,
                 aes(x = Species,
                     y = Petal.Width))

# Basic boxplot adding a `geom_boxplot()` layer
# Important note: the x value must be factored or grouped to have individual
# boxplots returned
plot_2 +
  geom_boxplot()
```



```
# Colorful boxplot adding more layers and aesthetic mappings to the data
plot_2 +
  geom_boxplot(aes(color = Species)) +
  geom_jitter(aes(x = Species),
              position = position_jitter(width = 0.05)) +
  labs(y = "Petal Width") +
  theme_bw()
```
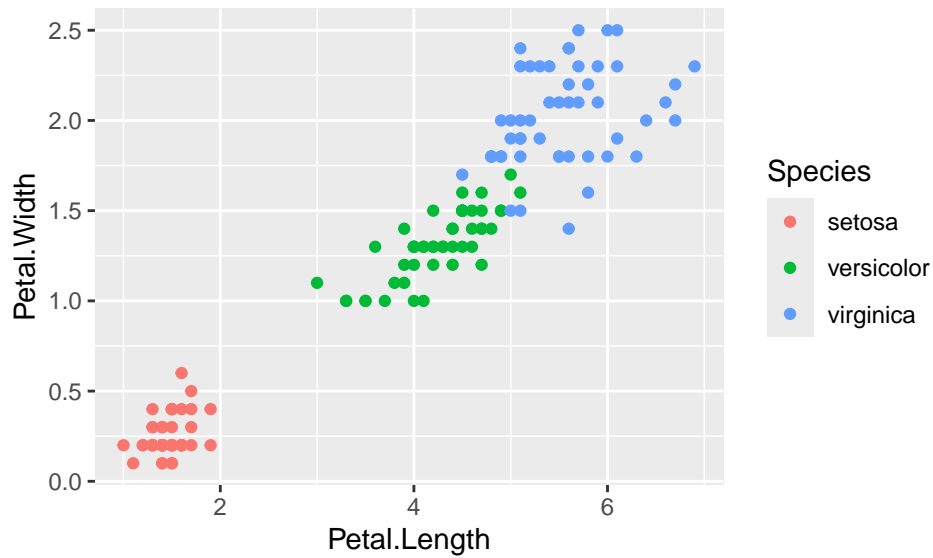
R comes with default colors and so does ggplot2. These are nice enough but sometimes it is nice to customize our plots. There are few ways we can do this.

- RColorBrewer: Such a helpful tool - contains color palettes and can make a range of colors of your choice. Must be installed and loaded prior to use:

  - install.packages("Rcolorbrewer")
  - library(Rcolorbrewer)

- viridisLite: This is part of ggplot2 that provides color palettes. This does not need to be installed and loaded as it comes with ggplot2.

  - Another resource

- Colors in R: A list of all the color names available in R

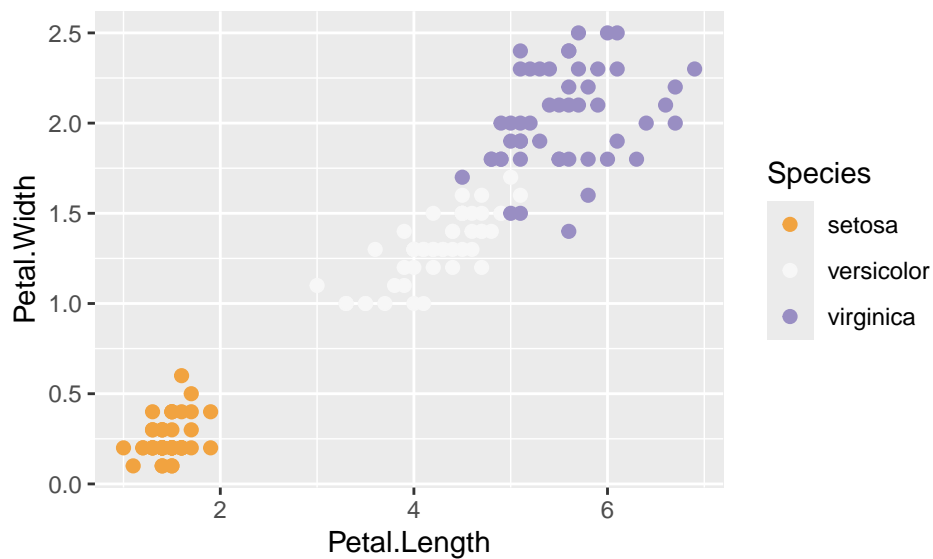- You can also provide hexidecimal values. HERE is a helpful tool.

```
# Creating the base plot. Here we have indicated the data set and its
# mappings.We have also added which column should be used for color
# and fill mappings.
plot_3 <- ggplot(data = iris,
                 aes(x = Petal.Length,
                     y = Petal.Width,
                     color = Species,
                     fill = Species))

# Basic boxplot adding a `geom_boxplot()` layer
plot_3 +
  geom_point()
```
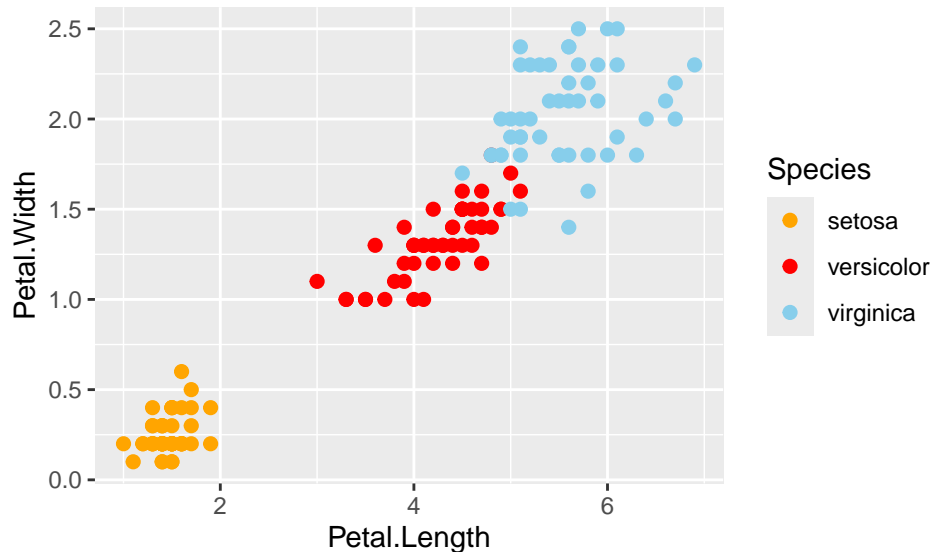
```
# RColorBrewer Examples
# install.packages("RColorBrewer")
library(RColorBrewer)

# RColorBrewer Example 1
plot_3 +
  geom_point(size = 2) +
  scale_color_brewer(palette = "PuOr")
```



```
# Using Your Own Colors, Example 2
plot_3 +
```

```
  geom_point(size = 2) +
  scale_color_manual(values = c("orange", "red", "skyblue"))
```



## Factors

This can be the cause of many issues when plotting data. Always check if the data is factored and determined whether is should or should not be to get the correct plot.

R has a special class for working with categorical data, called factors. Once created, factors can only contain a pre-defined set of values, known as levels. Factors are stored as integers associated with labels and they can be ordered or unordered. While factors look (and often behave) like character vectors, they are actually treated as integer vectors by R. So you need to be very careful when treating them as strings.

By default, R always sorts levels in alphabetical order. In the example below, R will assign 1 to the level "experimental" and 2 to the level "wild type" (because "e" comes before "w", even though the first element in this vector is "wild type").

```
sample <- factor(c("wild type", "experimental", "wild type", "experimental"))

# Different ways to see the levels of factored data
sample
levels(sample)
nlevels(sample)
```

### Reorder factors

```r
sample <- factor(sample,
                 levels = c("wild type", "experimental"))

sample
levels(sample)
nlevels(sample)
```

When working with data in a data frame, the columns that contain text are not automatically coerced (= converted) into the factor data type, but once we have loaded the data we can do the conversion using the factor() function

```r
head(iris)
iris$Species <- factor(iris$Species)
head(iris)

summary(iris$Species)
levels(iris$Species)
nlevels(iris$Species)
```

### Converting Factors

Sometimes it is necessary to convert data from one type to another

```r
sample <- factor(sample,
                 levels = c("wild type", "experimental"))

as.character(sample)
sample_char <- as.character(sample)
sample_char

as.numeric(sample)
```

### Renaming Factors

```r
# Inspect the contents of the column
levels(iris$Species)

# Replace level names with a new vector of names
```

```
levels(iris$Species) <- c("Setosa", "Versicolor", "Virginica")
ggplot(data = iris, aes(x = Sepal.Length, y = Sepal.Width, color = Species, shape = Spec
  geom_point() +
  labs(x = "Sepal Length", y = "Sepal Width", title = "Sepal Length & Width") +
  #scale_color_brewer(palette = "PuOr")
  scale_color_manual(values = c("orange", "red", "blue"))

# Check names of factors have been replaced within data frame
levels(iris$Species)
```

# Manipulating, analyzing and exporting data with tidyverse

## Data manipulation using dplyr and tidyr

dplyr is a package for helping with tabular data manipulation. It pairs nicely with tidyr which enables you to convert between different data formats for plotting and analysis.

The tidyverse package is an "umbrella-package" that installs tidyr, dplyr, and several other useful packages for data analysis, such as ggplot2, tibble, etc.

dplyr functions

- `select()` subset columns
- `filter()` subset rows on conditions
- `mutate()` create new columns by using information from other columns
- `group_by()` and `summarize()` create summary statistics on grouped data
- `arrange()` sort results
- `count()` count discrete values

### Selecting columns and filtering rows

- `select()` subset columns
  - The first argument is the data frame (iris), and the subsequent arguments are the columns to keep.
- `filter()` subset rows on conditions

```
# Install Tidyverse
# install.packages("tidyverse")

# Load the umbrella package Tidyverse
library(tidyverse)
```

```
## -- Attaching core tidyverse packages ---------------------- tidyverse 2.0.0 --
## v dplyr     1.1.4     v readr     2.1.5
## v forcats   1.0.0     v stringr   1.5.1
## v lubridate 1.9.3     v tibble    3.2.1
## v purrr     1.0.2     v tidyr     1.3.1
## -- Conflicts ------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
## i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts
```

```r
# Reloading data to ensure nothing has been manipulated
data(iris)
```

```r
# Select specific columns
select(iris,
       Species, Petal.Length, Petal.Width)
```

Select all columns except certain ones by using a "-" in front of the variable to exclude it

```r
select(iris,
       -Sepal.Length, -Sepal.Width)
```

Choose rows based on specific criterion with `filter()`

```r
filter(iris, Species == "virginica")
```

**Pipes**

What if you want to select and filter at the same time? There are three ways to do this: use intermediate steps, nested functions, or pipes.

Pipes let you take the output of one function and send it directly to the next, which is useful when you need to do many things to the same data set. Pipes in R look like `%>%` and are made available via the magrittr package, installed automatically with dplyr. If you use RStudio, you can type the pipe with `Ctrl` + `Shift` + `M` if you have a PC or `Cmd` + `Shift` + `M` if you have a Mac.

```r
# Filter rows and save as an object
small <- filter(iris,
                Petal.Length < 5)

# Select columns and save as another object
small_x <- select(small,
```

```
                  Species, Petal.Length, Petal.Width)
small_x

# Combine the two operations in one piece of code using pipes
iris %>%
  filter(Petal.Length < 5) %>%
  select(Species, Petal.Length, Petal.Width)
```

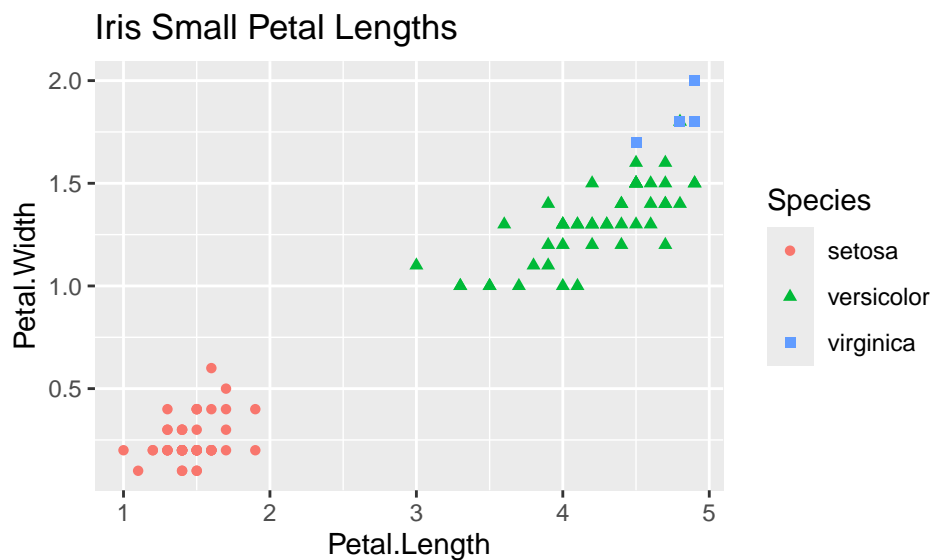Create a plot using the subset data

```
# Create a new object with the subset data
iris_petals <- iris %>%
  filter(Petal.Length < 5) %>%
  select(Species, Petal.Length, Petal.Width)

# Use new data to create a plot object
petal_plot <- ggplot(iris_petals,
                aes(x = Petal.Length,
                    y = Petal.Width))

# Add layers to plot object
petal_plot +
  geom_point(aes(shape = Species,
                 color = Species)) +
  labs(title = "Iris Small Petal Lengths")
```
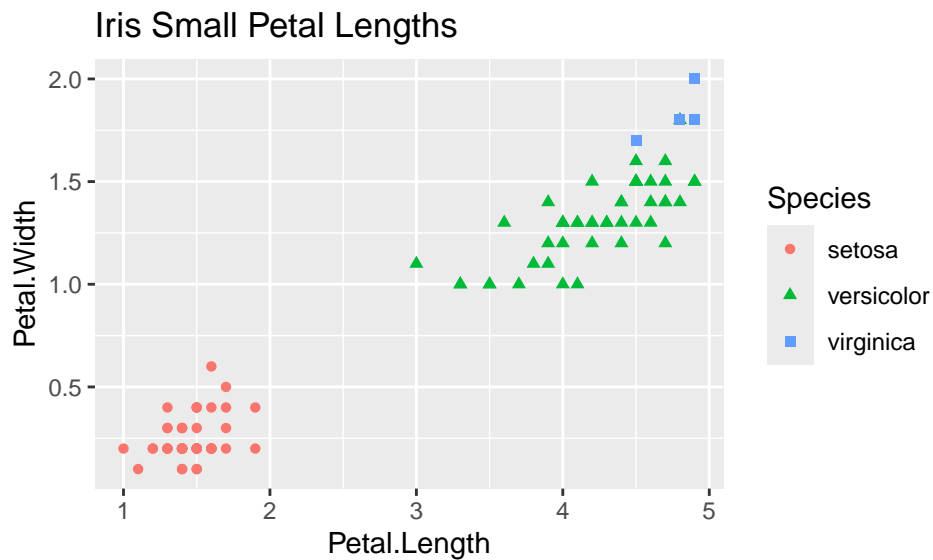


Write the subset data and plot in one piece of code

```
iris %>%
  filter(Petal.Length < 5) %>%
  select(Species, Petal.Length, Petal.Width) %>%
  ggplot(aes(x = Petal.Length,
             y = Petal.Width)) +
  geom_point(aes(shape = Species,
                 color = Species)) +
  labs(title = "Iris Small Petal Lengths")
```
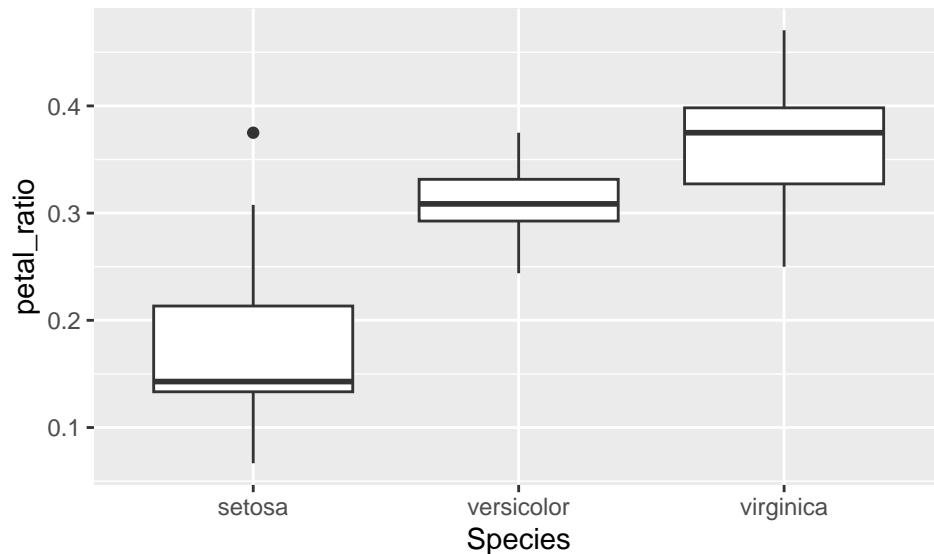


## Mutate

Create new columns based on the values in existing columns using `mutate()`

```
iris %>%
  mutate(petal_ratio = Petal.Width / Petal.Length,
         petal_percent = petal_ratio * 100)
```

Add a plot in one piece of code

```
iris %>%
  mutate(petal_ratio = Petal.Width / Petal.Length,
         petal_percent = petal_ratio * 100) %>%
  ggplot(aes(x = Species,
             y = petal_ratio)) +
  geom_boxplot()
```

**Split-apply-combine data analysis and the `summarize()` function**

Many data analysis tasks can be approached using the **split-apply-combine paradigm**

- Split the data into groups
- Apply some analysis to each group
- Combine the results

Key functions of dplyr for this workflow are `group_by()` and `summarize()`. `group_by()` is often used together with `summarize()`, which collapses each group into a single-row summary of that group.

```
iris %>%
  group_by(Species) %>%
  summarize(petal_width_mean = mean(Petal.Width, na.rm = TRUE))
```

Once the data are grouped, you can also summarize multiple variables at the same time. Add a column indicating the minimum and maximum width for each species.

```
iris %>%
  filter(!is.na(Petal.Width)) %>%
  group_by(Species) %>%
  summarize(petal_width_mean = mean(Petal.Width),
            petal_width_min = min(Petal.Width))
```

## Exporting plots

`ggsave()` allows you easily change the dimension and resolution of your plot by adjusting the appropriate arguments (width, height, and dpi)

```r
my_plot <- ggplot(data = iris, aes(x = Sepal.Length, y = Sepal.Width)) +
  geom_point()

ggsave("path/name_of_file.png", my_plot, width = 15, height = 10)
```