

STATS 266 Handout - Data Visualization

Qi Wang

2025-04-24

Contents

1	Introduction	3
2	Reading the Data	3
3	Numerical Data	4
3.1	Scatter Plots	5
3.2	Line Chart / Time Series Plot	7
3.3	Histogram and Density	10
3.4	Correlation Matrix Plot	11
4	Categorical Data	13
4.1	Grouping	13
4.2	Boxplot	15
4.3	Bar plot	16
5	Special Cases	18
5.1	Heatmap and Contour Plot	18
5.2	Simulated 2D Gaussian bump	20
6	Personalize	22
6.1	Customizing ggplot2 Plots	22
6.2	Summary	29

7	Exerciese: Iris Dataset	29
7.1	Introduction	29
7.2	Load the Dataset	30
7.3	Exercise 1: Scatter Plot for Relationship Between Two Numerical Variables	30
7.4	Exercise 2: Boxplot for Comparison Across Categories	31
7.5	Exercise 3: Pairwise Scatter Plot	32
7.6	Exercise 4: Violin Plot for Distribution Comparison	34
7.7	Exercise 5: Correlation Plot Between Numerical Variables	35
7.8	Summary	36
8	More..	36
9	Ackowledgement	36

1 Introduction

Welcome to **STATS 266: Introduction to R**. This handout provides an introduction about data visualization in R. By the end of this document, you should be able to:

- To be able to use `ggplot2` to generate publication-quality graphics.
- To apply geometry, aesthetic, and statistics layers to a `ggplot` plot.
- To manipulate the aesthetics of a plot using different colors, shapes, and lines.
- To improve data visualization through transforming scales and paneling by group.
- To save a plot created with `ggplot` to disk.

For this part, valuable materials to refer to include <https://ggplot2.tidyverse.org/> and <https://swcarpentry.github.io/r-novice-gapminder/08-plot-ggplot2.html>.

2 Reading the Data

Before doing any visualization, we should read the data into our environment. At the very beginning of the course, I introduced working for a project, and `here::here()` function. That's for the file path to the document. With this function, you will not need to specify the file path again.

R can read a variety of data types, including structured and unstructured formats. Below are the common types of data that R can read along with the functions used to import them:

- **CSV Files:** `read.csv("file.csv")` or `readr::read_csv("file.csv")`
- **Excel Files:** `readxl::read_excel("file.xlsx")`
- **Text Files (TSV, Fixed Width, etc.):**
 - `read.table("file.txt")`
 - `read.delim("file.txt")`
 - `readr::read_delim("file.txt", delim = "\t")`
- **JSON Files:** `jsonlite::fromJSON("file.json")`
- **XML Files:** `XML::xmlParse("file.xml")`
- **SPSS, SAS, and Stata Files:**
 - `haven::read_sav("file.sav")` (SPSS)
 - `haven::read_sas("file.sas7bdat")` (SAS)
 - `haven::read_dta("file.dta")` (Stata)
- **R Binary Files:**
 - `load("file.RData")`

```
– readRDS("file.rds")
```

- **Database Connections:**

```
– DBI::dbReadTable(con, "table_name")
```

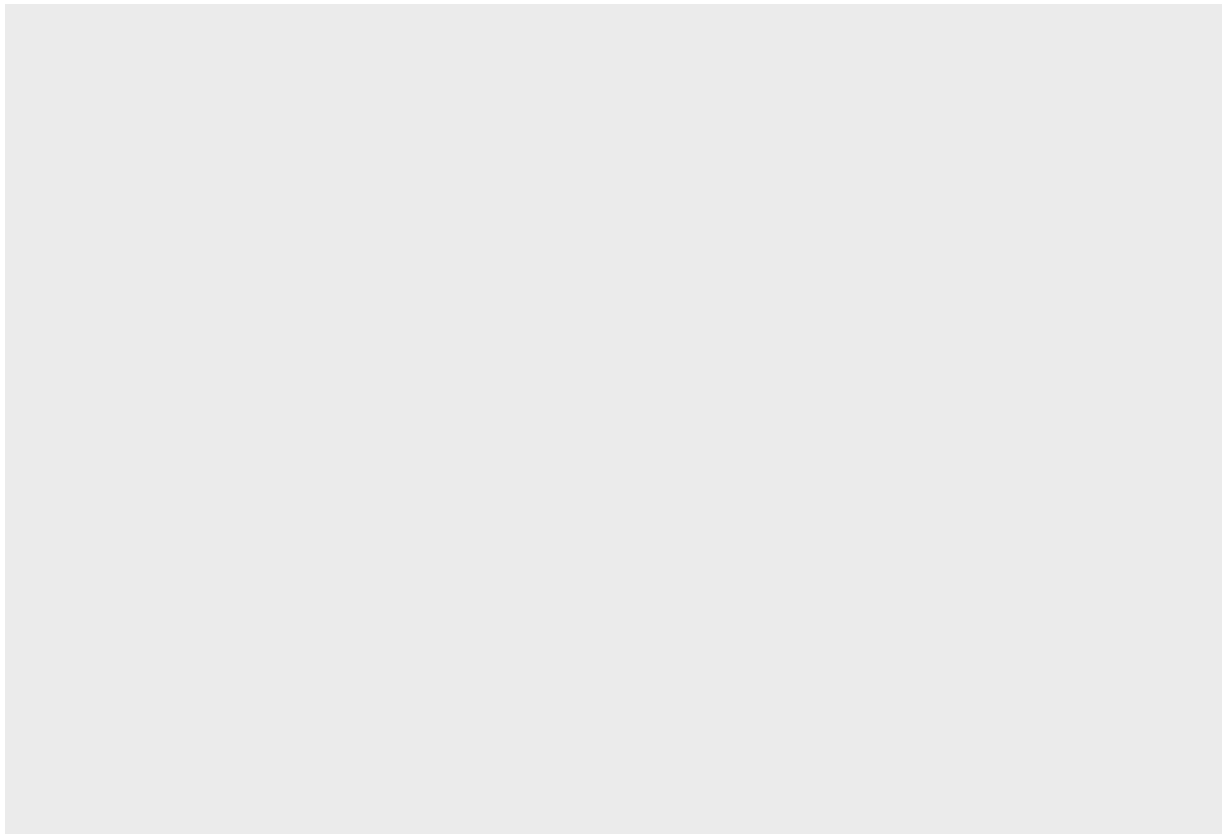
```
– dplyr::tbl(con, "table_name")
```

For efficient data handling, packages like `data.table`, `readr`, and `vroom` provide optimized functions for reading large datasets.

3 Numerical Data

`ggplot2` is a package in R that can create good looking figures in R. It works like putting a layer on another layer, so we use “+” to concatenate layers together.

```
library(ggplot2)
ggplot()
```

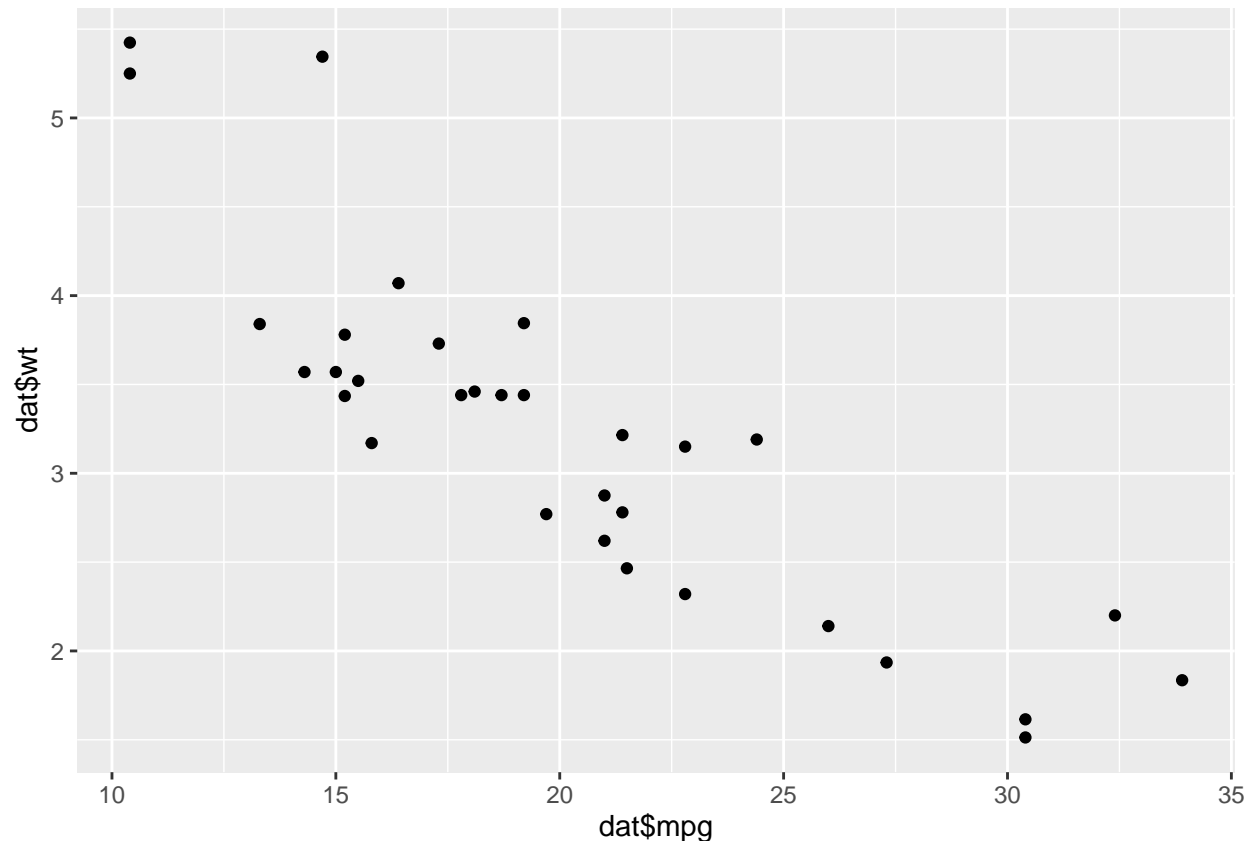


If we only run a `ggplot()`, it returns to a blank figure to us, since we only set up a background, no future plots are made. For the next steps, we will use the in-built dataset “mtcars” as examples.

3.1 Scatter Plots

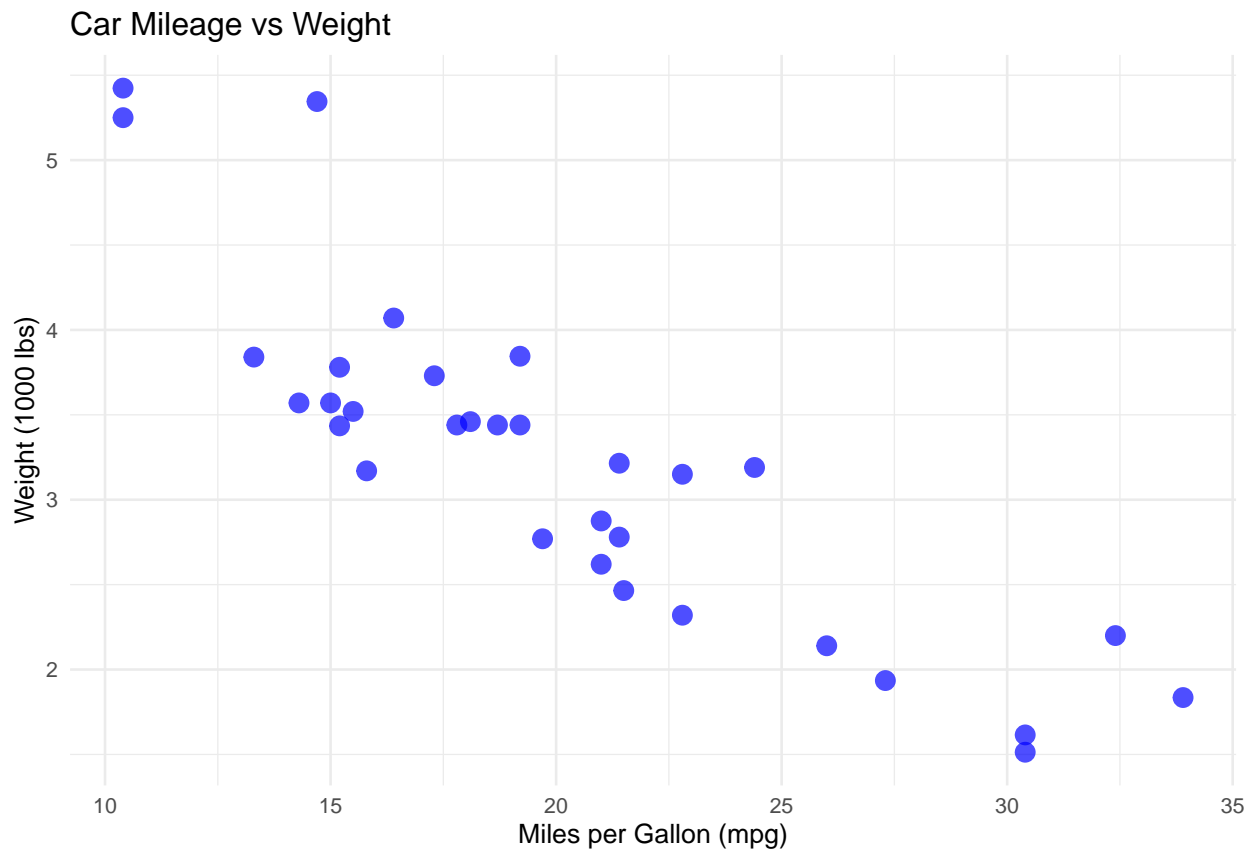
Scatter plots is one of the most commonly used plots in statistics. It visualizes the relationship between two variables. We can distinguish whether they are positively or negatively correlated based on scatter plots. It's also useful when we compare the values of two variables.

```
dat <- mtcars
ggplot() +
  geom_point(aes(x = dat$mpg, y = dat$wt))
```



Based on the blank figure, we added one more layer called `geom_point`. This returns to a scatter plot, with `aes(x = ..., y = ...)` being the x-axis values and y-axis values. Based on this, we can do something to make the plot more readable:

```
ggplot() +
  geom_point(data = mtcars, aes(x = mpg, y = wt), color = "blue", size = 3, alpha = 0.7) + #
  labs(
    title = "Car Mileage vs Weight",
    x = "Miles per Gallon (mpg)",
    y = "Weight (1000 lbs)"
  ) +
  theme_minimal(base_size = 10) # Using a clean theme
```



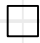

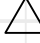



















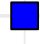



In the `labs()`, we can set the legend title, figure title, and xy axis labels. In the `geom_point` we can set the size of the point, shape of the point, and the color of the point. There are 25 different styles of points in R:

```
library(ggplot2)

# Create a dataframe with shape IDs
point_shapes <- data.frame(
  x = rep(1:6, length.out = 26), # Adjust x positions for 26 points
  y = rep(5:1, each = 6, length.out = 26), # Adjust y positions for 26 points
  shape = 0:25 # Shape IDs (0 to 25)
)

# Create the plot
ggplot(point_shapes, aes(x = x, y = y)) +
  geom_point(aes(shape = shape), size = 5, fill = "blue") + # Shape varies, filled for 21-25
  scale_shape_identity() + # Use shape IDs directly
  geom_text(aes(label = shape), nudge_y = -0.3, size = 5) + # Label each point with its shape
  labs(title = "Point shapes available in R") +
  theme_minimal(base_size = 14) +
  theme(axis.title = element_blank(), axis.text = element_blank(), axis.ticks = element_blank())
```

Point shapes available in R

					
0	1	2	3	4	5
					
6	7	8	9	10	11
					
12	13	14	15	16	17
					
18	19	20	21	22	23
					
24	25				

3.2 Line Chart / Time Series Plot

Sometimes, if the data is a time series:

```
library(ggplot2)
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(tidyr)
library(lubridate) # Ensure proper date handling

##
## Attaching package: 'lubridate'
```

```
## The following objects are masked from 'package:base':
##
##     date, intersect, setdiff, union
```

```
# Set English locale for months
# Sys.setlocale("LC_TIME", "C")

# Simulated time series data
set.seed(123)
time_series_data <- data.frame(
  date = seq(as.Date("2023-01-01"), by = "month", length.out = 12),
  series_A = cumsum(rnorm(12, mean = 5, sd = 2)),
  series_B = cumsum(rnorm(12, mean = 3, sd = 1.5)),
  series_C = cumsum(rnorm(12, mean = 4, sd = 1.8))
)
time_series_data
```

```
##           date  series_A  series_B  series_C
## 1 2023-01-01  3.879049  3.601157  2.874929
## 2 2023-02-01  8.418694  6.767181  3.838881
## 3 2023-03-01 16.536110  8.933420  9.346898
## 4 2023-04-01 21.677127 14.613789 13.622970
## 5 2023-05-01 26.935703 18.360565 15.574323
## 6 2023-06-01 35.365833 18.410639 21.831190
## 7 2023-07-01 41.287665 22.462673 26.598826
## 8 2023-08-01 43.757543 24.753486 30.067697
## 9 2023-09-01 47.383837 26.151750 35.678923
## 10 2023-10-01 51.492513 28.824788 41.259563
## 11 2023-11-01 58.940676 30.285781 46.738409
## 12 2023-12-01 64.660304 32.192445 51.977962
```

The data includes three time series. So we can use a line chart to plot them together. In this format of data, we call them wide format since they are combined together with a same date at each row. However, ggplot cannot recognize this type of data, we need to transform them to a long format.

```
# Transform data into long format for ggplot
long_data <- time_series_data %>%
  pivot_longer(cols = -date, names_to = "series", values_to = "value")
long_data
```

```
## # A tibble: 36 x 3
##   date      series  value
##   <date>    <chr>   <dbl>
## 1 2023-01-01 series_A  3.88
## 2 2023-01-01 series_B  3.60
## 3 2023-01-01 series_C  2.87
```

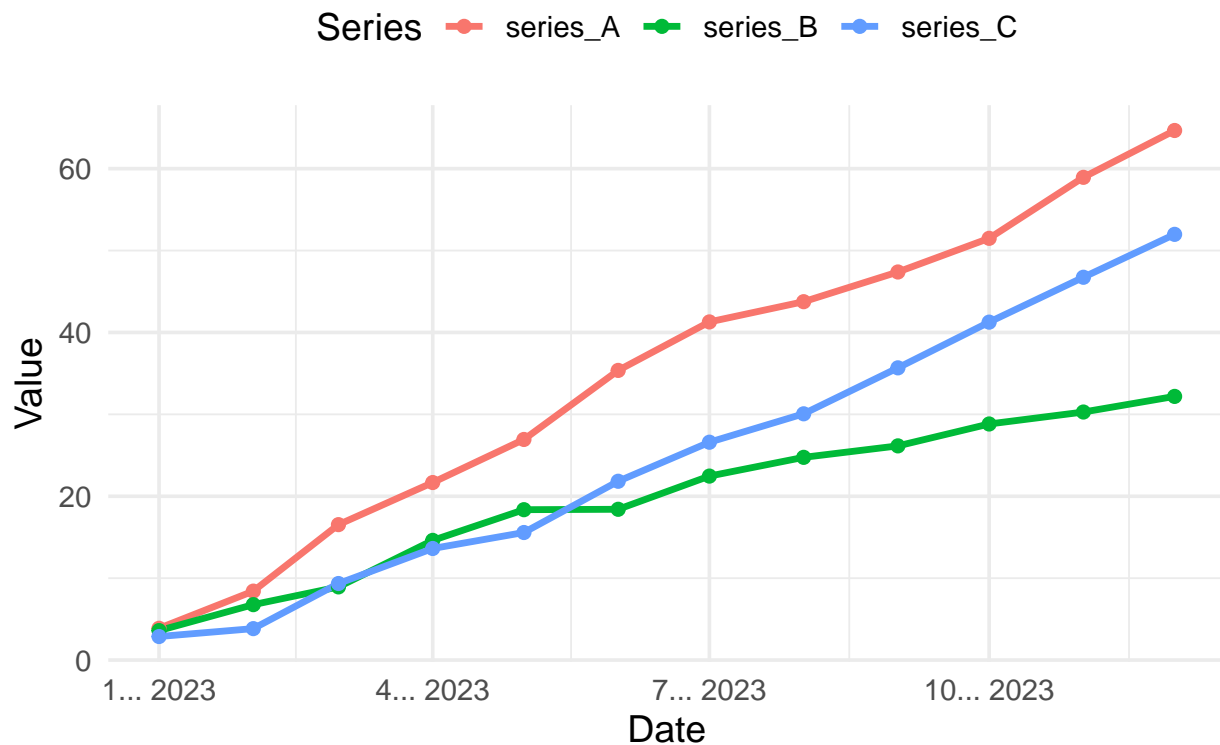


```
## 4 2023-02-01 series_A 8.42
## 5 2023-02-01 series_B 6.77
## 6 2023-02-01 series_C 3.84
## 7 2023-03-01 series_A 16.5
## 8 2023-03-01 series_B 8.93
## 9 2023-03-01 series_C 9.35
## 10 2023-04-01 series_A 21.7
## # i 26 more rows
```

We add one more column called “series”, it’s kind of an index of which series this observation is from. Based on this:

```
# Plot with English x-axis labels
ggplot(long_data, aes(x = date, y = value, color = series)) +
  geom_line(size = 1.2) +
  geom_point(size = 2) +
  scale_x_date(date_labels = "%b %Y") +
  labs(
    title = "Multiple Time Series Plot",
    x = "Date",
    y = "Value",
    color = "Series"
  ) +
  theme_minimal(base_size = 14) +
  theme(legend.position = "top")
```

Multiple Time Series Plot



This example adds the data directly to the `ggplot()` function, that means for the following `geom_` functions, by default we are using the same data, same mapping, same color settings. In the `aes()` function, x axis is the date, y axis is the value, and we set the color of the line/point is the series variable in the long format data. In the `theme()` function, we can adjust a lot of personalized preferences. Refer to: <https://www.rdocumentation.org/packages/ggplot2/versions/3.5.0/topics/theme>.

Time series plots can visualize the trends, like stock market and population.

3.3 Histogram and Density

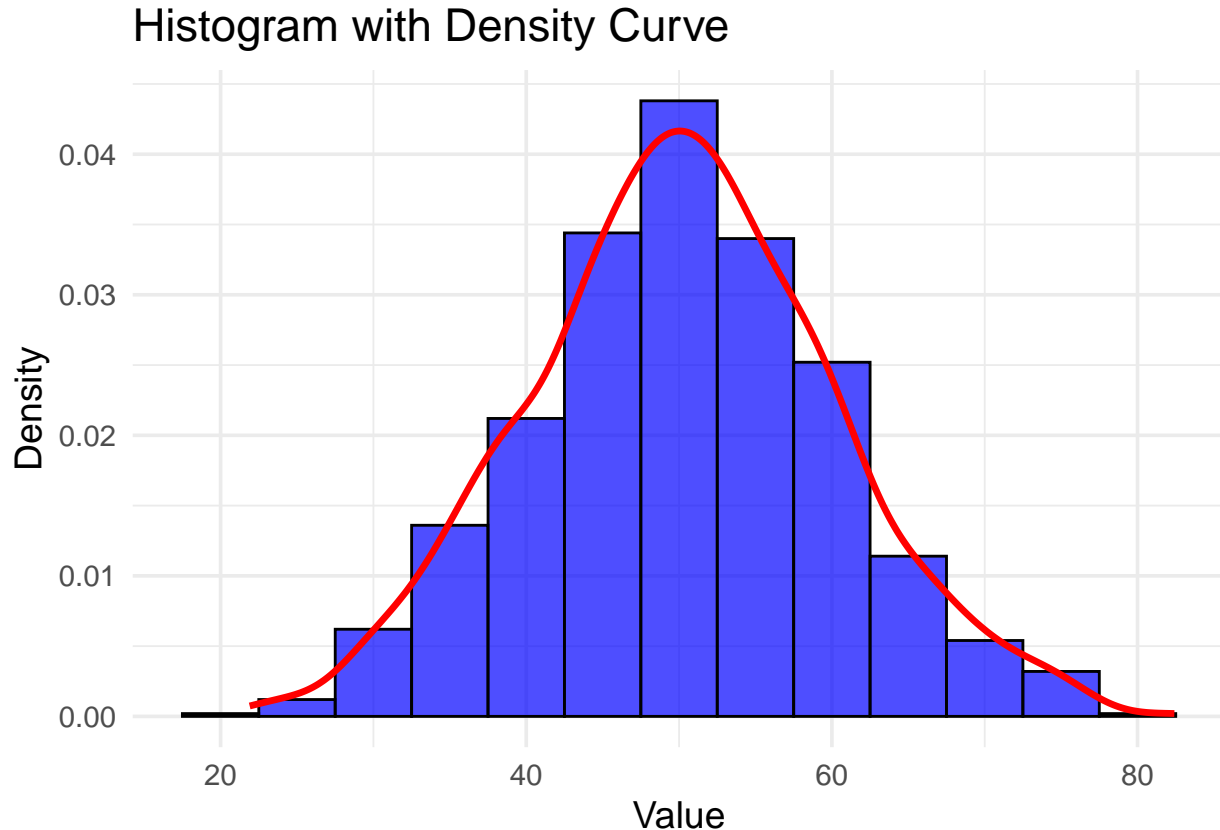
To visualize a distribution, or how the data are distributed, we can use a histogram in `ggplot2`. Below, we are using normal distribution samples as an example:

```
library(ggplot2)

# Generate random samples from a normal distribution
set.seed(123) # For reproducibility
data <- data.frame(value = rnorm(1000, mean = 50, sd = 10))

# Create histogram with a density line
ggplot(data, aes(x = value)) +
  geom_histogram(aes(y = ..density..), binwidth = 5, fill = "blue", color = "black", alpha = 0.5) +
  geom_density(color = "red", size = 1.2) + # Add a smooth density line
  labs(
    title = "Histogram with Density Curve",
```

```
x = "Value",  
y = "Density"  
) +  
theme_minimal(base_size = 14)
```



We first drew a blue histogram, and then set a density curve above it.

3.4 Correlation Matrix Plot

Correlation is a statistical measure that describes the strength and direction of a relationship between two numeric variables. The most common measure is the Pearson correlation coefficient, which ranges from -1 to 1:

- A value of **1** implies a perfect positive linear relationship
- A value of **0** implies no linear relationship
- A value of **-1** implies a perfect negative linear relationship

Visualizing correlation can help us better understand relationships between variables in a dataset. In this tutorial, we'll use R and the `ggplot2` package to create a correlation plot.

```
library(ggplot2)  
data(mtcars)  
head(mtcars)
```

```
##           mpg cyl disp  hp drat   wt  qsec vs am gear carb
## Mazda RX4      21.0   6  160 110 3.90 2.620 16.46 0  1   4   4
## Mazda RX4 Wag  21.0   6  160 110 3.90 2.875 17.02 0  1   4   4
## Datsun 710     22.8   4  108  93 3.85 2.320 18.61 1  1   4   1
## Hornet 4 Drive  21.4   6  258 110 3.08 3.215 19.44 1  0   3   1
## Hornet Sportabout 18.7   8  360 175 3.15 3.440 17.02 0  0   3   2
## Valiant        18.1   6  225 105 2.76 3.460 20.22 1  0   3   1
```

3.4.1 Step 1: Compute Correlation Matrix

We first calculate the correlation matrix using `cor()`.

```
library(ggcorrplot)
```

```
## Warning: package 'ggcorrplot' was built under R version 4.3.3
```

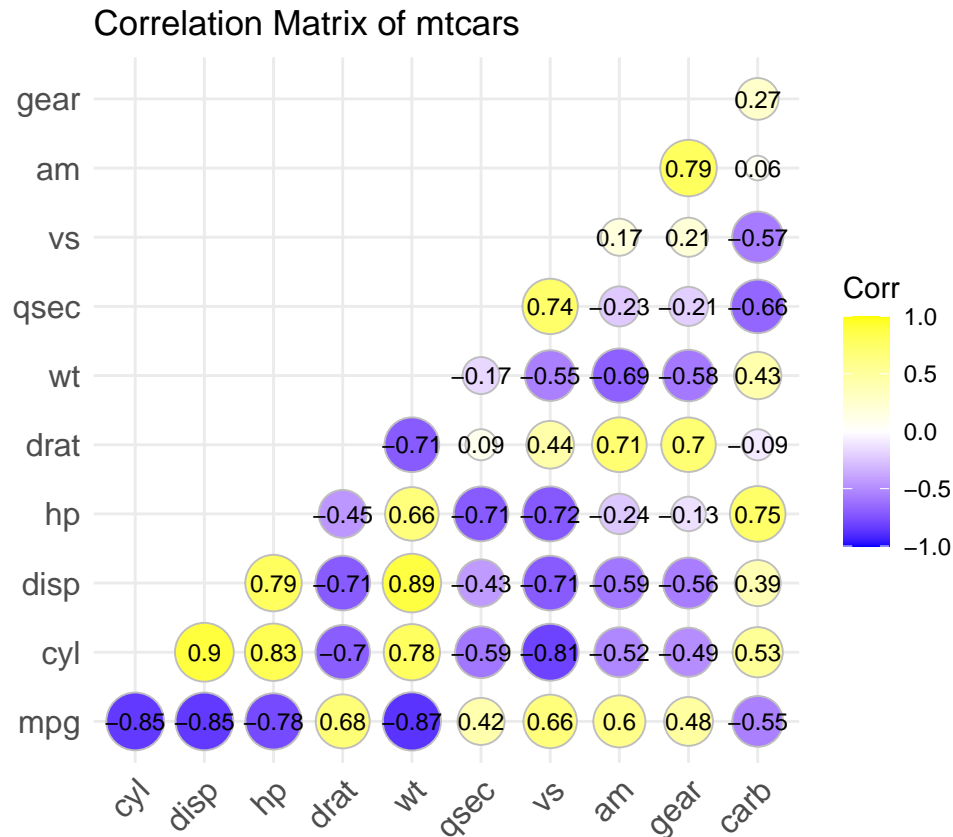
```
corr_matrix <- round(cor(mtcars), 2)
corr_matrix
```

```
##           mpg   cyl  disp    hp  drat    wt  qsec    vs   am gear carb
## mpg      1.00 -0.85 -0.85 -0.78  0.68 -0.87  0.42  0.66  0.60  0.48 -0.55
## cyl     -0.85  1.00  0.90  0.83 -0.70  0.78 -0.59 -0.81 -0.52 -0.49  0.53
## disp    -0.85  0.90  1.00  0.79 -0.71  0.89 -0.43 -0.71 -0.59 -0.56  0.39
## hp       -0.78  0.83  0.79  1.00 -0.45  0.66 -0.71 -0.72 -0.24 -0.13  0.75
## drat     0.68 -0.70 -0.71 -0.45  1.00 -0.71  0.09  0.44  0.71  0.70 -0.09
## wt      -0.87  0.78  0.89  0.66 -0.71  1.00 -0.17 -0.55 -0.69 -0.58  0.43
## qsec     0.42 -0.59 -0.43 -0.71  0.09 -0.17  1.00  0.74 -0.23 -0.21 -0.66
## vs       0.66 -0.81 -0.71 -0.72  0.44 -0.55  0.74  1.00  0.17  0.21 -0.57
## am       0.60 -0.52 -0.59 -0.24  0.71 -0.69 -0.23  0.17  1.00  0.79  0.06
## gear     0.48 -0.49 -0.56 -0.13  0.70 -0.58 -0.21  0.21  0.79  1.00  0.27
## carb    -0.55  0.53  0.39  0.75 -0.09  0.43 -0.66 -0.57  0.06  0.27  1.00
```

3.4.2 Step 2: Visualize with ggcorrplot

We use the `ggcorrplot` package to make the visualization more intuitive. This function creates a heatmap where the color and size of the tiles represent the correlation strength.

```
ggcorrplot(corr_matrix,
  method = "circle",
  type = "lower",
  lab = TRUE,
  lab_size = 3,
  colors = c("blue", "white", "yellow"),
  title = "Correlation Matrix of mtcars",
  ggtheme = theme_minimal())
```



the Plot

In the correlation plot:

- **Darker blue circles** indicate strong positive correlation (e.g., weight and horsepower)
- **Darker red circles** indicate strong negative correlation (e.g., mpg and weight)
- **Smaller or white circles** indicate weak or no correlation

This plot helps you identify which variables are strongly related and may influence one another, which is useful in regression modeling, feature selection, and exploratory data analysis.

3.4.3 Conclusion

Correlation plots provide a quick and effective way to explore relationships in your data. With `ggcorrplot` and `ggplot2`, it's easy to produce publication-quality visualizations to support your analysis.

4 Categorical Data

4.1 Grouping

We simulate two bivariate normal distributions using `MASS::mvrnorm()` and visualize them using `ggplot2`. Each group is colored differently to show separation in the 2D space.

```
# Load required packages
```

```
library(MASS)
```

```
##
```

```
## Attaching package: 'MASS'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      select
```

```
library(ggplot2)
```

```
# Set seed for reproducibility
```

```
set.seed(123)
```

```
# Parameters for Group 1
```

```
mu1 <- c(2, 3)
```

```
sigma1 <- matrix(c(1, 0.5, 0.5, 1), nrow = 2)
```

```
group1 <- mvrnorm(n = 100, mu = mu1, Sigma = sigma1)
```

```
# Parameters for Group 2
```

```
mu2 <- c(6, 7)
```

```
sigma2 <- matrix(c(1, -0.3, -0.3, 1), nrow = 2)
```

```
group2 <- mvrnorm(n = 100, mu = mu2, Sigma = sigma2)
```

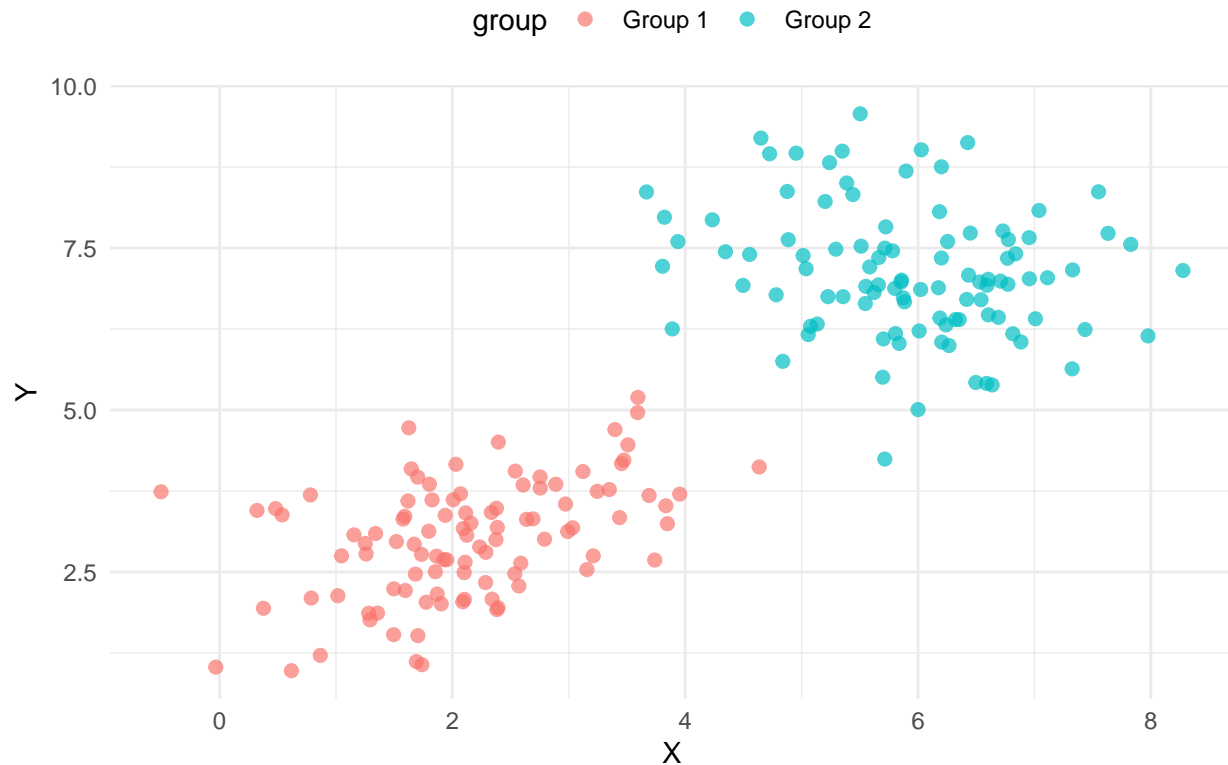
```
# Combine into one data frame
```

```
df <- data.frame(  
  x = c(group1[,1], group2[,1]),  
  y = c(group1[,2], group2[,2]),  
  group = factor(c(rep("Group 1", 100), rep("Group 2", 100)))  
)
```

```
# Plot
```

```
ggplot(df, aes(x = x, y = y, color = group)) +  
  geom_point(alpha = 0.7, size = 2) +  
  labs(title = "Simulated Bivariate Normals",  
        x = "X", y = "Y") +  
  theme_minimal() +  
  theme(legend.position = "top")
```

Simulated Bivariate Normals



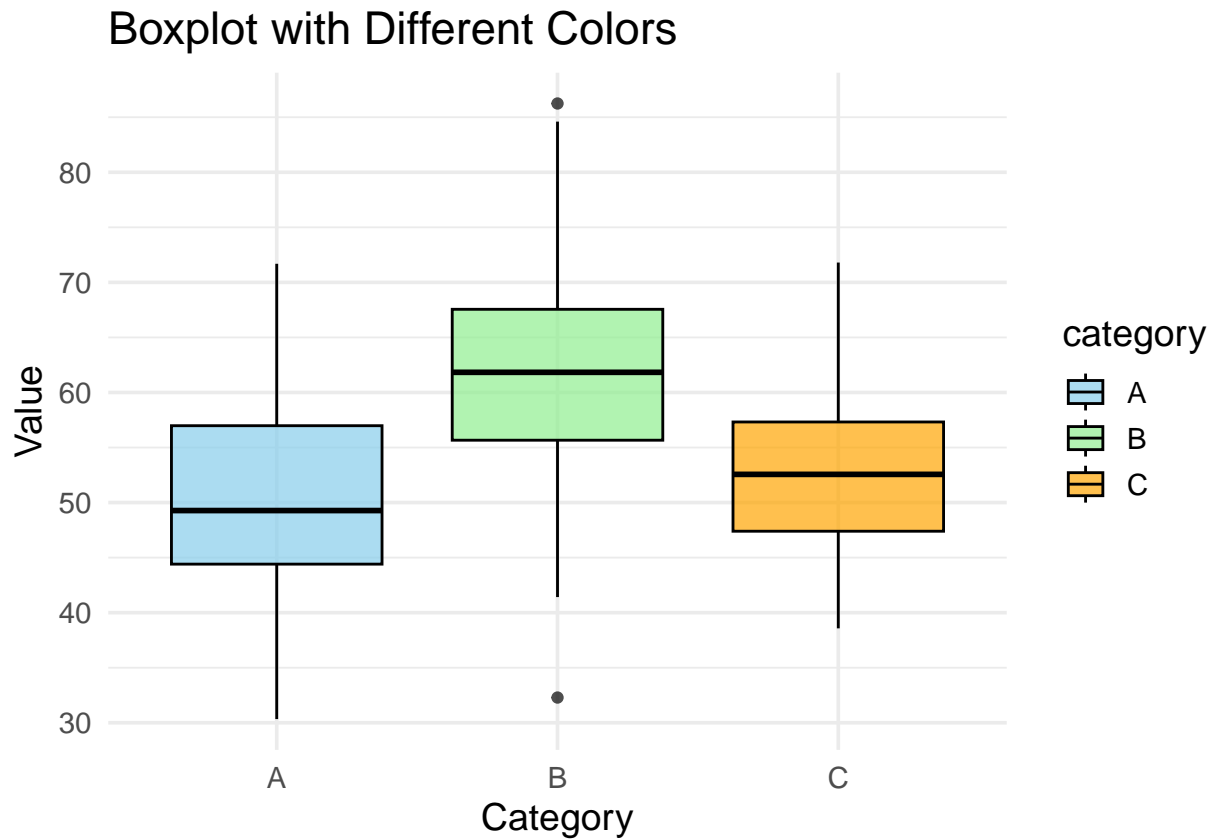
4.2 Boxplot

Boxplots are used when we have multiple groups to compare. We want to compare the distribution of them, for example:

```
library(ggplot2)

# Create sample data
set.seed(123)
data <- data.frame(
  category = rep(c("A", "B", "C"), each = 50),
  value = c(rnorm(50, mean = 50, sd = 10),
            rnorm(50, mean = 60, sd = 12),
            rnorm(50, mean = 55, sd = 8))
)

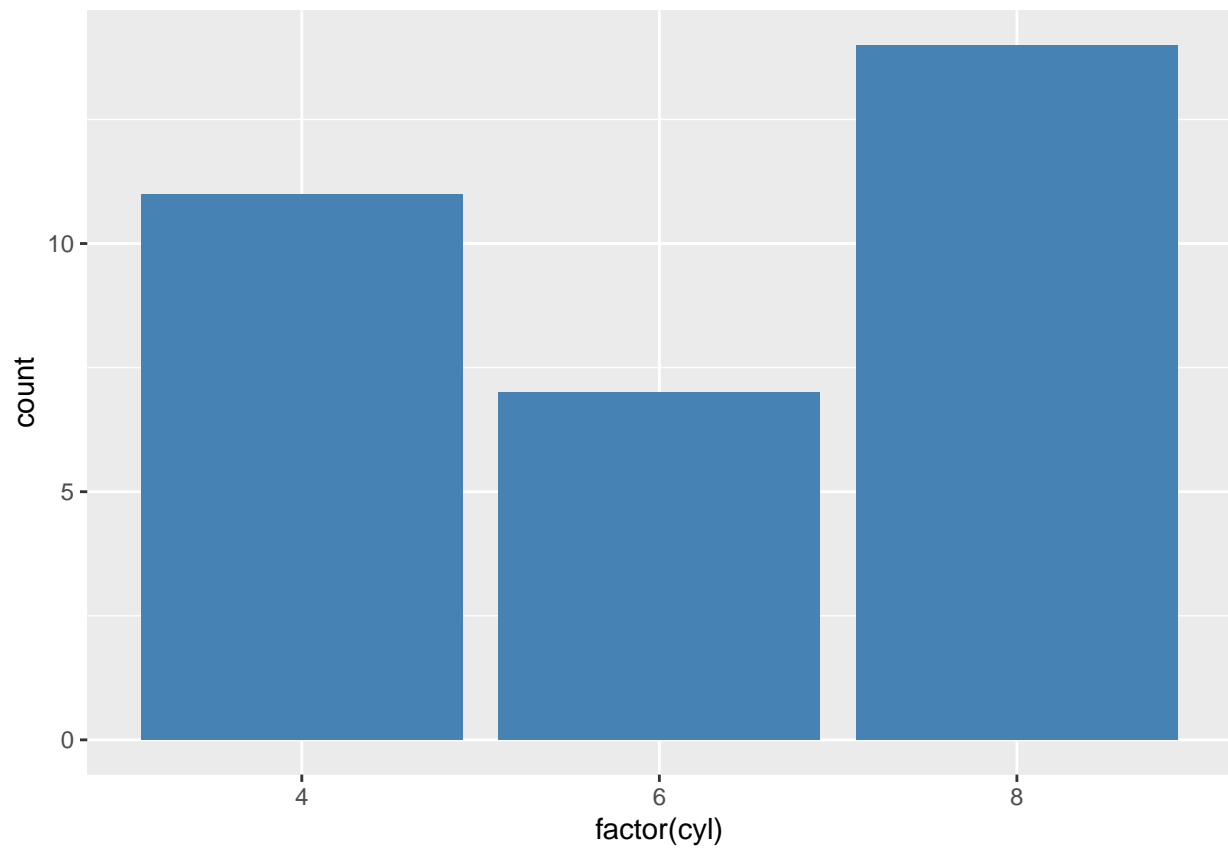
# Boxplot with different colors for each category
ggplot(data, aes(x = category, y = value, fill = category)) + # Color by category
  geom_boxplot(color = "black", alpha = 0.7) + # Add black borders
  scale_fill_manual(values = c("A" = "skyblue", "B" = "lightgreen", "C" = "orange")) + # Custom colors
  labs(title = "Boxplot with Different Colors", x = "Category", y = "Value") +
  theme_minimal(base_size = 14)
```



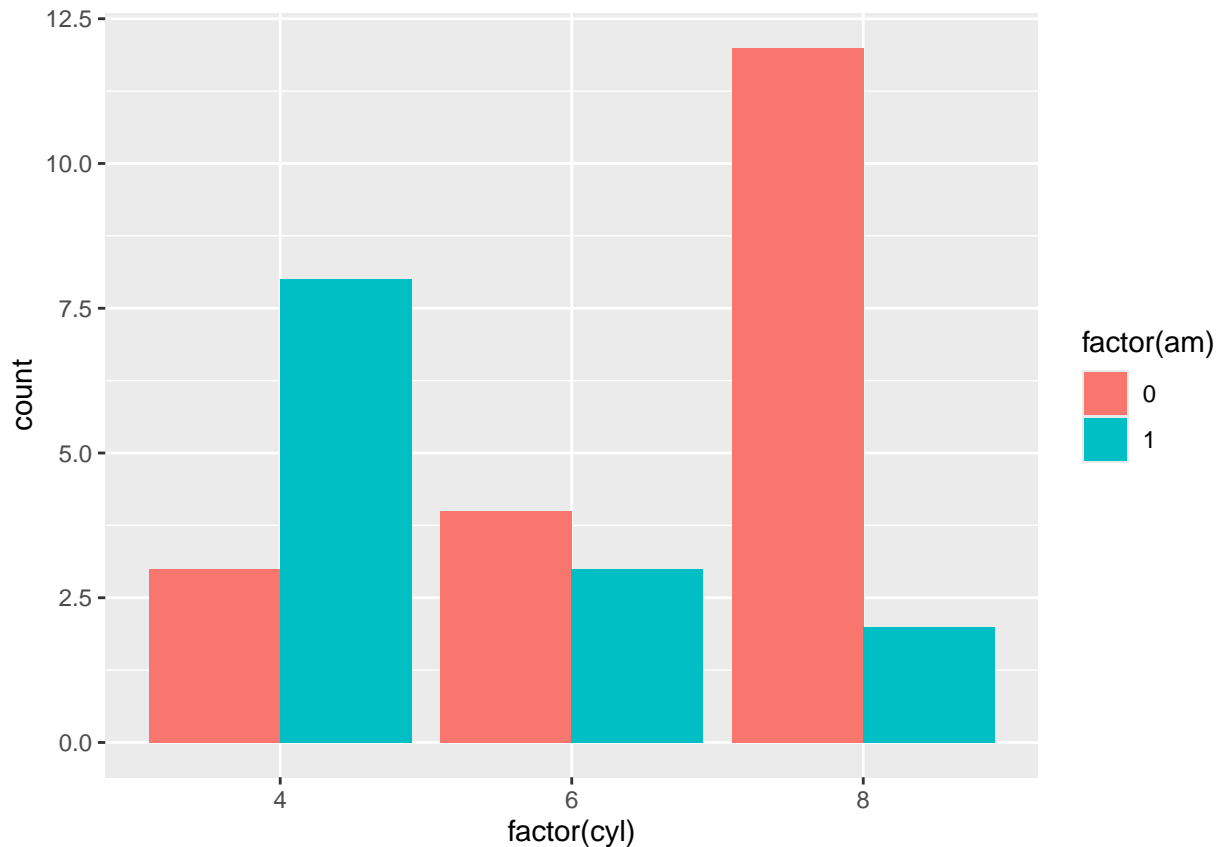
The top edge of the box describes the third quartile of the data, the thick black line in the middle describes the median of the data, and the bottom edge of the box is the first quartile. From the plot, we can see the group B has an overall higher value than the group A and C.

4.3 Bar plot

```
ggplot(mtcars, aes(x = factor(cyl))) +  
  geom_bar(fill = "steelblue")
```

```
ggplot(mtcars, aes(x = factor(cyl), fill = factor(am))) +  
  geom_bar(position = "dodge")
```



5 Special Cases

5.1 Heatmap and Contour Plot

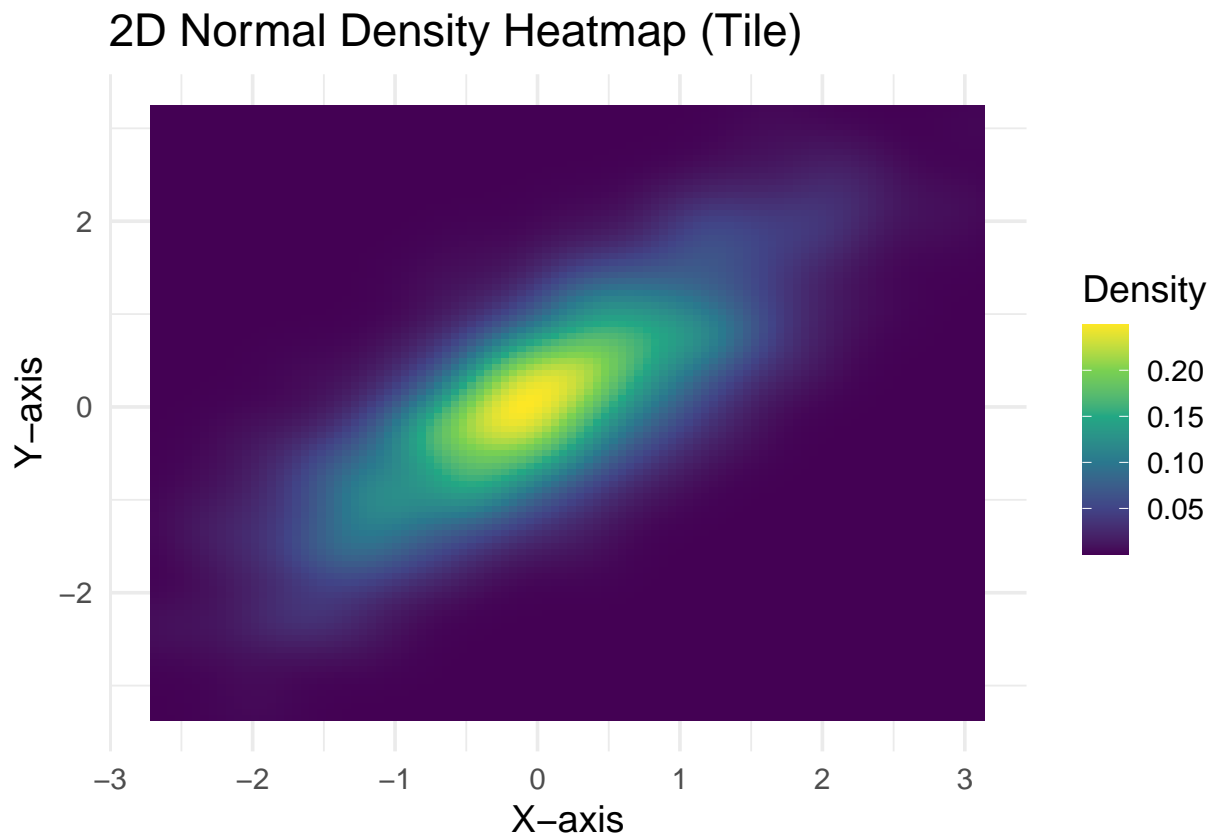
Some spatial dataset, has a location together with a value. The location is usually given in a 2-d coordinate. In this case, we can use a heatmap or a contour plot.

```
library(ggplot2)
library(MASS) # For multivariate normal distribution

# Generate 2D normal distribution data
set.seed(123)
mu <- c(0, 0) # Mean vector
sigma <- matrix(c(1, 0.8, 0.8, 1), ncol = 2) # Covariance matrix
samples <- mvrnorm(n = 1000, mu = mu, Sigma = sigma) # Generate samples
data <- data.frame(x = samples[,1], y = samples[,2]) # Convert to dataframe

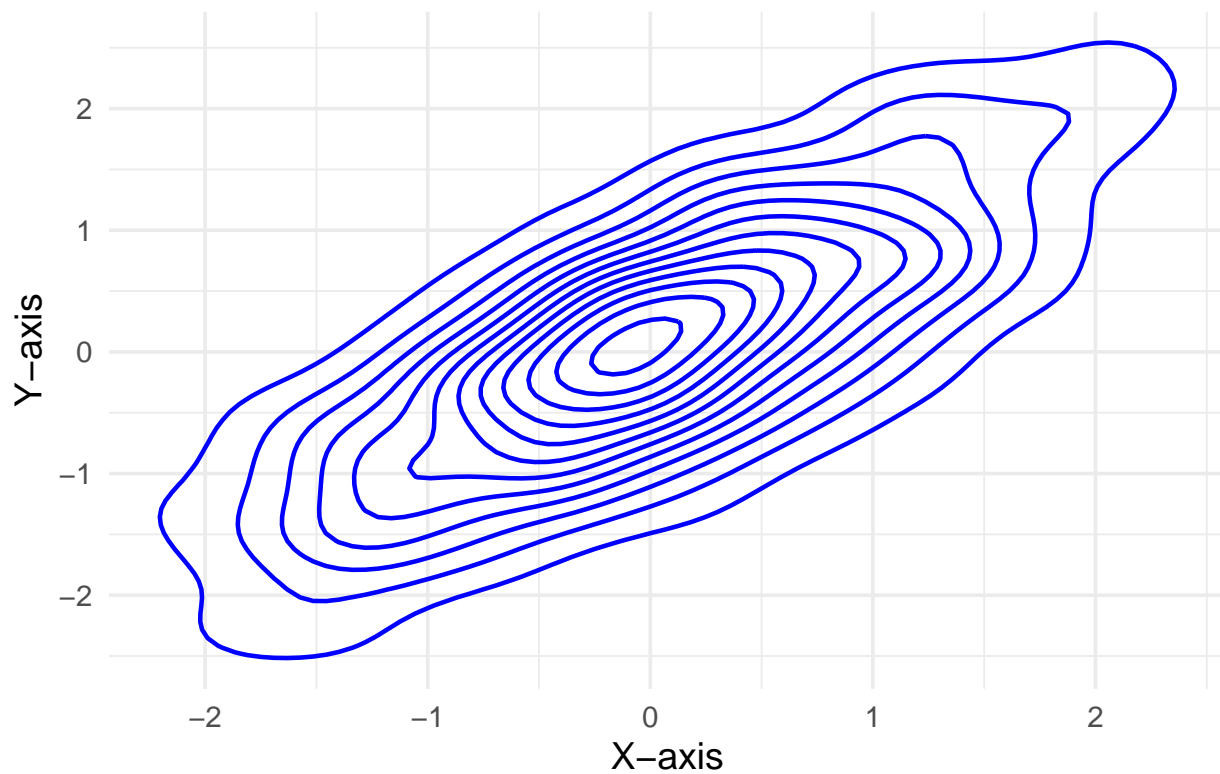
# Compute 2D density
density_data <- with(data, MASS::kde2d(x, y, n = 100)) # Estimate density
density_df <- data.frame(expand.grid(x = density_data$x, y = density_data$y),
                          density = as.vector(density_data$z)) # Convert to dataframe
```

```
# 2D heatmap with geom_tile()
ggplot(density_df, aes(x = x, y = y, fill = density)) +
  geom_tile() + # Heatmap using tiles
  scale_fill_viridis_c() + # Color scale
  labs(title = "2D Normal Density Heatmap (Tile)",
        x = "X-axis", y = "Y-axis", fill = "Density") +
  theme_minimal(base_size = 14)
```



```
ggplot(density_df, aes(x = x, y = y, z = density)) +
  geom_contour(color = "blue", size = 0.8) + # Contour lines
  labs(title = "2D Normal Density Contour Plot",
        x = "X-axis", y = "Y-axis") +
  theme_minimal(base_size = 14)
```

2D Normal Density Contour Plot



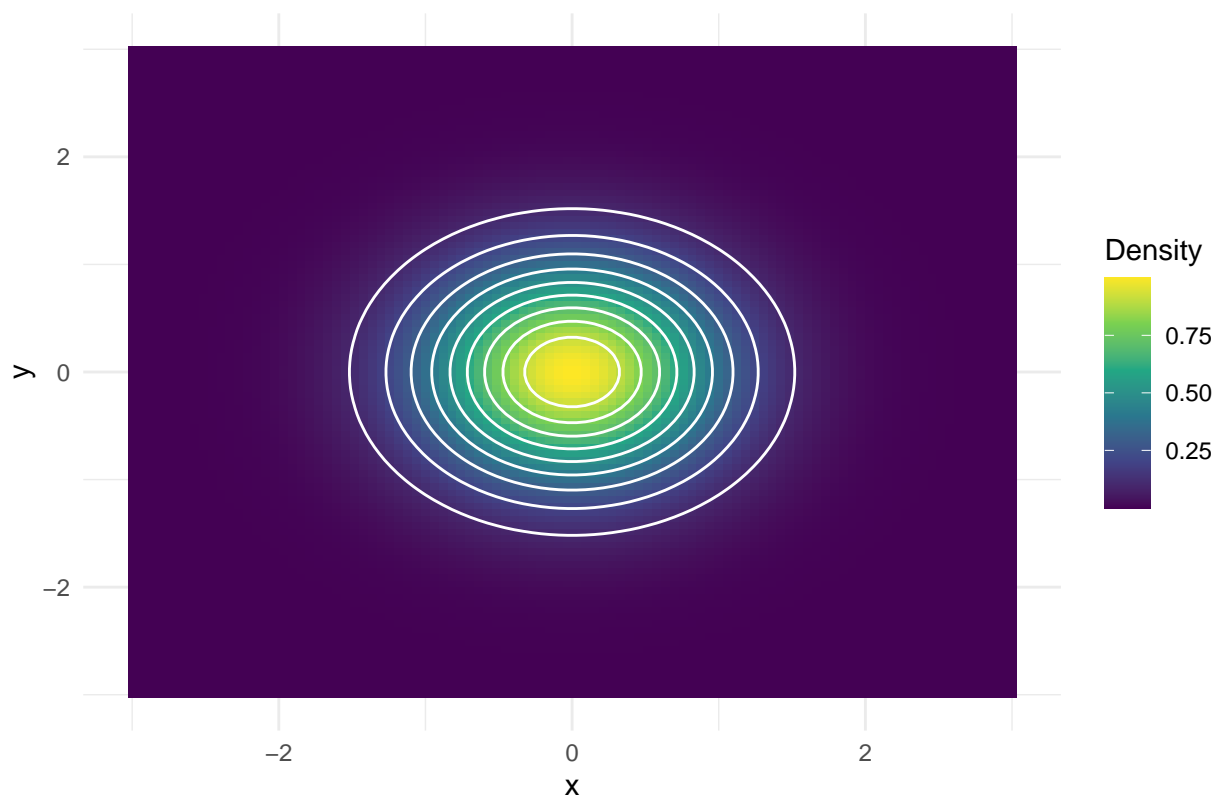
5.2 Simulated 2D Gaussian bump

```
library(ggplot2)

x <- seq(-3, 3, length.out = 100)
y <- seq(-3, 3, length.out = 100)
grid <- expand.grid(x = x, y = y)
grid$z <- with(grid, exp(-x^2 - y^2))

ggplot(grid, aes(x = x, y = y, z = z)) +
  geom_tile(aes(fill = z)) +
  geom_contour(color = "white") +
  scale_fill_viridis_c() +
  theme_minimal() +
  labs(title = "Contour + Heatmap", fill = "Density")
```

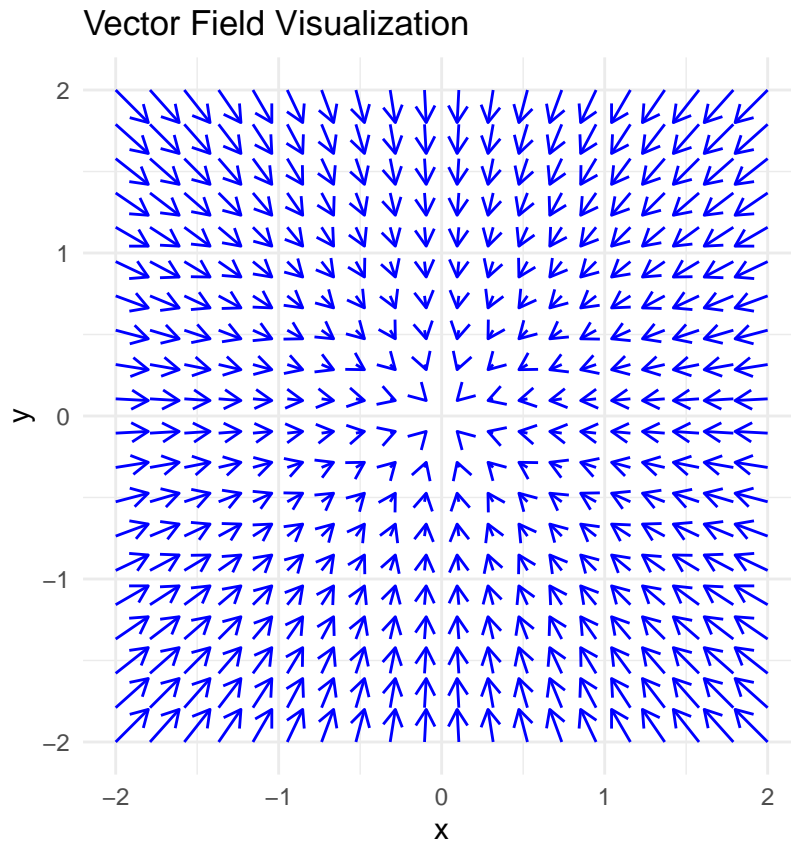
Contour + Heatmap



```
library(ggplot2)

# Simulate a vector field
x <- y <- seq(-2, 2, length.out = 20)
grid <- expand.grid(x = x, y = y)
grid$u <- -grid$x
grid$v <- -grid$y

ggplot(grid, aes(x = x, y = y)) +
  geom_segment(aes(xend = x + u * 0.1, yend = y + v * 0.1),
    arrow = arrow(length = unit(0.1, "inches")),
    color = "blue") +
  coord_fixed() +
  theme_minimal() +
  labs(title = "Vector Field Visualization")
```



6 Personalize

6.1 Customizing ggplot2 Plots

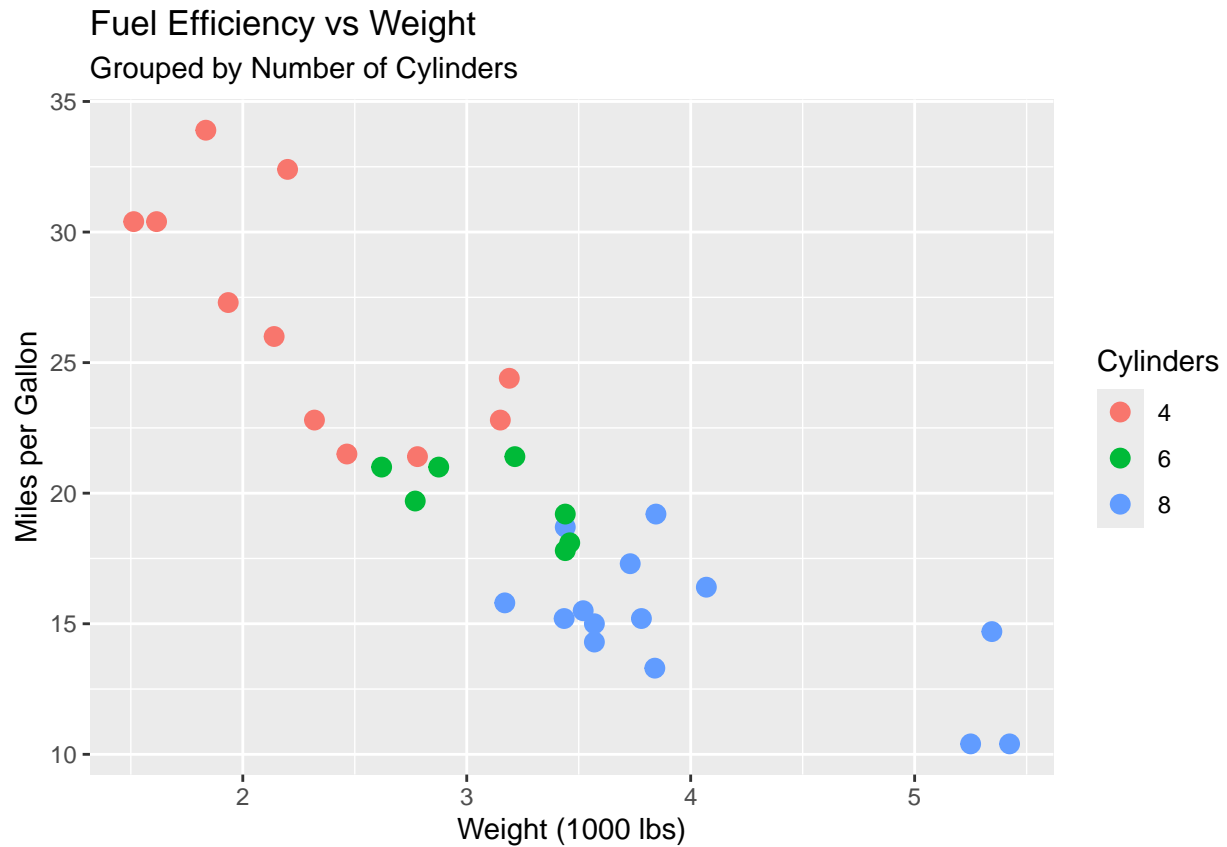
ggplot2 provides powerful tools for customizing the appearance of your plots. Personalization can help highlight important patterns, improve clarity, and make your visuals publication-ready.

6.1.1 Titles, Axis Labels, and Legends

Use `labs()` to change plot title, axis names, and legend title.

```
library(ggplot2)

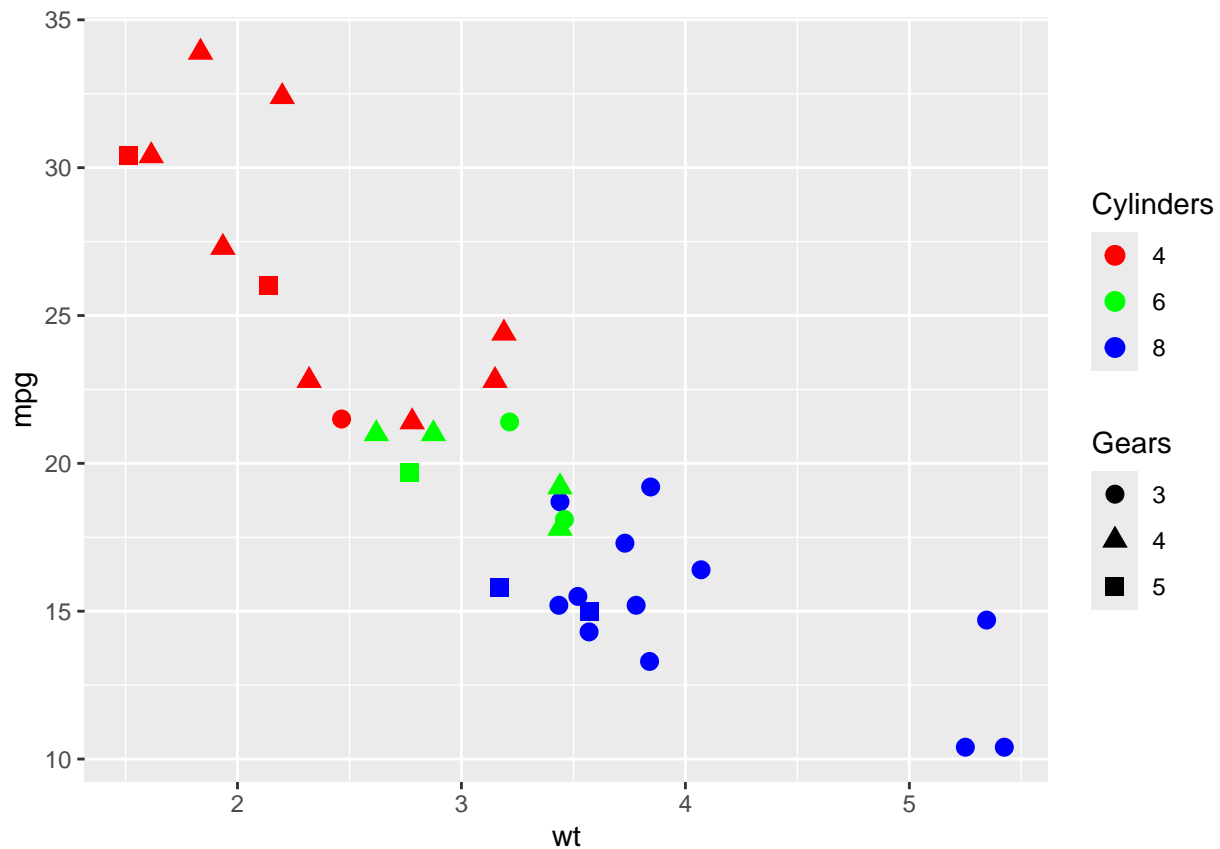
ggplot(mtcars, aes(x = wt, y = mpg, color = factor(cyl))) +
  geom_point(size = 3) +
  labs(title = "Fuel Efficiency vs Weight",
       subtitle = "Grouped by Number of Cylinders",
       x = "Weight (1000 lbs)",
       y = "Miles per Gallon",
       color = "Cylinders")
```



6.1.2 Colors and Shapes

Use `scale_color_manual()` to manually define colors; use `shape =` for custom point types.

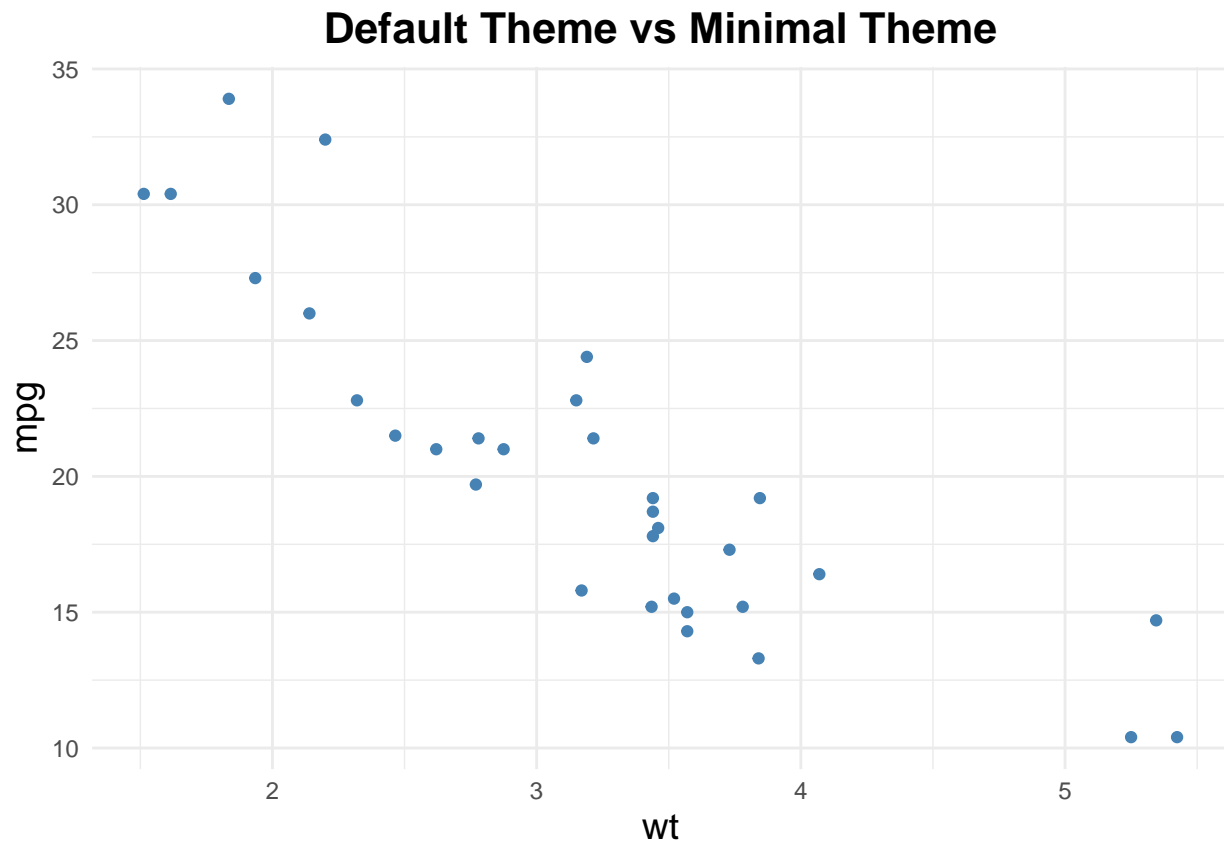
```
ggplot(mtcars, aes(x = wt, y = mpg, color = factor(cyl), shape = factor(gear))) +  
  geom_point(size = 3) +  
  scale_color_manual(values = c("red", "green", "blue")) +  
  labs(color = "Cylinders", shape = "Gears")
```



6.1.3 Themes and Fonts

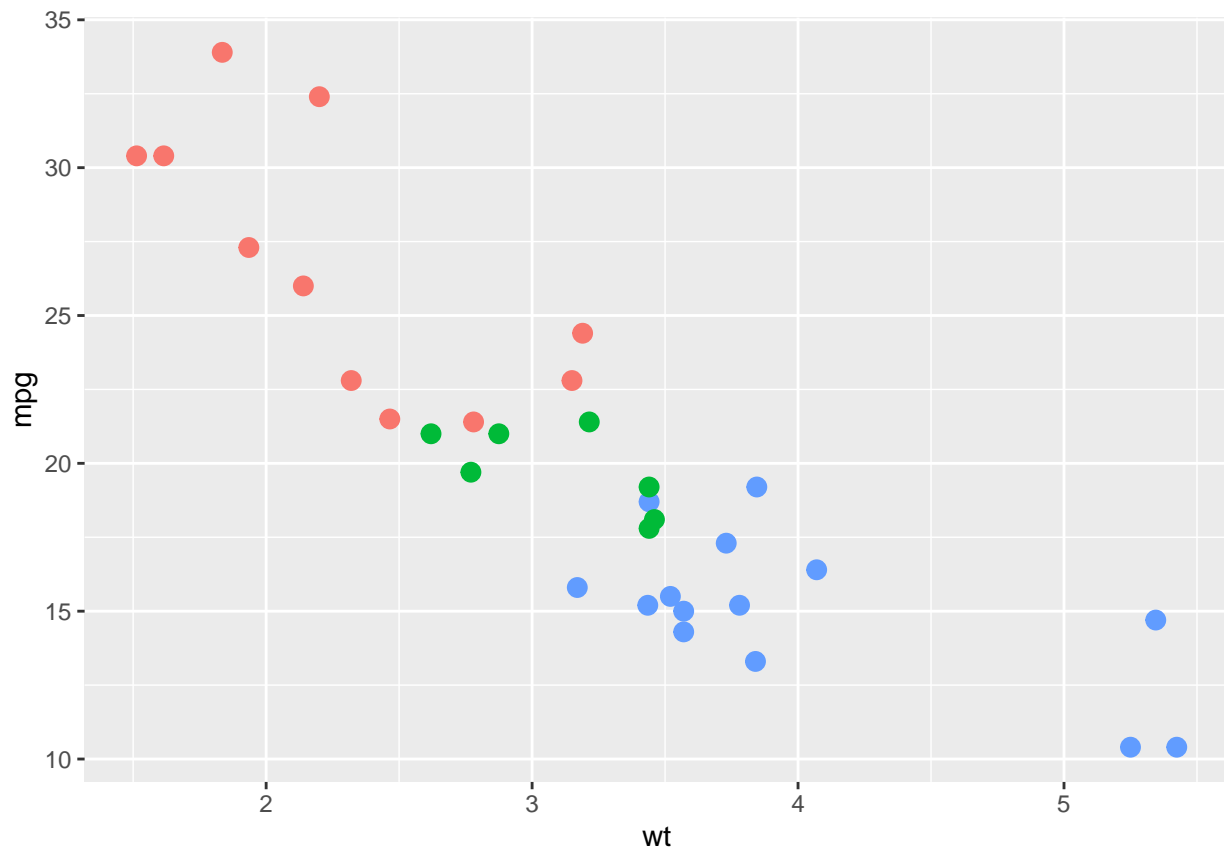
Use built-in themes like `theme_minimal()`, `theme_classic()` or customize with `theme()`.

```
ggplot(mtcars, aes(x = wt, y = mpg)) +
  geom_point(color = "steelblue") +
  labs(title = "Default Theme vs Minimal Theme") +
  theme_minimal() +
  theme(
    plot.title = element_text(hjust = 0.5, size = 16, face = "bold"),
    axis.title = element_text(size = 14),
    legend.position = "bottom"
  )
```

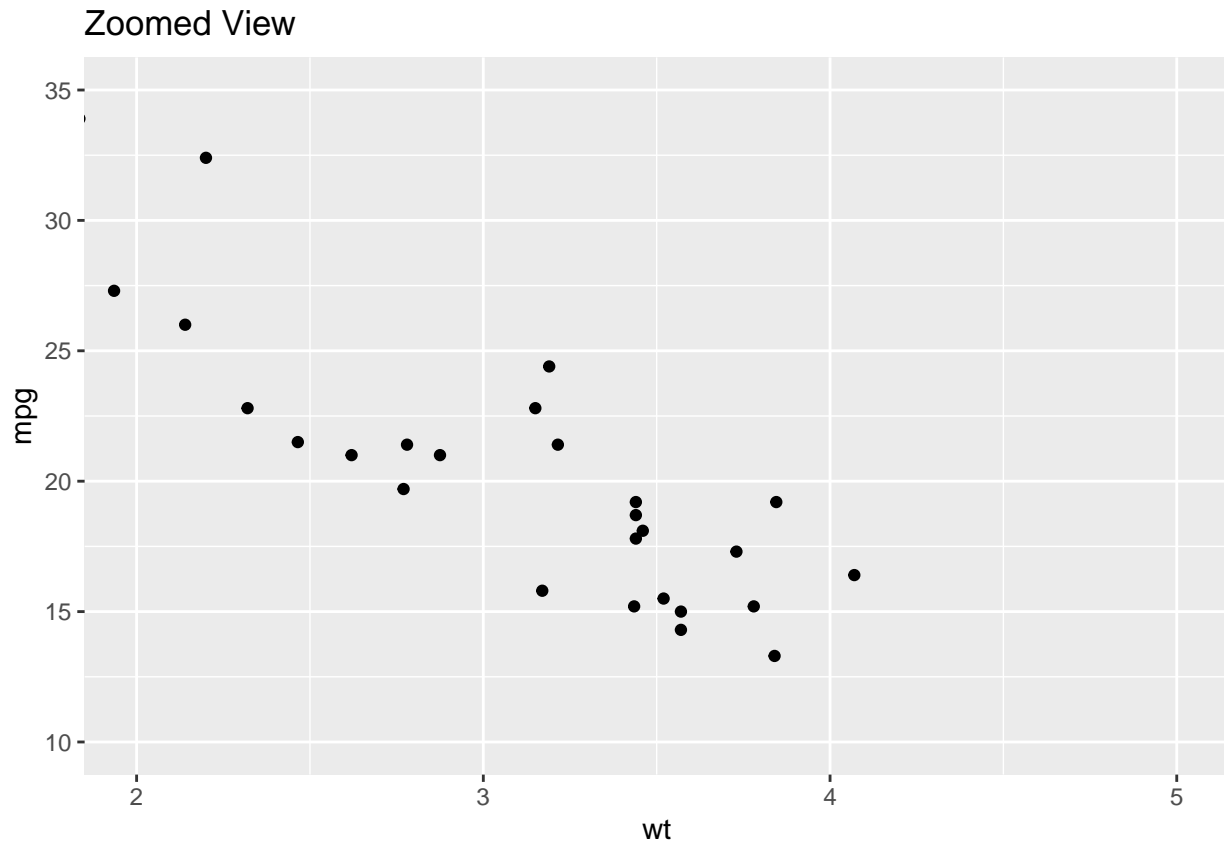
6.1.4 Legend Control

```
ggplot(mtcars, aes(x = wt, y = mpg, color = factor(cyl))) +  
  geom_point(size = 3) +  
  theme(legend.position = "top") # Other options: "bottom", "left", "none"
```

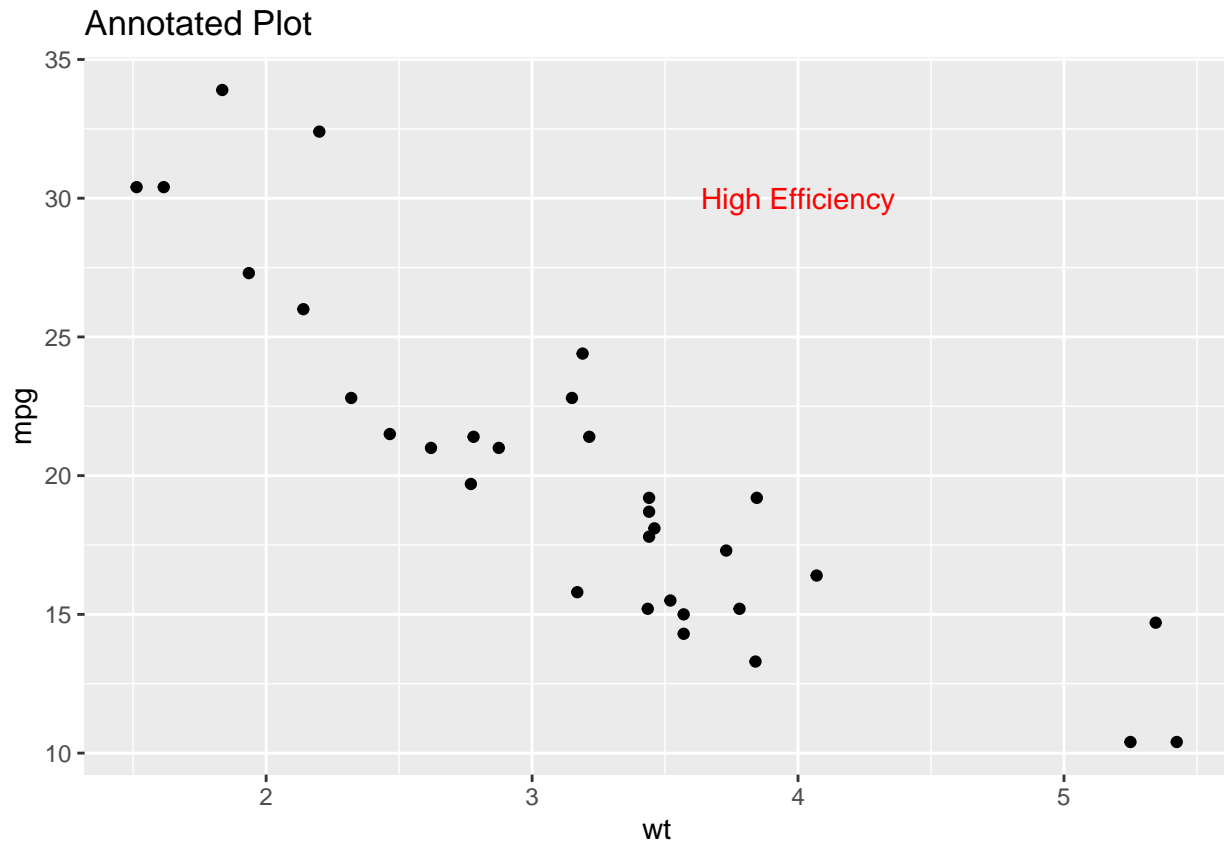
6.1.5 Coordinate Control

```
ggplot(mtcars, aes(x = wt, y = mpg)) +  
  geom_point() +  
  coord_cartesian(xlim = c(2, 5), ylim = c(10, 35)) +  
  labs(title = "Zoomed View")
```



6.1.6 Text and Annotation

```
ggplot(mtcars, aes(x = wt, y = mpg)) +  
  geom_point() +  
  annotate("text", x = 4, y = 30, label = "High Efficiency", color = "red") +  
  labs(title = "Annotated Plot")
```



6.2 Summary

Task	Function(s)
Change title/label	<code>labs()</code> , <code>ggtitle()</code> , <code>xlab()</code>
Manual color palette	<code>scale_color_manual()</code>
Font/size/position	<code>theme()</code> with <code>element_text()</code>
Theme style	<code>theme_minimal()</code> , <code>theme_classic()</code>
Legend control	<code>theme(legend.position = ...)</code>
Axis limits	<code>coord_cartesian()</code>
Add annotation	<code>annotate()</code> , <code>geom_text()</code>

With these tools, your ggplot2 plots can be effectively customized for clarity, impact, and communication.

7 Exercise: Iris Dataset

7.1 Introduction

In this exercise, we will explore the **iris** dataset, which contains measurements of 150 iris flowers from three different species. The goal is to visualize and analyze the relationships between the variables using various types of plots.

The **iris** dataset is a classic and widely used dataset in statistics and machine learning. It consists of 150 observations of iris flowers from three different species: *setosa*, *versicolor*, and *virginica*. Each observation includes four numeric measurements:

- **Sepal.Length**: Length of the sepal (in centimeters)
- **Sepal.Width**: Width of the sepal (in centimeters)
- **Petal.Length**: Length of the petal (in centimeters)
- **Petal.Width**: Width of the petal (in centimeters)

In addition to these measurements, the dataset includes a **Species** column, which indicates the species to which each flower belongs. The dataset is often used to demonstrate classification, clustering, and data visualization techniques.

In this exercise, we will explore various ways to visualize and analyze the relationships between these variables.

7.2 Load the Dataset

```
# Load necessary packages
library(ggplot2)

# Load the iris dataset
data(iris)
head(iris)
```

```
##   Sepal.Length Sepal.Width Petal.Length Petal.Width Species
## 1          5.1         3.5          1.4          0.2   setosa
## 2          4.9         3.0          1.4          0.2   setosa
## 3          4.7         3.2          1.3          0.2   setosa
## 4          4.6         3.1          1.5          0.2   setosa
## 5          5.0         3.6          1.4          0.2   setosa
## 6          5.4         3.9          1.7          0.4   setosa
```

7.3 Exercise 1: Scatter Plot for Relationship Between Two Numerical Variables

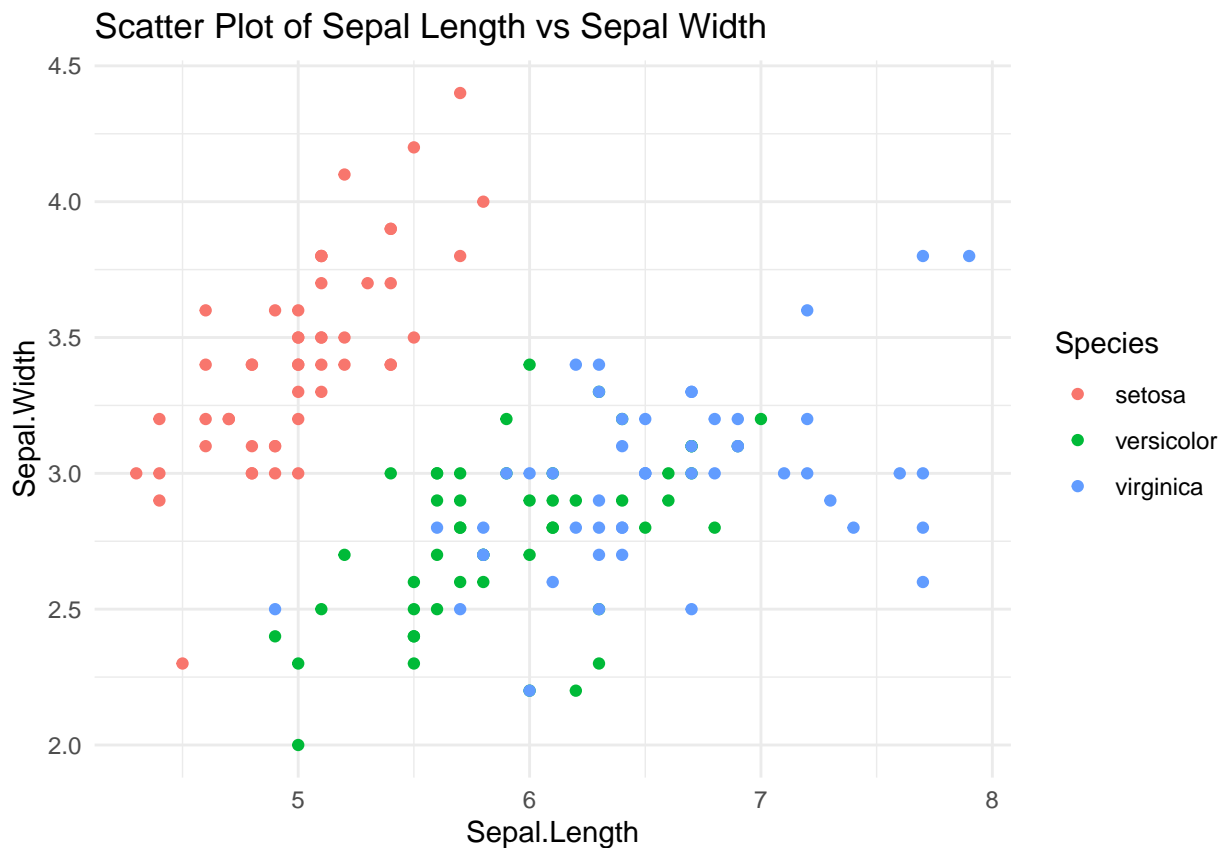
7.3.1 Question:

- You are interested in how **Sepal.Length** relates to **Sepal.Width**. Which type of plot would you use to visualize this relationship, and why?

7.3.2 Solution:

Suppose you want to study the relationship between `Sepal.Length` and `Sepal.Width`. A scatter plot is a good way to visualize how these two variables are related.

```
# Scatter plot of Sepal.Length vs Sepal.Width
ggplot(iris, aes(x = Sepal.Length, y = Sepal.Width)) +
  geom_point(aes(color = Species)) +
  labs(title = "Scatter Plot of Sepal Length vs Sepal Width") +
  theme_minimal()
```



7.4 Exercise 2: Boxplot for Comparison Across Categories

7.4.1 Question:

- If you want to compare the distribution of `Petal.Length` between different species, which plot is the most suitable, and why?

7.4.2 Solution:

Now, suppose you want to compare `Petal.Length` across the different species. A boxplot is a good option to show the distribution and potential outliers of `Petal.Length` for each species.

```
# Boxplot of Petal.Length by Species
ggplot(iris, aes(x = Species, y = Petal.Length, fill = Species)) +
  geom_boxplot() +
  labs(title = "Boxplot of Petal Length by Species") +
  theme_minimal()
```



7.5 Exercise 3: Pairwise Scatter Plot

7.5.1 Question:

- What can you learn from the pairwise scatter plot about the relationships between the numerical variables in the dataset? Which pair of variables show the strongest correlation?

7.5.2 Solution:

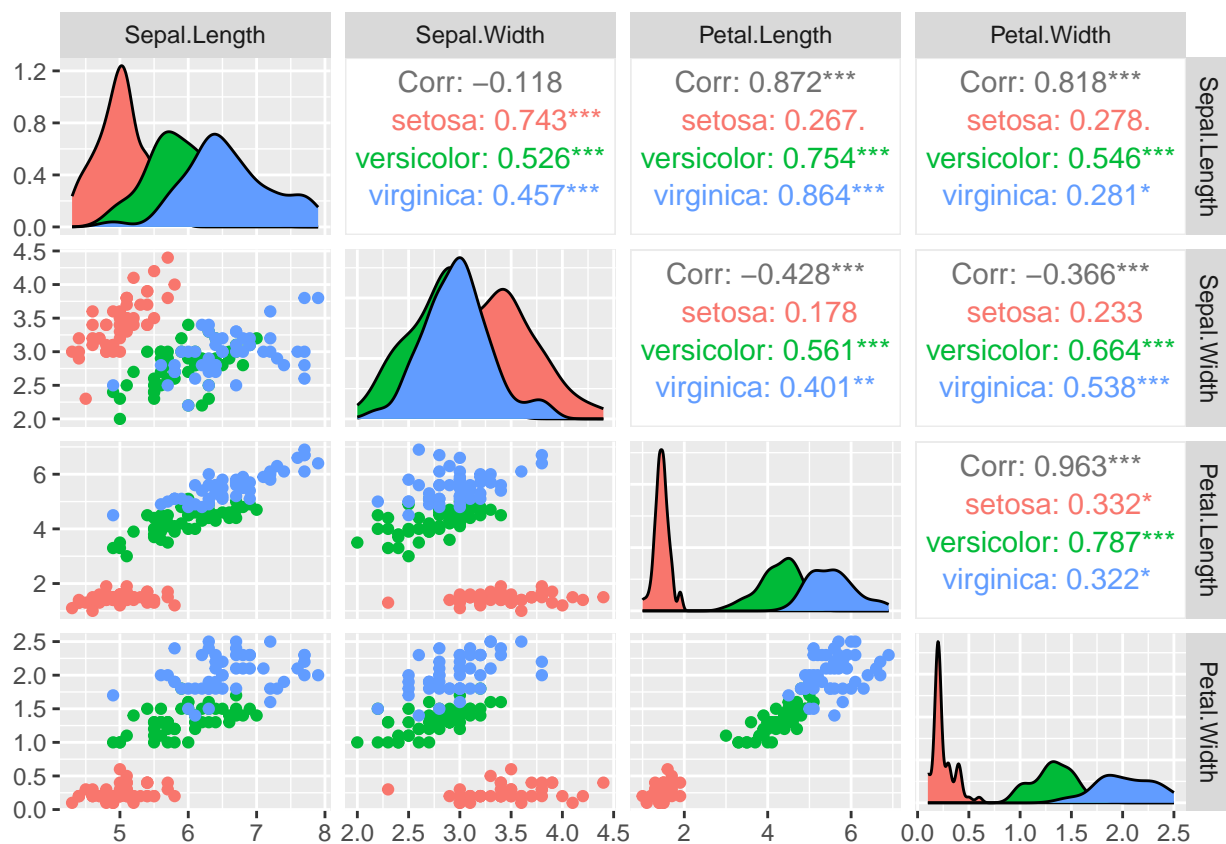
You want to investigate the relationship between multiple numerical variables, say `Sepal.Length`, `Sepal.Width`, `Petal.Length`, and `Petal.Width`. A pairwise scatter plot matrix is useful to examine the relationships between all combinations of variables.

```
# Pairwise scatter plot matrix  
library(GGally)
```

```
## Warning: package 'GGally' was built under R version 4.3.3
```

```
## Registered S3 method overwritten by 'GGally':  
##   method from  
##   +.gg      ggplot2
```

```
ggpairs(iris, columns = 1:4, aes(color = Species))
```



7.6 Exercise 4: Violin Plot for Distribution Comparison

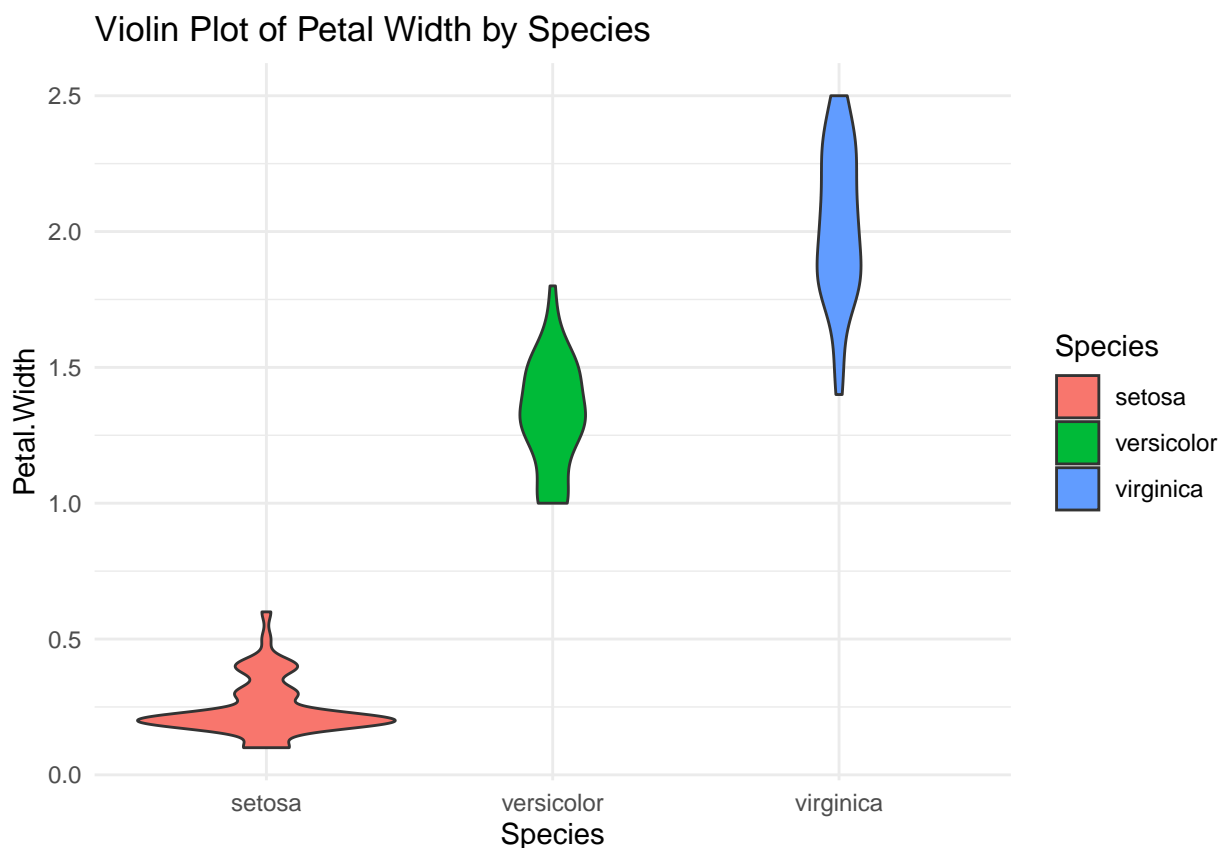
7.6.1 Question:

- When comparing distributions across categories, how does a violin plot differ from a boxplot, and when might you prefer one over the other?

7.6.2 Solution:

If you want to compare the distribution of `Petal.Width` for each species, a violin plot can be an excellent choice, as it combines aspects of boxplots and density plots.

```
# Violin plot of Petal.Width by Species  
ggplot(iris, aes(x = Species, y = Petal.Width, fill = Species)) +  
  geom_violin() +  
  labs(title = "Violin Plot of Petal Width by Species") +  
  theme_minimal()
```



7.7 Exercise 5: Correlation Plot Between Numerical Variables

7.7.1 Question:

- Which pairs of variables are most strongly correlated in the dataset? How does this help you understand the relationships between measurements?

7.7.2 Solution:

If you're interested in the correlation between numerical variables in the dataset, a correlation plot is a helpful visualization. It shows the strength and direction of relationships between variables.

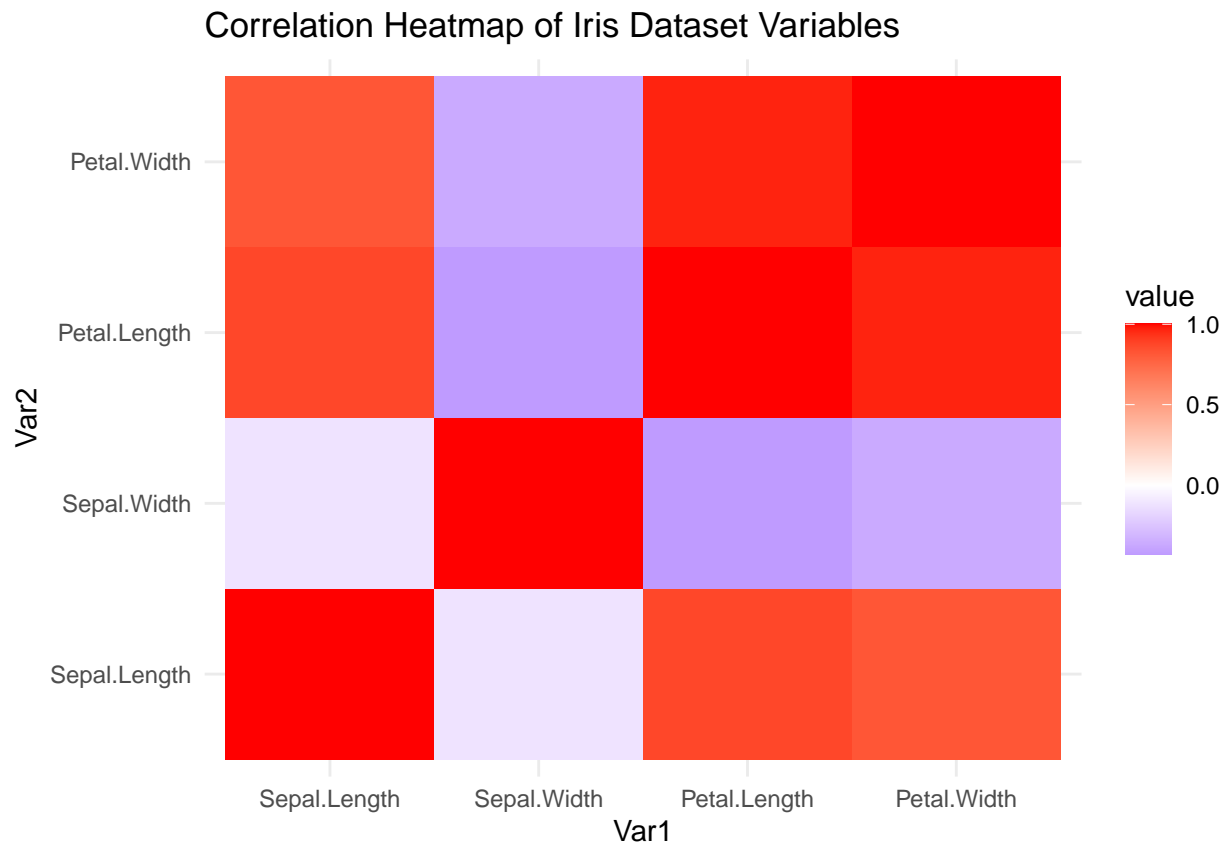
```
# Correlation plot using ggplot2
cor_matrix <- cor(iris[,1:4])
library(reshape2)
```

```
##
## Attaching package: 'reshape2'

## The following object is masked from 'package:tidyr':
##
##      smiths
```

```
# Convert the correlation matrix to long format
cor_melt <- melt(cor_matrix)

# Plot the correlation heatmap
ggplot(cor_melt, aes(Var1, Var2, fill = value)) +
  geom_tile() +
  scale_fill_gradient2(midpoint = 0, low = "blue", high = "red", mid = "white") +
  theme_minimal() +
  labs(title = "Correlation Heatmap of Iris Dataset Variables")
```



7.8 Summary

In this exercise, you explored different visualization techniques using the **iris** dataset to understand the relationships between variables. Depending on the type of relationship (numerical vs categorical, pairs of variables, etc.), different plots such as scatter plots, boxplots, and pairwise plots are useful for drawing insights.

8 More..

ggplot2 is a very powerful package in R that can almost visualize any kind of data. Please go to <https://ggplot2.tidyverse.org/articles/ggplot2.html> for further information if in the future you want to do some fancy plots in R.

9 Acknowledgement

This teaching material is adapted from the previous material of this course made by [Marcela Alfaro-Córdoba](#) and [Sheng Jiang](#).