

# **Data Visualization for Social Science**

**Using R!**

Alfredo Hernandez Sanchez, PhD

2024-04-05

# Table of contents

<b>Welcome</b>	<b>3</b>
Contact . . . . .	3
License . . . . .	3
<b>1 Introduction</b>	<b>4</b>
Why Visualize? . . . . .	5
About this Book . . . . .	6
Recommended Readings . . . . .	6
<b>2 The Grammar of Graphics</b>	<b>7</b>
2.1 The <code>tidyverse</code> Package . . . . .	7
2.2 The <code>ggplot2</code> Package . . . . .	7
2.3 Example . . . . .	8
<b>References</b>	<b>11</b>

# Welcome

This book offers a gentle introduction to data visualization using R and – occasionally – python. In this book you will learn how to make beautiful, informative, and reproducible visualizations with examples from different social sciences. It is inspired by a course on Data Visualization that I taught at the Barcelona Institute of International Studies. Thus, a special thanks goes out to the many graduate students whose efforts, questions, and feedback helped greatly improve the content of this book!

## Contact

This book is in open review. If you have any questions, comments or suggestions; please contact me by [email](#) or report an issue on GitHub.

## License

Data Visualization for Social Science by Alfredo Hernandez Sanchez is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License. Based on a work at [https://github.com/alhdzsz/dviz\\_book](https://github.com/alhdzsz/dviz_book).

# 1 Introduction

In an experiment conducted by researchers from University College London (Mcmanus and Gesiak 2014), 277 participants were asked to look at several pairs of paintings: one of the pairs was an original by abstract painter Piet Mondrian, and the other was fake version that closely resembled it.<sup>1</sup> The participants were asked:

When looking at the pictures you should decide overall which you think looks better, in that it looks nicer, it looks better organised, or it looks better balanced.

The results suggested that people could identify the originals with some degree of accuracy ( $\mu$  54.7%, SE .40). In other words, reliably better than chance! The experiment aimed to compare two methods in *Empirical Aesthetics*: the method of choice and the method of production. This choice experiment “implies people know something about what makes a real Mondrian.” In other words, we have an *intuition* of proportion and beauty.

---

<sup>1</sup>The pseudo-Mondrians were created by jittering all the lines in the original but keeping the same relative positions.

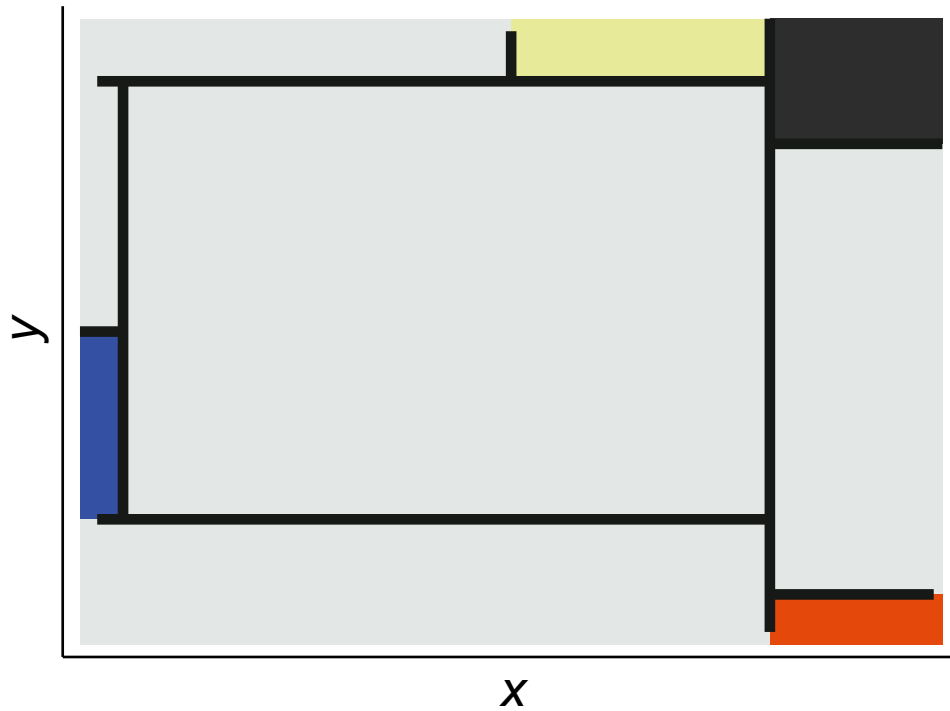


Figure 1.1: An *attempt* at a Mondrian using `ggplot2`

## Why Visualize?

Data visualization plays numerous roles in the social sciences, from summarizing large amounts of information in a small space, to supporting claims about patterns and relationships among a vast array of indicators of human behavior.

Zinovyev (Zinovyev 2010) identifies four types of visualizations in political science:

- *Statistical graphics and infographics* with extensive use of color, form, size, shape and style to superimpose many quantitative variables in the same chart or diagram
- *Geographical information systems (GIS)* to visualize geographically-linked data
- *Graph visualization or network maps* for representing relations between objects
- Projection of multidimensional data on low-dimensional screens with further visualization, *data cartography*

Throughout this book, we will cover examples of from all of these types of visuals.

## About this Book

If you are in the [Data Visualization](#) course at IBEL, you will need it. If you are not, it might be useful anyway!

To keep things as simple as possible, the book follows this syntax:

- `packages` are placed inside a shaded box (e.g. `ggplot2`),
- common `functions()` are also inside a shaded box, and followed by parentheses (e.g. `mutate()` from `tidyverse` or `class()` from base R),
- less common `functions()` are the same, but the package is explicitly called `::` (e.g. `reshape2::melt()`),
- short R commands (e.g. `%in%`), are also shaded, non-R commands are in bold (e.g. **Ctrl + P**),
- the common *pipe* operator `%>%` will be used when possible in the code<sup>2</sup> (i.e., we will mostly use the `tidyverse` syntax over base R).

## Recommended Readings

You are not expected to have any familiarity with R at the beginning of the course, though some knowledge of statistics will be very helpful. We will cover the basics of working with R and [RStudio](#) during the first few sessions. Some tutorial videos on the basics of working with RStudio are available here. Similarly, you may also consult the following open-source books on R:<sup>3</sup>

- [R Cookbook](#) (Long and Teetor 2019)
- [R for Data Science](#) (Grolemund and Wickham 2016)
- [R Graphics Cookbook](#) (Chang 2018)
- [Efficient R Programming](#) (Gillespie and Lovelace 2016)
- [Hands-on Programming with R](#) (Grolemund 2014)
- [Fundamentals of Data Visualization](#) (Wilke 2019)
- [Text Mining with R](#) (Silge and Robinson 2017)
- [An Introduction to R](#) (Venables, Smith, and R Core Team 2021)
- [R Markdown: The Definitive Guide](#) (Xie, Allaire, and Grolemund 2018)
- [R Markdown Cookbook](#) (Xie, Dervieux, and Riederer 2020)

---

<sup>2</sup>For Windows users, the `%>%` shortcut in RStudio is **Ctrl + Shift + M** and for Mac users it is **Cmd + Shift + M**.

<sup>3</sup>For a comprehensive list of R-related books, consult the *R-Project Website*

## 2 The Grammar of Graphics

### 2.1 The tidyverse Package

Throughout this course, we will be using tidy data principles<sup>1</sup> to create several types of visualizations. The main package we will use is the `tidyverse`, which includes several useful tools for data wrangling, analysis and visualization. The first step then is to install the package! You can do this from the packages vignette in *explorer pane* in RStudio, or by writing `install.packages("tidyverse")` into the *console pane*.

Once the package has been installed, the next step will be to load the library so that we can start using it! Simply write the command below in a script the *editor pane* and click *run*, or directly in the *console pane* and press *enter*.

```
library(tidyverse)
```

After loading the `tidyverse` package from the library, we will get access to two very important functions which we will be using extensively. The first is the the command `ggplot()` which will allow us to make plots based on the *grammar of graphics*. The second is the *pipe operator* or `%>%`, which translates loosely to the phrase “and then”, and which we will use to put several commands and functions together in a pipeline.<sup>2</sup>

### 2.2 The ggplot2 Package

The `ggplot2` package is installed and loaded alongside the `tidyverse` package, though it can also be called on separately. This is a very powerful tool to make print-quality graphs and all sorts of visual outputs. To do this, it draws on *the grammar of graphics*, which is a concept developed by Leland Wilkinson (Wilkinson 2005). The main idea behind this complex book is that plots can be divided into several elements, each with a specific role to play. `ggplot2` has 7 such elements:

- **Data**

---

<sup>1</sup>These principles are: a) each variable should have its own column, b) each observation should have its own row, and c) each value should have its own cell.

<sup>2</sup>For Windows users, you can use the RStudio short cut `ctrl + shift + m` to write this pipe `%>%` operator.

- **Aesthetics**
- **Layers**
- Scales
- Coordinates
- Facets
- Themes

Throughout this chapter, we will focus on the first three (**Data**, **Aesthetics**, **Layers**) which are the minimum requirements to make a basic plot. The element **data** tells R which vector(s) from your environment are going to be used to draw the plot. The **aesthetics** element determines which variable(s) will be used and in what capacity. The **layers** element tells R which type of geometry you wish to draw and in which order.

```
df %>%
  ggplot(aes(x=var1,y=var2))+
  geom_point()
```

In the example above, we are telling R that there is an object **df** in our environment which has at least two vectors (columns), one called **var1** and another **var2**. We are also telling it that we want **var1** to be our **x** axis and **var2** to be our **y** axis, we define this inside the **aes()** command either globally for the plot (i.e. inside the **ggplot()** command) or specifically for a layer (i.e. inside **geom\_point()**). Finally, we are telling R that we want to make a scatter plot by defining the layer **geom\_point()**. Notice that after the **ggplot()** command and until the end of the graph, we use a **+** sign.

## 2.3 Example

To make our first **ggplot** plot, we will use the **mtcars** data set as an example.

```
data("mtcars")
```

The cars data set has 32 observations and 11 variables. Once the data has been loaded, let's use the **pipe operator** to do some cleaning. In the code below, we are creating a new object called **df** - a common way of naming data frames - and filling it with the **mtcars** data with some modifications. We are asking R to a) take the **mtcars** data, b) *and then* **%>%** select four variables, c) *and then* **%>%** give them new name. This pipeline is saved into the new object **df**.



```
df <- mtcars %>%
  select(cyl, mpg, hp, am) %>%
  rename(cylinders=cyl,
         mileage=mpg,
         horsepower=hp,
         transmission=am)
```

With this `df` stored in our environment, we can start making plots. Let's begin with a histogram that shows the distribution of mileage across the 32 variables in our data set. For this we will use `geom_histogram`.

```
df %>% #Our Data
  ggplot(aes(x= mileage))+ #Our Aesthetics
  geom_histogram() #Our Layer
```

``stat_bin()`` using ``bins = 30``. Pick better value with ``binwidth``.

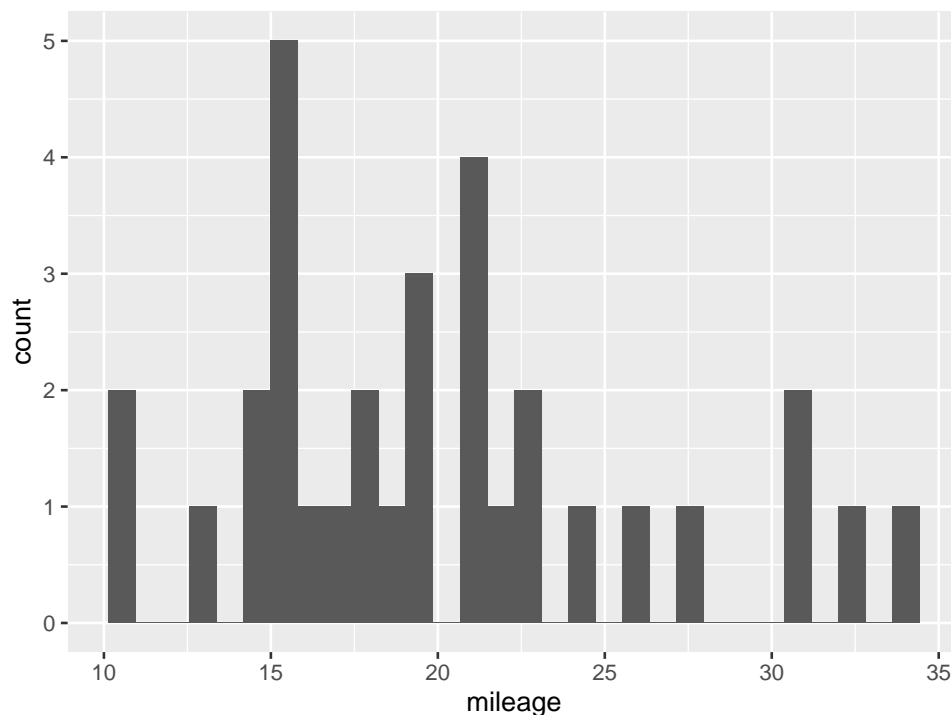


Figure 2.1: A histogram with default settings

Figure @ref(fig:hist-fig) shows us our very first **ggplot**, which shows the number of observations at each of the binned levels. From this plot we know that most cars in our data set do around

15 miles per gallon. However, it is not very nice looking! We can improve this by adding more parameters.

You might notice that below the code R is giving us a **warning: stat\_bin() using bins = 30. Pick better value with binwidth.** Here the software is hinting that we might want to change the number of bars (**bins**) or their width (**binwidth**) in our plot to make it more informative.<sup>3</sup> In figure @ref(fig:hist2-fig) we change the number of bins to 5 inside our `geom_histogram` layer, and also declare the color of the column fill (darkgray) and the outline (black).

```
df %>%  
  ggplot(aes(x= mileage))+  
  geom_histogram(bins = 5, fill="darkgray", color="black")
```

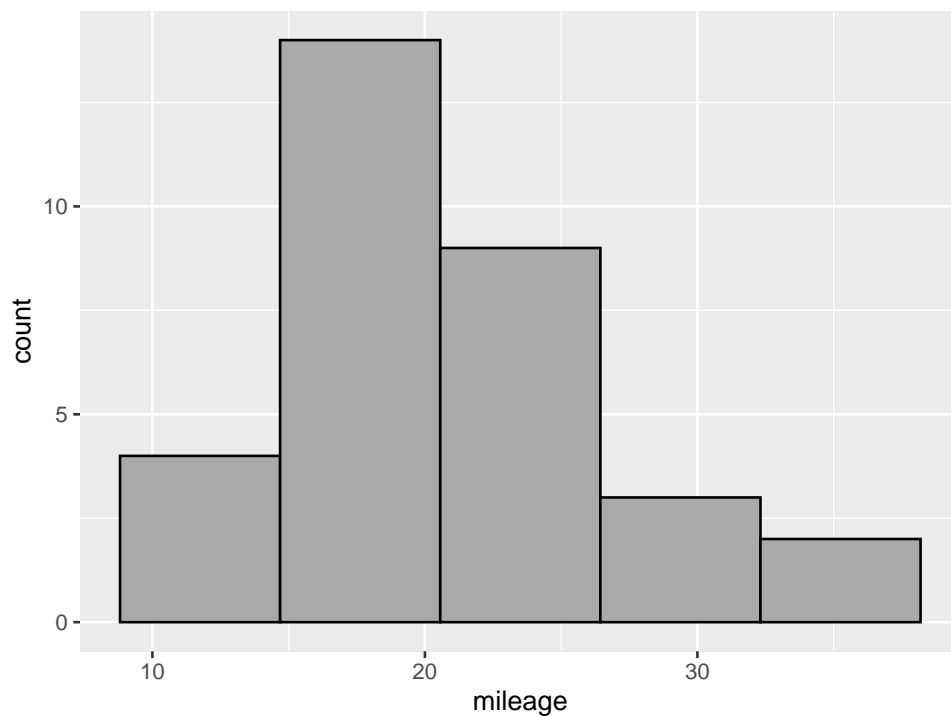


Figure 2.2: A nicer looking histogram

---

<sup>3</sup>Most other software will give you a default based on some parameter such as the Freedman-Diaconis rule, `ggplot` does not do this, forcing you to experiment with different parameters that best reflect your data.

# References

- Chang, Winston. 2018. *R Graphics Cookbook: Practical Recipes for Visualizing Data*. Second. Sebastopol, California: O'Reilly Media. <https://r-graphics.org/>.
- Gillespie, Colin, and Robin Lovelace. 2016. *Efficient r Programming: A Practical Guide to Smarter Programming*. Sebastopol, California: O'Reilly Media. <https://csgillespie.github.io/efficientR/>.
- Grolemund, Garrett. 2014. *Hands-on Programming with r: Write Your Own Functions and Simulations*. Sebastopol, California: O'Reilly Media. <https://rstudio-education.github.io/hopr/>.
- Grolemund, Garrett, and Hadley Wickham. 2016. *R for Data Science*. Sebastopol, California: O'Reilly Media. <https://r4ds.had.co.nz/>.
- Long, JD, and Paul Teetor. 2019. *R Cookbook: Proven Recipes for Data Analysis, Statistics, and Graphics*. Second. Sebastopol, California: O'Reilly Media. <https://rc2e.com/>.
- Mcmanus, Ian, and Paul Gesiak. 2014. "Experimenting with Mondrian: Comparing the Method of Production with the Method of Choice." In. <https://doi.org/10.13140/2.1.1561.2967>.
- Silge, Julia, and David Robinson. 2017. *Text Mining with r: A Tidy Approach*. Sebastopol, California: O'Reilly Media. <https://www.tidytextmining.com/>.
- Venables, W. N., D. M. Smith, and the R Core Team. 2021. *An Introduction to r*. <https://cran.r-project.org/doc/manuals/R-intro.pdf>.
- Wilke, Claus O. 2019. *Fundamentals of Data Visualization: A Primer on Making Informative and Compelling Figures*. Sebastopol, California: O'Reilly Media. <https://clauswilke.com/dataviz/>.
- Wilkinson, Leland. 2005. *The Grammar of Graphics*. 2nd ed. New York: Springer-Verlag. <https://www.springer.com/gp/book/9780387245447>.
- Xie, Yihui, J. J. Allaire, and Garrett Grolemund. 2018. *R Markdown: The Definitive Guide*. Boca Raton, Florida: Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown>.
- Xie, Yihui, Christophe Dervieux, and Emily Riederer. 2020. *R Markdown Cookbook*. Boca Raton, Florida: Chapman; Hall/CRC. <https://bookdown.org/yihui/rmarkdown-cookbook>.
- Zinovyev, iAndrei. 2010. "Data Visualization in Political and Social Sciences." In. <https://arxiv.org/abs/1008.1188>.