

Writing a reproducible paper in R Markdown*

Paul C. Bauer

Mannheim Centre for European Social Research

5 November, 2018

Feedback?

mail@paulcbauer.eu

Abstract

The present paper provides a template for a reproducible scientific paper written in R Markdown. Below I outline some of the “tricks”/code (e.g., referencing tables, sections etc.) I had to figure out to produce this document. The underlying files which produce this document can be downloaded [here](#). I think I got pretty far but there is always room for improvement and more automatization, in parallel to the incredible developments in R and Rstudio (bookdown etc.). I intend to update this file when I discover more convenient code.

*Corresponding address: mail@paulcbauer.eu. Acknowledgments: I am grateful to all those generous people that invest their time into open-source software.

Contents

1	Why reproducible research (in R)?	3
2	Prerequisites	3
3	Basics: Input and output files	4
4	Referencing within your document	5
5	Software versioning	5
6	Data	6
6.1	Import	6
6.2	Putting your entire data into the .rmd file	7
7	Tables	7
7.1	stargazer(): Summary and regression tables	8
7.2	kable() and kable_styling()	9
8	Inline code & results	10
9	Figures	11
9.1	R base graphs	11
9.2	ggplot2 graphs	11
9.3	Plotly graphs	12
10	Good practices	13
11	Citation styles	14
12	Appendix	15
12.1	All the code in the paper	15
	References	17

1 Why reproducible research (in R)?

Some arguments...

- **Access:** Research is normally funded by taxpayers (researchers are also taxpayers). Hence, it should be freely accessible to everyone without any barriers, e.g., without requiring commercial software. Importantly, researchers from developing countries are even more dependent on free access to knowledge (Kirsop and Chan 2005).
- **Reproducibility:** Even if you have written a study and analyzed the data yourself you will forget what you did after a few months. A fully reproducible setup will help you to trace back your own steps. Obviously, the same is true for other researchers who may want to understand your work and built on it. It may sound like a joke but why not aim for a document that can be used to reproduce your findings in 500 years.
- **Errors:** Manual steps in data analysis (e.g., manually copy/pasting values into a table etc.) may introduce errors. R Markdown allows you to **automatize** such steps and/or avoid them.
- **Revisions:** Revising a paper takes much less time if you have all the code you need in one place, i.e., one .rmd file. For instance, if you decide to exclude a subset of your data you simply need to insert one line of your code at the beginning and everything is rebuilt/re-estimated automatically.

2 Prerequisites

I assume that you are using R on a day-to-day basis. You may have even started to work a little in R Markdown but you don't write your complete paper in R Markdown. If you don't know what R Markdown is watch [this short video](#). Then...

- ...install [R](#) and [Rstudio](#) (most recent versions) (R Core Team 2017; RStudio Team 2015).
- ...install [tinytex](#), a lightweight version of Tex Live (Allaire et al. 2017; Xie 2018b).

```
install.packages(c('tinytex', 'rmarkdown'))
tinytex::install_tinytex()
```

- ...install the packages below using the code below (Sievert et al. 2017; Xie 2014, 2015, 2016, 2017, 2018a; Zhu 2017).

```
install.packages(c("rmarkdown", "knitr", "kableExtra",  
                  "stargazer", "plotly", "knitr",  
                  "bookdown"))
```

- ...download the 5 input files I created — `paper.rmd`, `header.tex`, `references.bib`, `data.csv` and `american-sociological-association.csl` — from [this folder](#). Ignore the other files.
- ...learn R and read about the other underlying components namely [Markdown](#), [R Markdown](#) and [Latex](#).

3 Basics: Input and output files

All the files you need to produce the present PDF file are the input files...

- ...a `paper.rmd` file (the underlying R Markdown file).
- ...a `header.tex` file (a tex file for document structure).
- ...a `references.bib` file (the bibliography).
 - I use `paperpile` to manage my references and export the `.bib` file into the folder that contains my `.rmd` file.
- ...a `data.csv` file (some raw data).
- ... a `american-sociological-association.csl` file that defines the style of your bibliography.¹

[Download these files](#) and save them into a folder. Close R/Rstudio and directly open `paper.rmd` with RStudio. Doing so assures that the working directory is set to the folder that contains `paper.rmd` and the other files.²

Once you run/compile the `paper.rmd` file in Rstudio it creates mainly two output files:

- `paper.tex`
- `paper.pdf` (the one you are reading right now)

¹You can download various citation style files from this webpage: <https://github.com/citation-style-language/styles>.

²You can always check your working directory in R with `getwd()`.

In addition, there may be files that you generate and store locally in the folder during the compilation process. This is the case for some of the Plotly graphs below.

Ideally, we can simply provide others with a zip folder that contains both our input files and our output files. Then it's possible to reproduce the process from managing/analyzing some raw data to producing the final scientific article.

Below we always display the R code in the chunks that produce the output. In your paper you will normally only present outputs (e.g., tables, figures etc.) by choosing the chunk option “Show output only” in R Studio. The chunk commands itself are not displayed but they do matter for referencing etc. So simply orient yourself at the underlying `paper.rmd` file.

4 Referencing within your document

To see how referencing works simply see the different examples for figures, tables and sections below. For instance in Section 7 you can find different ways of referencing tables. The code of the underlying `paper.rmd` will show you how I referenced Section 7 right here namely with `'Section \@ref(sec:tables)'`.

5 Software versioning

Software changes and gets updated, especially with an active developer community like that of R. Luckily you can always access [old versions of R](#) and old version of R packages in [the archive](#). In the archive you need to choose a particular package, e.g dplyr and search for the right version, e.g., `dplyr_0.2.tar.gz`. Then insert the path in the following function: `install.packages("https://...../dplyr_0.2.tar.gz", repos=NULL, type="source")`. Ideally, however, results will be simply reproducible in the most current R and package versions.

I would recommend to use the command below and copy/paste the output into the first R chunk in your `paper.rmd` file. Normally we would do this when we compile/generate the paper for the last time. For more advanced tools see [packrat](#).

```
cat(paste("#", capture.output(sessionInfo()), "\n", collapse = ""))
```

```
## # R version 3.5.1 (2018-07-02)
## # Platform: x86_64-w64-mingw32/x64 (64-bit)
## # Running under: Windows 10 x64 (build 17134)
```

```
## #
## # Matrix products: default
## #
## # locale:
## # [1] LC_COLLATE=English_Germany.1252 LC_CTYPE=English_Germany.1252
## # [3] LC_MONETARY=English_Germany.1252 LC_NUMERIC=C
## # [5] LC_TIME=English_Germany.1252
## #
## # attached base packages:
## # [1] stats      graphics  grDevices  utils      datasets  methods    base
## #
## # loaded via a namespace (and not attached):
## # [1] compiler_3.5.1 backports_1.1.2 magrittr_1.5    bookdown_0.7
## # [5] rprojroot_1.3-2 tools_3.5.1    htmltools_0.3.6 yaml_2.2.0
## # [9] Rcpp_0.12.19   stringi_1.1.7 rmarkdown_1.10 knitr_1.20
## # [13] stringr_1.3.1  xfun_0.3      digest_0.6.18  evaluate_0.12
```

```
# or use message() instead of cat()
```

6 Data

6.1 Import

```
data <- read.csv("data.csv")
head(data)
```

```
##   X speed dist
## 1 1     4    2
## 2 2     4   10
## 3 3     7    4
## 4 4     7   22
## 5 5     8   16
## 6 6     9   10
```

6.2 Putting your entire data into the .rmd file

Applying the function `dput()` to an object gives you the code needed to reproduce that object. So you could paste that code into your .rmd file if you don't want to have extra data files. This makes sense where data files are small.

```
dput(data)
```

```
## structure(list(X = 1:50, speed = c(4L, 4L, 7L, 7L, 8L, 9L, 10L,  
## 10L, 10L, 11L, 11L, 12L, 12L, 12L, 12L, 13L, 13L, 13L, 13L, 14L,  
## 14L, 14L, 14L, 15L, 15L, 15L, 16L, 16L, 17L, 17L, 17L, 18L, 18L,  
## 18L, 18L, 19L, 19L, 19L, 20L, 20L, 20L, 20L, 20L, 22L, 23L, 24L,  
## 24L, 24L, 24L, 25L), dist = c(2L, 10L, 4L, 22L, 16L, 10L, 18L,  
## 26L, 34L, 17L, 28L, 14L, 20L, 24L, 28L, 26L, 34L, 34L, 46L, 26L,  
## 36L, 60L, 80L, 20L, 26L, 54L, 32L, 40L, 32L, 40L, 50L, 42L, 56L,  
## 76L, 84L, 36L, 46L, 68L, 32L, 48L, 52L, 56L, 64L, 66L, 54L, 70L,  
## 92L, 93L, 120L, 85L)), class = "data.frame", row.names = c(NA,  
## -50L))
```

You can then insert the `dput` output in your .rmd as below.

```
data <- structure(list(X = 1:50, speed = c(4L, 4L, 7L, 7L, 8L, 9L, 10L,  
10L, 10L, 11L, 11L, 12L, 12L, 12L, 12L, 13L, 13L, 13L, 13L, 14L,  
14L, 14L, 14L, 15L, 15L, 15L, 16L, 16L, 17L, 17L, 17L, 18L, 18L,  
18L, 18L, 19L, 19L, 19L, 20L, 20L, 20L, 20L, 20L, 20L, 22L, 23L, 24L,  
24L, 24L, 24L, 25L), dist = c(2L, 10L, 4L, 22L, 16L, 10L, 18L,  
26L, 34L, 17L, 28L, 14L, 20L, 24L, 28L, 26L, 34L, 34L, 46L, 26L,  
36L, 60L, 80L, 20L, 26L, 54L, 32L, 40L, 32L, 40L, 50L, 42L, 56L,  
76L, 84L, 36L, 46L, 68L, 32L, 48L, 52L, 56L, 64L, 66L, 54L, 70L,  
92L, 93L, 120L, 85L)),  
class = "data.frame", row.names = c(NA,  
-50L))
```

7 Tables

Producing good tables and referencing these tables within a R Markdown PDF has been a hassle but got much better. Examples that you may use are shown below. The way you reference tables

is slightly different, e.g., for `stargazer` the label is contained in the function, for `kable` it's contained in the chunk name.

7.1 `stargazer()`: Summary and regression tables

Table 1 shows summary stats of your data.³ I normally use `stargazer()` (Hlavac 2013) which offers extreme flexibility regarding table output (see `?stargazer`).

```
library(stargazer)
stargazer(cars,
  title = "Summary table with stargazer",
  label="tab1",
  table.placement = "H",
  header=FALSE)
```

Table 1: Summary table with stargazer

Statistic	N	Mean	St. Dev.	Min	Pctl(25)	Pctl(75)	Max
speed	50	15.400	5.288	4	12	19	25
dist	50	42.980	25.769	2	26	56	120

Table 2 shows the output for a regression table. Make sure you name all your models and explicitly refer to model names (M1, M2 etc.) in the text.

```
library(stargazer)
model1 <- lm(speed ~ dist, data = cars)
model2 <- lm(speed ~ dist, data = cars)
model3 <- lm(dist ~ speed, data = cars)
stargazer(model1, model2, model3,
  title = "Regression table with stargazer",
  label="tab2",
  table.placement = "H",
  column.labels = c("M1", "M2", "M3"),
  model.numbers = FALSE,
  header=FALSE)
```

³To reference the table where you set the identifier in the `stargazer` function you only need to use the actual label, i.e., `tab1`.

Table 2: Regression table with stargazer

	<i>Dependent variable:</i>		
	speed		dist
	M1	M2	M3
dist	0.166*** (0.017)	0.166*** (0.017)	
speed			3.932*** (0.416)
Constant	8.284*** (0.874)	8.284*** (0.874)	-17.579** (6.758)
Observations	50	50	50
R ²	0.651	0.651	0.651
Adjusted R ²	0.644	0.644	0.644
Residual Std. Error (df = 48)	3.156	3.156	15.380
F Statistic (df = 1; 48)	89.567***	89.567***	89.567***

Note: *p<0.1; **p<0.05; ***p<0.01

7.2 kable() and kable_styling()

Another great function is `kable()` (knitr package) in combination with `kableExtra`. Table 3 provides an example.⁴ Again you can modify so many things in both the `kable()` and the `kable_styling()` function. See [this overview](#) of all the kable stylings that are possible provided by the package author himself.

```
library(knitr)
library(kableExtra)
kable(cars[1:10,], row.names = TRUE,
      caption = 'Table with kable() and kablestyling()',
      format = "latex", booktabs = T) %>%
  kable_styling(full_width = T,
```

⁴To reference the table produced by the chunk you need to add `tab:` to the chunk name, i.e., `tab:tab3`.

```

latex_options = c("striped",
                  "scale_down",
                  "HOLD_position"),
font_size = 10)

```

Table 3: Table with kable() and kablestyling()

	speed	dist
1	4	2
2	4	10
3	7	4
4	7	22
5	8	16
6	9	10
7	10	18
8	10	26
9	10	34
10	11	17

8 Inline code & results

Reproduction reaches new heights when you work with inline code. For instance, you can automatize the display of certain coefficients within the text. An example is to include estimates, e.g., the coefficient of `dist` of the model we ran above. ``r round(coef(model1)[2], 2)`` will insert the coefficient as follows: 0.17. Or ``r 3 + 7`` will insert a 10 in the text.

Inline code/results that depend on earlier objects in your document will automatically be updated once you change those objects. For instance, imagine a reviewer asks you to omit certain observations from your sample. You can simply do so in the beginning of your code and push play subsequently.. at time you might have to set `cache = FALSE` at the beginning so that all the code chunks are rerun.

Researchers often avoid referring to results in-text etc. because you easily forget to change them when revising a manuscript. At the same it can make an article much more informative and easier to read, e.g., if you discuss a coefficient in the text you can directly show it in the section in which you discuss it. Inline code allows you to do just that. R Markdown allows you to that do so in a reproducible and automatized manner.

9 Figures

9.1 R base graphs

Inserting figures can be slightly more complicated. Ideally, we would produce and insert them directly in the .rmd file. It's relatively simple to insert R base graphs as you can see in Figure 1.

```
plot(cars$speed, cars$dist)
```

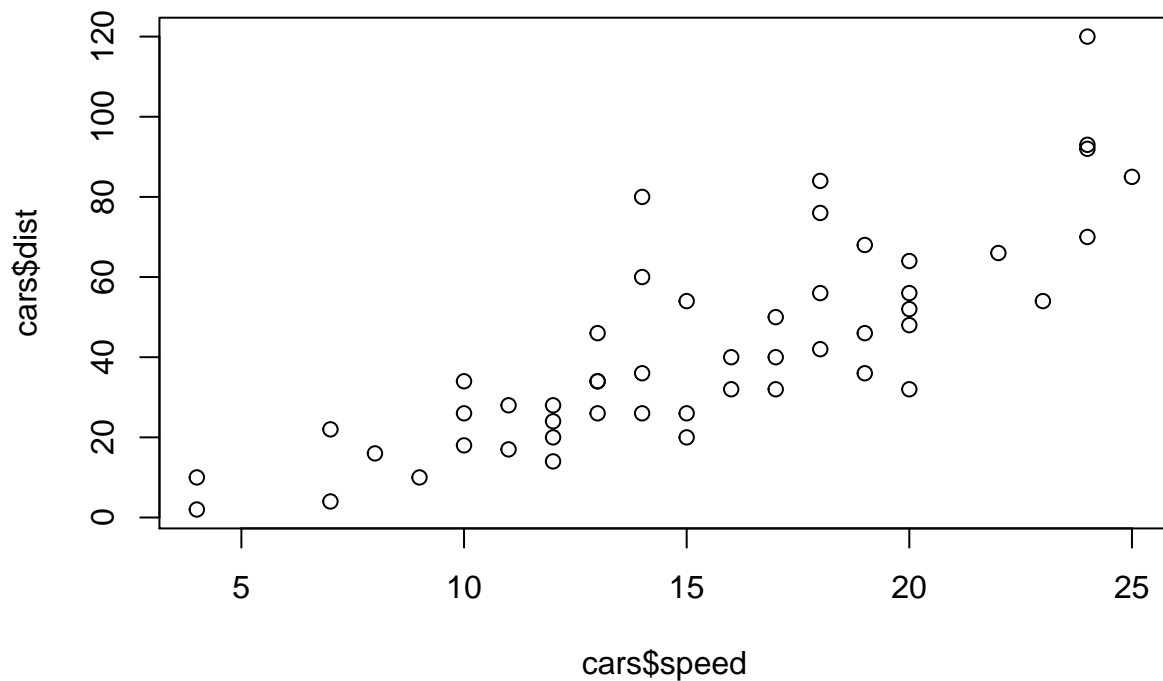


Figure 1: Scatterplot of Speed and Distance

But it turns out that it doesn't always work so well.

9.2 ggplot2 graphs

Same is true for ggplot2 as you can see in Figure 2.

```
mtcars$cyl <- as.factor(mtcars$cyl) # Convert cyl to factor
library(ggplot2)
ggplot(mtcars, aes(x=wt, y=mpg, shape=cyl)) + geom_point() +
  labs(x="Weight (lb/1000)", y = "Miles/(US) gallon",
       shape="Number of \n Cylinders") + theme_classic()
```

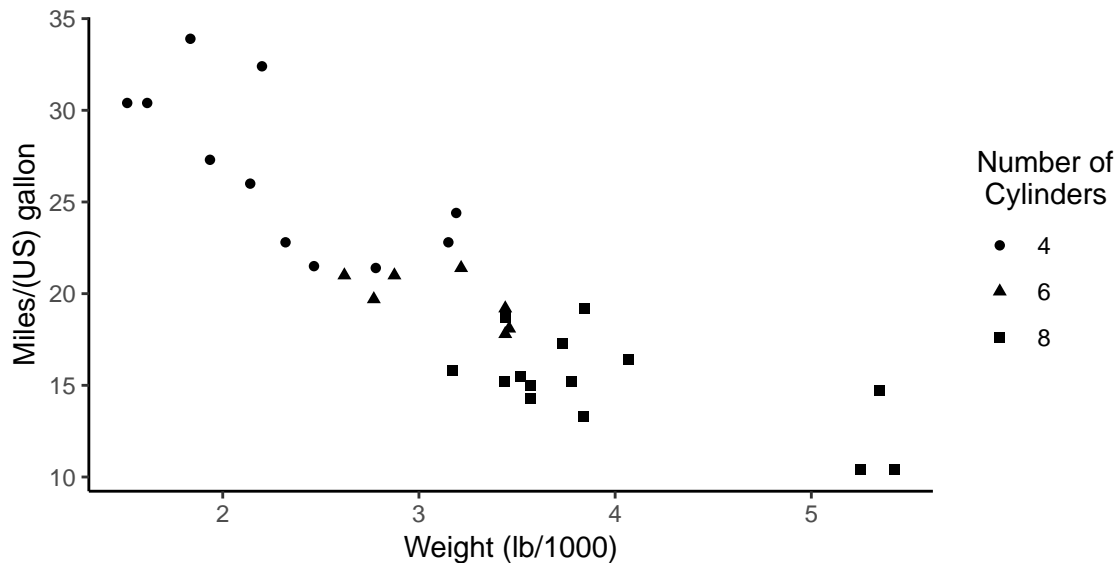


Figure 2: Miles per gallon according to the weight

9.3 Plotly graphs

Plotly is a popular graph engine that let's you also produce interactive graphs that you can embed in html webpages or documents (e.g., see [here](#)). I am a big fan, however, right now there is no easy, automatic way to insert high resolution Plotly graphs into your R Markdown PDF. See Figure 3 below.

There seems to be a new engine for exporting graphs locally without relying on the plotly server — which you wanna avoid for a properly reproducible PDF — called orca (?orca). However, the [installation](#) is too challenging for the average user at the moment.

```
library(plotly)
plot_ly(cars, type = "scatter", mode="markers",
        x=~speed,
        y=~dist)
```

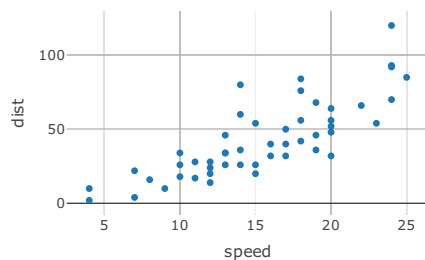


Figure 3: Scatterplot of Speed and Distance

10 Good practices

Every researcher has his own optimized setup. Currently I would recommend the following:

- Keep all files of your project (that matter for producing the PDF) in one folder without subfolders. You can zip and directly upload that folder to the [Harvard dataverse](#)).
- Make sure that filenames have a logic to them.
 - Main file with text/code: “paper.rmd”, “report.rmd”
 - Data files: “data_XXXXXX.*”
 - Image files: “fig_XXXXXX.*”
 - Tables files: “table_XXXX.*”
 - etc.

- Ideally, your filenames will correspond to the names in the paper. For instance, Figure 1 in the paper may have a corresponding file called `fig_1_xxxxx.pdf`.
- Use document the document outline in R studio (Ctrl + Shift + O) when you work with R Markdown.
- Name rchunks according to what they do or produce:
 - “fig-...” for chunks producing figures
 - “table-...” for chunks producing tables
 - “model-...” for chunks producing model estimates
 - “import-...” for chunks importing data
 - “recoding-...” for chunks in which data is recoded
- Use “really” informative variable names:
 - Q: What do you think does the variable *trstep* measure? It actually measures trust in the European parliament.
 How could we call this variable instead? Yes, `trust.european.parliament` which is longer but will probably be understood by one of your descendants in 100 years.
 - If your setup is truly reproducible you will probably re-use the variable names that you generate as variable names in the tables you produce. Hence, there is an incentive to use good names.
- Use unique identifiers in the final document:
 - e.g., name the models you estimate “M1”, “M2” etc.
 - These unique names should also appear in the published paper.
 - Think of someone who wants to produce Figure 1/Model 1 in your paper but doesn’t find it in your code...

11 Citation styles

If your study needs to follow a particular citation style, you can set the corresponding style in the header of your `.rmd` document. To do so you have to download the corresponding `.csl` file.

In the present document we use the style of the American Sociological Association and set it in the preamble with `csl: american-sociological-association.csl`. However, you also need to download the respective `.csl` file from the following github page: <https://github.com/citation-style-language/styles> and copy it into your working directory for it to work.

The github directory contains a wide variety of citation style files depending on what discipline you work in.

12 Appendix

12.1 All the code in the paper

To simply attach all the code you used in the PDF file in the appendix see the R chunk in the underlying .rmd file:

```
knitr::opts_chunk$set(cache = FALSE)
# Use chache = TRUE if you want to speed up compilation

# A function to allow for showing some of the inline code
rinline <- function(code){
  html <- '<code class="r">```` `r CODE` ````</code>'
  sub("CODE", code, html)
}
install.packages(c('tinytex', 'rmarkdown'))
tinytex::install_tinytex()
install.packages(c("rmarkdown", "knitr", "kableExtra",
                  "stargazer", "plotly", "knitr",
                  "bookdown"))
cat(paste("#", capture.output(sessionInfo()), "\n", collapse = ""))
# or use message() instead of cat()
data <- read.csv("data.csv")
head(data)
dput(data)
data <- structure(list(X = 1:50, speed = c(4L, 4L, 7L, 7L, 8L, 9L, 10L,
10L, 10L, 11L, 11L, 12L, 12L, 12L, 12L, 13L, 13L, 13L, 13L, 14L,
14L, 14L, 14L, 15L, 15L, 15L, 16L, 16L, 17L, 17L, 17L, 18L, 18L,
18L, 18L, 19L, 19L, 19L, 20L, 20L, 20L, 20L, 20L, 22L, 23L, 24L,
24L, 24L, 24L, 25L), dist = c(2L, 10L, 4L, 22L, 16L, 10L, 18L,
26L, 34L, 17L, 28L, 14L, 20L, 24L, 28L, 26L, 34L, 34L, 46L, 26L,
36L, 60L, 80L, 20L, 26L, 54L, 32L, 40L, 32L, 40L, 50L, 42L, 56L,
76L, 84L, 36L, 46L, 68L, 32L, 48L, 52L, 56L, 64L, 66L, 54L, 70L,
92L, 93L, 120L, 85L)),
```

```

class = "data.frame", row.names = c(NA,
-50L))
library(stargazer)
stargazer(cars,
          title = "Summary table with stargazer",
          label="tab1",
          table.placement = "H",
          header=FALSE)
library(stargazer)
model1 <- lm(speed ~ dist, data = cars)
model2 <- lm(speed ~ dist, data = cars)
model3 <- lm(dist ~ speed, data = cars)
stargazer(model1, model2, model3,
          title = "Regression table with stargazer",
          label="tab2",
          table.placement = "H",
          column.labels = c("M1", "M2", "M3"),
          model.numbers = FALSE,
          header=FALSE)
library(knitr)
library(kableExtra)
kable(cars[1:10,], row.names = TRUE,
      caption = 'Table with kable() and kablestyling()',
      format = "latex", booktabs = T) %>%
  kable_styling(full_width = T,
                latex_options = c("striped",
                                "scale_down",
                                "HOLD_position"),
                font_size = 10)
plot(cars$speed, cars$dist)
mtcars$cyl <- as.factor(mtcars$cyl) # Convert cyl to factor
library(ggplot2)
ggplot(mtcars, aes(x=wt, y=mpg, shape=cyl)) + geom_point() +
  labs(x="Weight (lb/1000)", y = "Miles/(US) gallon",
       shape="Number of \n Cylinders") + theme_classic()
library(plotly)
plot_ly(cars, type = "scatter", mode="markers",
        x=~speed,

```



```
y=~dist)
```

References

- Allaire, JJ, Jeffrey Horner, Vicent Marti, and Natacha Porte. 2017. *Markdown: 'Markdown' Rendering for R*.
- Hlavac, Marek. 2013. "Stargazer: LaTeX Code and Ascii Text for Well-Formatted Regression and Summary Statistics Tables." URL: [Http://CRAN.R-Project. Org/Package= Stargazer](http://CRAN.R-project.org/package=Stargazer).
- Kirsop, Barbara and Leslie Chan. 2005. "Transforming Access to Research Literature for Developing Countries." *Serials Review* 31(4):246–55.
- R Core Team. 2017. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing.
- RStudio Team. 2015. *RStudio: Integrated Development Environment for R*. Boston, MA: RStudio, Inc.
- Sievert, Carson, Chris Parmer, Toby Hocking, Scott Chamberlain, Karthik Ram, Marianne Corvellec, and Pedro Despouy. 2017. *Plotly: Create Interactive Web Graphics via 'Plotly.js'*.
- Xie, Yihui. 2014. "Knitr: A Comprehensive Tool for Reproducible Research in R." in *Implementing reproducible computational research*, edited by V. Stodden, F. Leisch, and R. D. Peng. Chapman; Hall/CRC.
- Xie, Yihui. 2015. *Dynamic Documents with R and Knitr*. 2nd ed. Boca Raton, Florida: Chapman; Hall/CRC.
- Xie, Yihui. 2016. *Bookdown: Authoring Books and Technical Documents with R Markdown*. Boca Raton, Florida: Chapman; Hall/CRC.
- Xie, Yihui. 2017. *Bookdown: Authoring Books and Technical Documents with R Markdown*.
- Xie, Yihui. 2018a. *Knitr: A General-Purpose Package for Dynamic Report Generation in R*.
- Xie, Yihui. 2018b. *Tinytex: Helper Functions to Install and Maintain 'Tex Live', and Compile 'Latex' Documents*.
- Zhu, Hao. 2017. *KableExtra: Construct Complex Table with 'Kable' and Pipe Syntax*.