

# Introduction

## Analytically reproducible documents

Hélène Langet

2025-03-03

## Table of contents

<b>1 Research = a dynamic process</b>	<b>2</b>
<b>2 Reproducibility in research</b>	<b>2</b>
<b>3 Analytically reproducible documents</b>	<b>2</b>
3.1 Formatted text . . . . .	4
3.1.1 HTML . . . . .	4
3.1.2 LaTeX . . . . .	4
3.1.3 Markdown . . . . .	4
3.2 Code commands . . . . .	4
3.2.1 R . . . . .	4
3.2.2 Python . . . . .	5
3.2.3 Observable JS . . . . .	5
3.3 Code outputs . . . . .	5
<b>4 Existing tools for writing analytically reproducible documents</b>	<b>8</b>
<b>5 Quarto</b>	<b>8</b>
5.1 Source document . . . . .	9
5.2 Rendered output . . . . .	9
5.3 Quarto rendered outputs . . . . .	9
5.4 Engines . . . . .	10
5.4.1 knitr . . . . .	10
5.4.2 Jupyter . . . . .	11
<b>6 To go further on reproducibility</b>	<b>11</b>

## 1 Research = a dynamic process

- Research insights are typically disseminated through reports (e.g., scientific presentations, publications, etc), including a textual narrative detailing the research context, methods, and key findings, often supplemented with figures and tables to summarize results, and a final discussion, with findings serving as evidence to support conclusions and recommendations ;
- Research is an **iterative** and **dynamic** process, meaning there are no **final** or **definitive** results or reports ;
- In addition, we continuously build upon the work of others to generate new insights and discoveries.

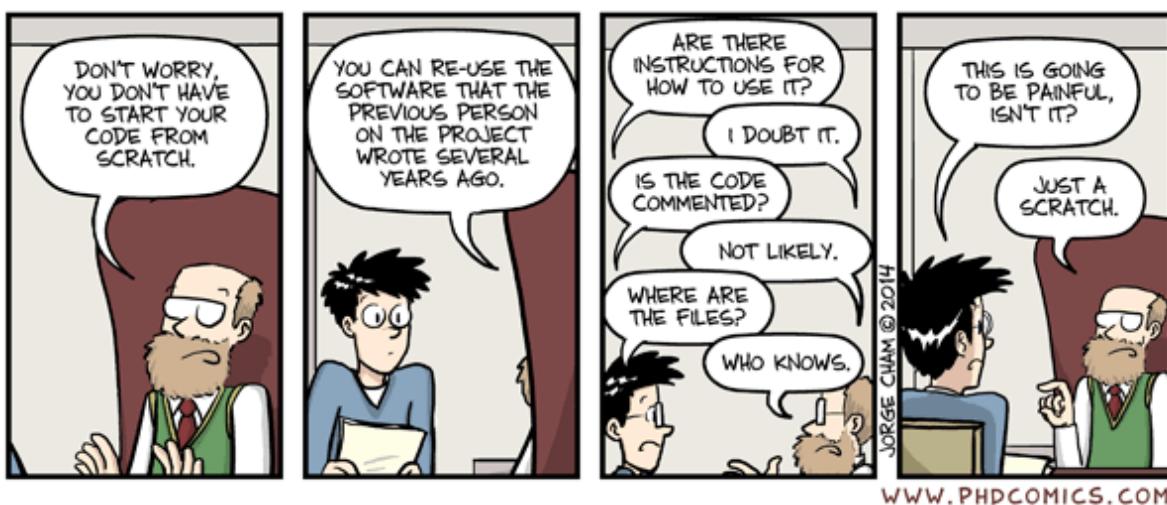


Figure 1: Image reproduced from Jorge Cham at [PhDComics](http://www.phdcomics.com).

## 2 Reproducibility in research

Achieving reproducibility requires clear access to the underlying **data**, the **code** used for analysis, and the **results** produced. It also depends on documenting the **tools**, such as software and libraries, alongside the **computational environment**, including hardware configurations and operating systems (1).

## 3 Analytically reproducible documents

Analytically reproducible documents typically contain 3 main types of content, integrating code and natural language in a way that is called “*literate programming*” (2).

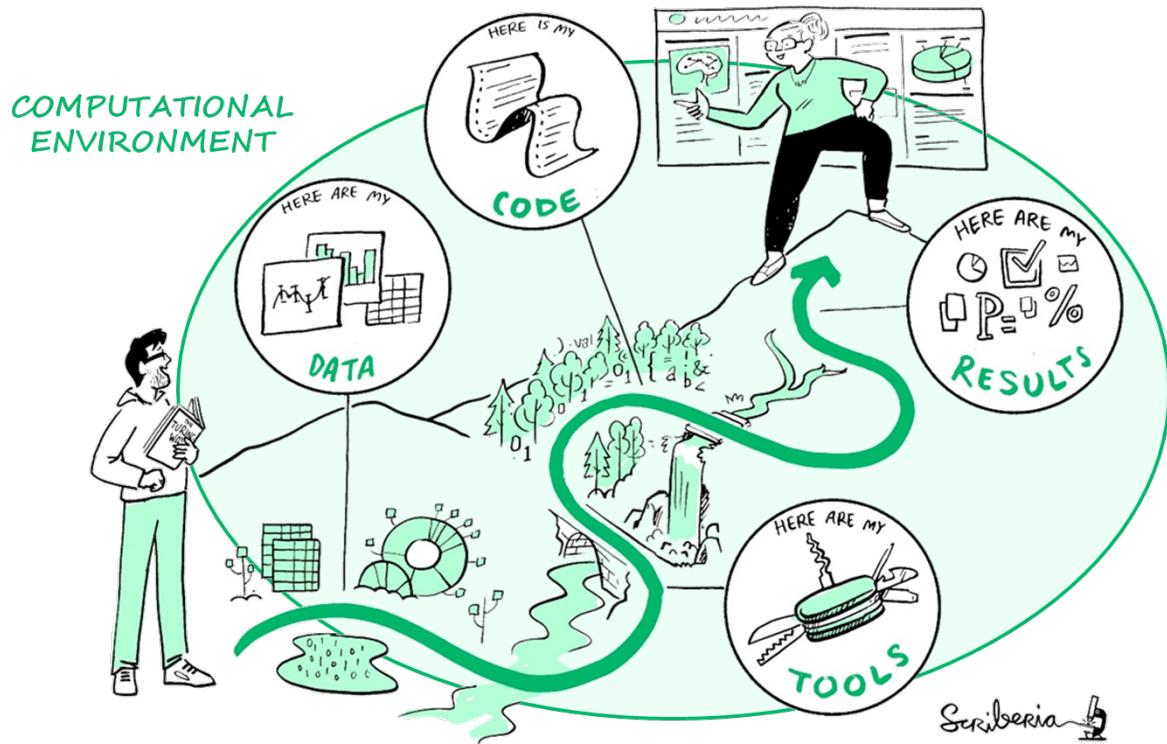


Figure 2: Illustration highlighting the key components of reproducibility in research, including data, code, results, tools, and the computational environment. This illustration is adapted from the one created by Scriberia with The Turing Way community, and used under a CC-BY 4.0 licence. DOI: [10.5281/zenodo.3332807](https://doi.org/10.5281/zenodo.3332807).

### 3.1 Formatted text

These are languages that can be written using any **plain text editor**. They use *markup* elements to define how text should be displayed or printed.

#### 3.1.1 HTML

HTML is used to structure content on the web.

```
<b>This text will be displayed in bold</b>
```

#### 3.1.2 LaTeX

LaTeX is used for academic and technical documents.

```
\textbf{This text will be displayed in bold}
```

#### 3.1.3 Markdown

Markdown is a lightweight markup language.

```
**This text will be displayed in bold**
```

### 3.2 Code commands

Different programming languages allow us to execute code to generate results or perform tasks.

#### 3.2.1 R

```
```{r}
library(ggplot2)
data.frame(country=c("Nigeria", "Kenya", "India"), prevalence=c(14.5, 9.2, 3.5)) |>
  ggplot(aes(x=country, y=prevalence)) +
  geom_bar(stat="identity", fill="steelblue")
```
```

### 3.2.2 Python

```
```{r}
library(reticulate)
Sys.setenv(RETICULATE_PYTHON = "C:/ProgramData/anaconda3/python.exe")
```

```{python}
import matplotlib.pyplot as plt
plt.bar(['Nigeria', 'Kenya', 'India'], [14.5, 9.2, 3.5], color='steelblue')
plt.show()
````
```

### 3.2.3 Observable JS

```
```{ojs}
BarChart({x: ["Nigeria", "Kenya", "India"], y: [14.5, 9.2, 3.5], yLabel: "Prevalence (%)"}
````
```

## 3.3 Code outputs

The output from executing code often results in visualizations or printed results. Below are the corresponding outputs for each language:

Call: `glm(formula = confirmed ~ age, family = binomial, data = df)`

Coefficients:

| (Intercept) | age      |
|-------------|----------|
| 1.312275    | 0.001292 |

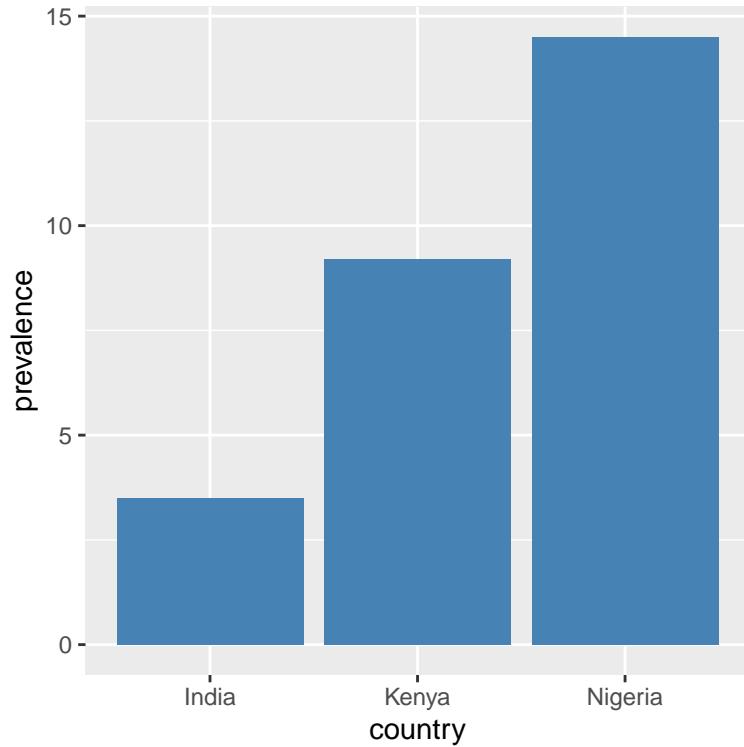
Degrees of Freedom: 65668 Total (i.e. Null); 65667 Residual

Null Deviance: 66000

Residual Deviance: 66000 AIC: 66000

|        | c1  | c2 |
|--------|-----|----|
| setosa | 5.1 |    |
| setosa | 4.9 |    |
| setosa | 4.7 |    |

|        | c1  | c2 |
|--------|-----|----|
| setosa | 4.6 |    |
| setosa | 5.0 |    |



Documents can be rendered in different type of outputs e.g., MS Word, PDF, HTML, Power-Point, etc.

**Pandoc** is a powerful open-source tool that enables seamless **conversion** between various document formats, making it an essential resource for working with markup languages like Markdown, HTML, and LaTeX. By using Pandoc, you can easily transform a document written in one markup language into a wide range of output formats, including MS Word, PDF, HTML, PowerPoint, and more, without needing to manually adjust formatting or structure.

This versatility makes Pandoc particularly valuable in workflows that involve literate programming. Pandoc also supports templates and extensions, allowing users to customize the output to meet specific stylistic or formatting requirements, simplifying the process of producing polished, professional documents.

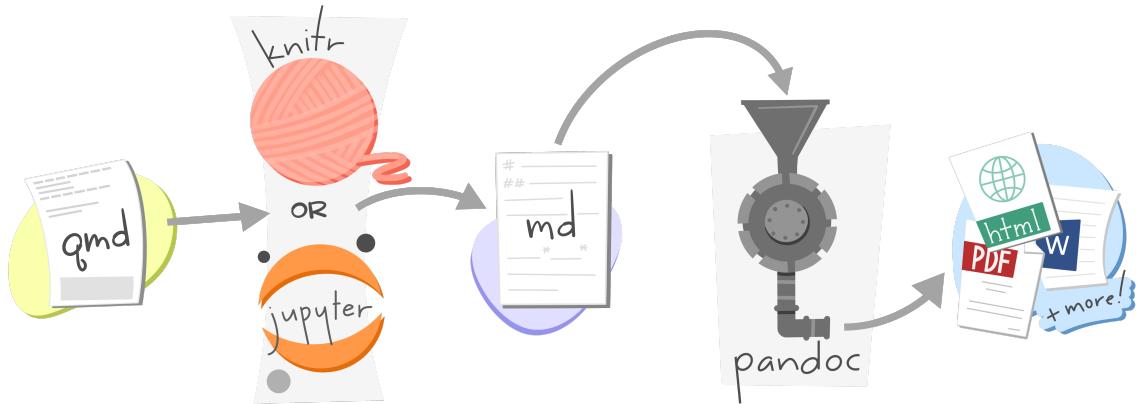


Figure 3: Process of transforming a Quarto document from its source format to the final rendered output. Artwork by Allison Horst.

| Language | Commands / File types   |
|----------|---|
| Stata    |  <p><b>dyndoc command</b><br/> putpdf command<br/> putdocx command<br/> putwrap command<br/> stmd command<br/> markstat command</p> <p>Limited capacity compared to Quarto</p> |
| R        |  <p>R Markdown (R + Markdown)<br/> R Sweave (R + LateX)<br/> R HTML (R + HTML)<br/> <b>Quarto (*.qmd)</b></p>  |
| Python   |  <p>Jupyter Notebook<br/> <b>Quarto (*.qmd)</b></p>  |

Figure 4: Existing tools for writing analytically reproducible documents

## 4 Existing tools for writing analytically reproducible documents

### 5 Quarto

- Quarto is the successor to R Markdown, but is not tied to the R language.
- Quarto files have a .qmd extension.

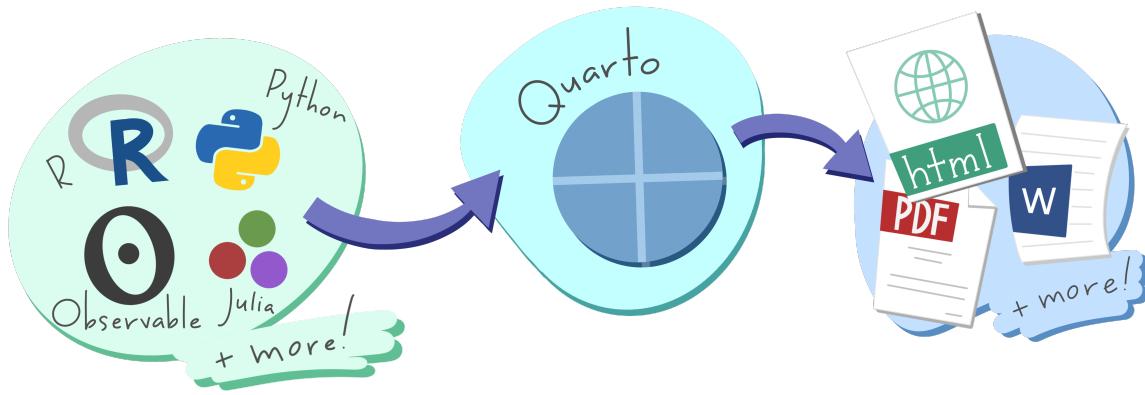


Figure 5: Artwork by Allison Horst.

## 5.1 Source document

The screenshot shows a code editor with R Markdown syntax. The code includes a title chunk, a note about running in notebook mode, and an R code chunk that prints the first five rows of the iris dataset using knitr::kable(). The output is a table with columns: Sepal.Length, Sepal.Width, Petal.Length, Petal.Width, and Species.

```
10+ ### Task 2
11
12 Display data. To run the code in *notebook* mode, click the green
   arrow on the right side of the code chunk.
13
14+ `r`
15 df <- iris
16 df |>
17 head(5) |>
18 knitr::kable()
19+ ```



Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa


```

## 5.2 Rendered output

The screenshot shows the rendered output of the R Markdown code. It includes a title chunk, a note about running in notebook mode, and an R code chunk that prints the first five rows of the iris dataset using knitr::kable(). The output is a table with columns: Sepal.Length, Sepal.Width, Petal.Length, Petal.Width, and Species.

**Task 2**

Display data. To run the code in `notebook` mode, click the green arrow on the right side of the code chunk.

```
df <- iris
df |>
  head(5) |>
  knitr::kable()
```

| Sepal.Length | Sepal.Width | Petal.Length | Petal.Width | Species |
|--------------|-------------|--------------|-------------|---------|
| 5.1          | 3.5         | 1.4          | 0.2         | setosa  |
| 4.9          | 3.0         | 1.4          | 0.2         | setosa  |
| 4.7          | 3.2         | 1.3          | 0.2         | setosa  |
| 4.6          | 3.1         | 1.5          | 0.2         | setosa  |
| 5.0          | 3.6         | 1.4          | 0.2         | setosa  |

## 5.3 Quarto rendered outputs

- Quarto documents can be rendered into many report formats including HTML, Word document and many more
- [List of supported formats](#)

The screenshot shows the Quarto Gallery page. At the top, there's a navigation bar with links for Overview, Get Started, Guide, Extensions, Reference, Gallery, Blog, and Help. Below the navigation, the word "Gallery" is prominently displayed in a large, bold font. A text block below it says: "Quarto can produce a wide variety of output formats. Here are some examples:" followed by a bulleted list: Articles & Reports, Presentations, Dashboards, Websites, Books, and Interactive Docs. To the right of this text is a data visualization titled "MPG and Horsepower". It's a scatter plot with a blue smoothing line. The plot area has a light gray background with a color scale bar on the right. Below the plot, there's a snippet of R code: 

```
library(ggplot2)
ggplot(mpg, aes(hp, mpg, color = am)) +
  geom_point() + geom_smooth(formula = y ~ x, method = "loess")
```

 and a link "Learn more: Executable Code". At the bottom right of the main content area is the Quarto logo.

## 5.4 Engines

An engine refers to the software or system that executes the embedded code within a document. The engine takes the code chunks written in a specific programming language (e.g., R, Python, or Julia), runs them, and returns the output, which is then incorporated into the rendered document.

Both knitr and Jupyter serve as engines to execute code embedded within a document, but they work in different programming environments.

### 5.4.1 knitr

This R package will read the code chunks, execute it, and ‘knit’ it back into the document. This is how tables and graphs are included alongside the text.



### 5.4.2 Jupyter

Jupyter is a popular engine for running Python code interactively. It supports multiple programming languages, but Python is the most common.



## 6 To go further on reproducibility

- MOOC Reproducible research I - Methodological principles for transparent science (in French)
  - MOOC Reproducible research II - Practices and tools for managing computations and data
1. The Turing Way Community. [The Turing Way: A handbook for reproducible, ethical and collaborative research \(1.0.2\)](https://book.the-turing-way.org/index.html). <https://book.the-turing-way.org/index.html>; 2022.
  2. Knuth DE. Literate Programming. The Computer Journal [Internet]. 1984 Jan;27(2):97–111. Available from: <https://doi.org/10.1093/comjnl/27.2.97>
  3. Batra, Neale and Mousset, Mathilde and Spina, Alex and Florence, Isaac and Coyer, Liza and others. [The Epidemiologist R Handbook](https://epirhandbook.com). <https://epirhandbook.com>; 2021.
  4. The Graph Courses Team. Websites and dashboards with Quarto. <https://thegraphcourses.org/courses/websites-and-dashboards-with-r/>; 2023.