# Creating targets-based scientific workflows

## Reproducible Scientific Workflows in R - Part 2

Ernest Guevarra

2024-02-12

# Outline

- Concepts on scientific workflows

- The `{targets}` package

- Practical session

# Concepts on scientific workflows

# Concept #1: Reproducibility, reproducibility, reproducibility!

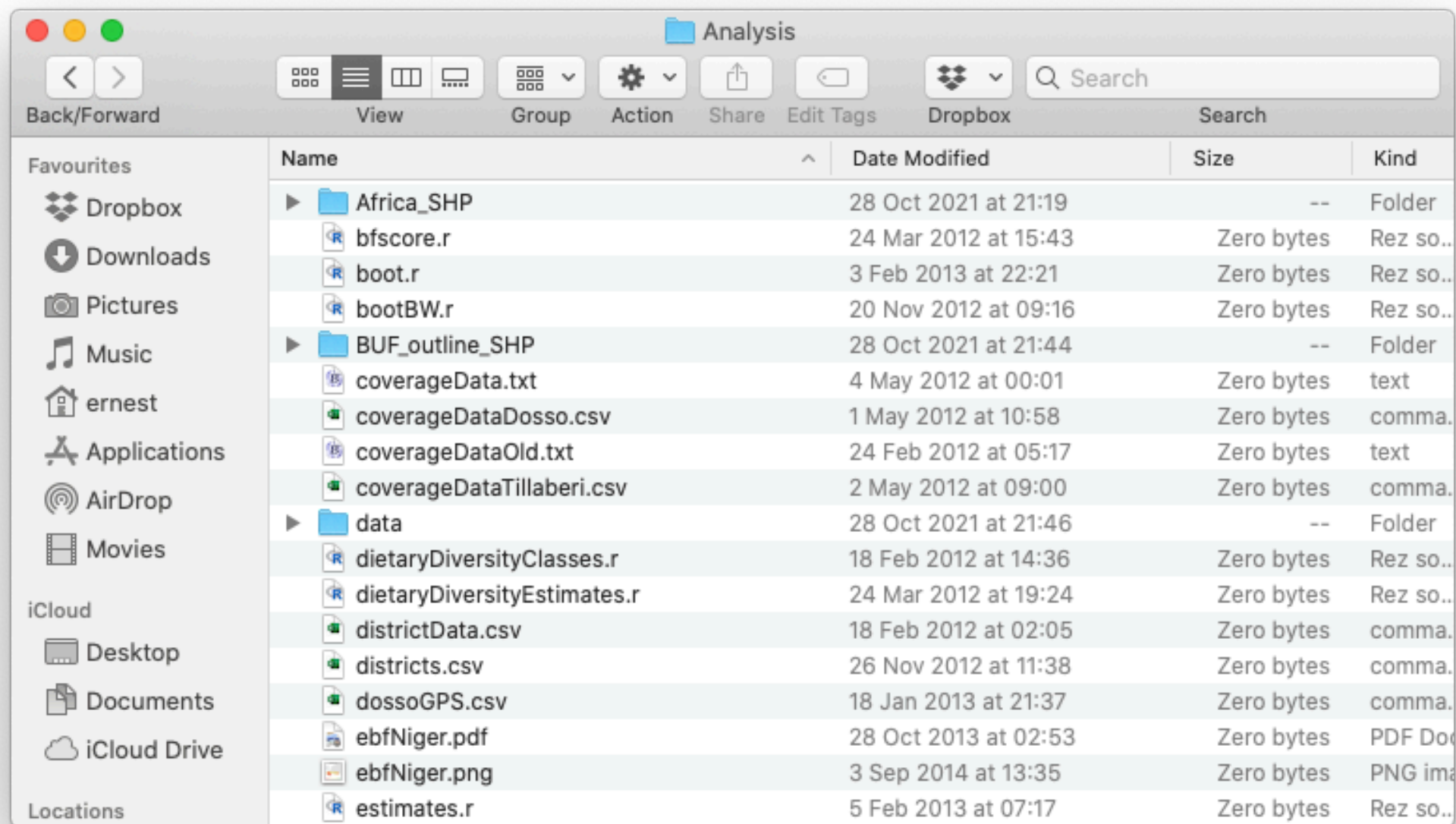# Most important tool for Reproducible Research is the mindset, when starting, taht the end product will be reproducible.

Keith Baggerly, via @kwbroman tweet

# Concept #2: Organisation

# File organization and naming are powerful weapons against chaos.

@JennyBryan

# Concept #3: DRY - Don't repeat yourself

# Don't repeat yourself. It's not only repetitive, it's redundant, and people have heard it before.

Lemony Snicket

```r
# Overlay maps of Niger and Nigeria to clean-up the map
par(new=TRUE)
plot(nigeria, axes = FALSE, xlim = mapXLimits, ylim = mapYLimits, border = "white", col = "white")

par(new = TRUE)
plot(boundaries, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.5, border = "black")

par(new = TRUE)
plot(n1, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n4, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n5, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n6, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n6.27, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n7, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")

par(new = TRUE)
plot(n14, axes = FALSE, xlim = mapXLimits, ylim=mapYLimits, lwd = 0.25, col = "blue")
```
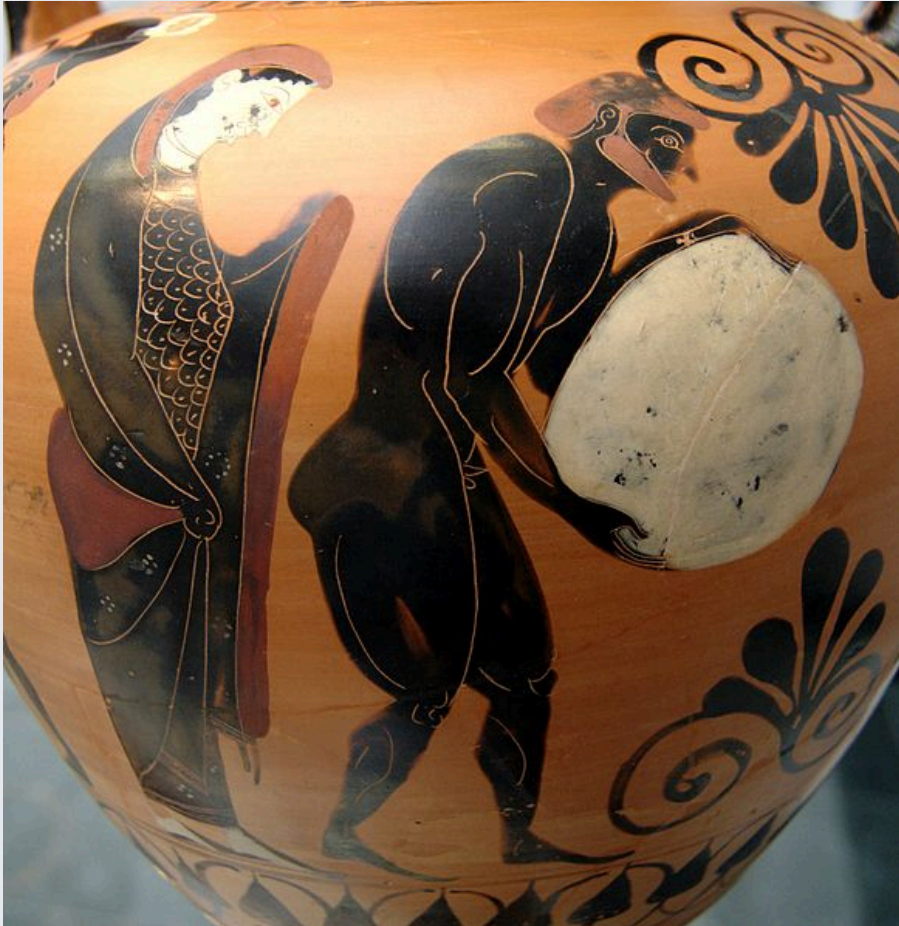
# Sisyphean loop



1. Launch the code.
2. Wait while it runs.
3. Discover an issue.
4. Restart from scratch.

# Concept #3: Frequency reduces difficulty

# If it hurts, do it more often.

@martinfowler, via @JennyBryan tweet

# The {targets} package



github.com/ropensci/targets

- a pipeline toolkit for Statistics and data science in R
- maintain a reproducible workflow without repeating yourself
- learns how your workflow fits together
- skips costly runtime for tasks that are already up-to-date
- runs only the necessary computation
- supports implicit parallel computing
- abstracts files as R objects
- shows tangible evidence that the results match the underlying code and data

# {targets} file organisation



- this is a typical file structure with user-defined components of any project-oriented workflow
- the `_targets.R` file, however, is special and specific to a {targets} workflow – it is the target script file
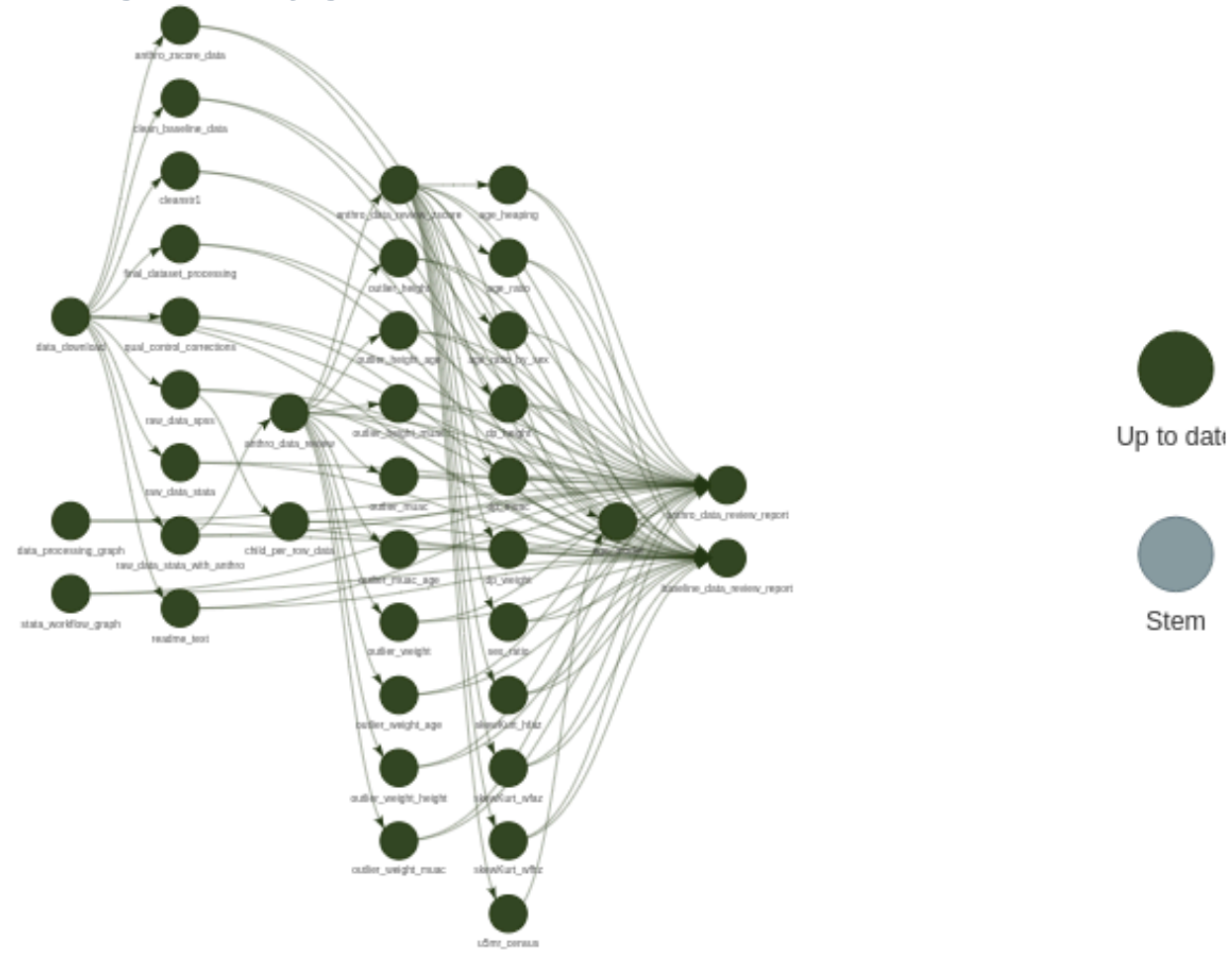- the target script file should be in the project's root directory

# {targets} script file



- Load the packages required (1)
- Load custom functions (1)
- Define individual targets - intermediate step of the workflow (2)
- End with a list of targets objects

# {targets} workflow

# Questions?

# Practical session

We will all continue to go through Exercise #1 in the Practical R for Epidemiologists book

# Questions?

# Thank you!

Slides can be viewed at https://oxford-ihtm.io/open-reproducible-science/session9.html

PDF version of slides can be downloaded at https://oxford-ihtm.io/open-reproducible-science/pdf/session9-reproducible-scientific-workflows.pdf

R scripts for slides available here