Workflow for the Week

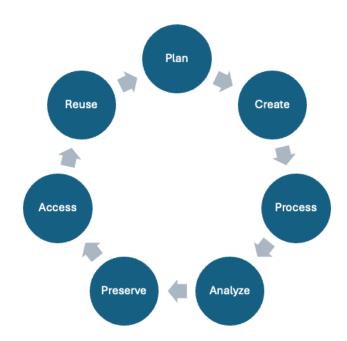
Nick Rochlin, University of Victoria RDM Jumpstart May 12, 2025

Key Tools & Concepts

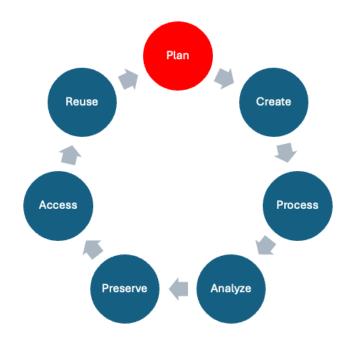
Overview

- Before collecting data
 - O Data management planning
 - O Open Science Framework (OSF)
- After collecting data
 - O OSF
 - O R, RStudio, & Rmarkdown
- At project completion
 - O Borealis/Dataverse

The Research Data Lifecycle

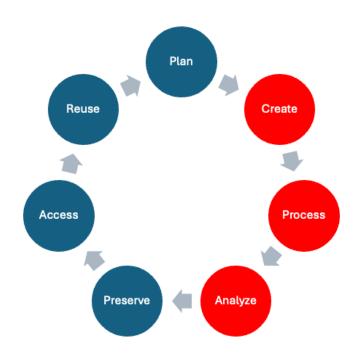


The Research Data Lifecycle

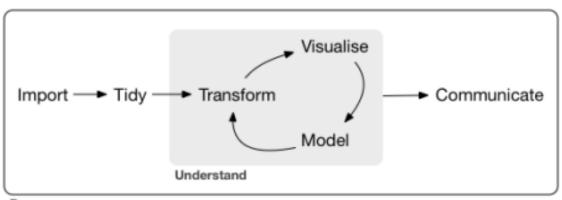


Data Management Plans - Covered in day 1

- Questions and considerations about your research data prior to starting a project.
- Will be covered in more detail this afternoon!



- A framework to support transparent and open research processes.
- A granular roadmap to support the "active" stages of the research data lifecycle.
- While the research data lifecycle is focused on making research processes and outputs transparent and reproducible, the data science workflow focuses on manipulating and asking questions from data in an open and reproducible way.



R

• An open-source programming language designed for data manipulation, statistical analysis, and visualization.

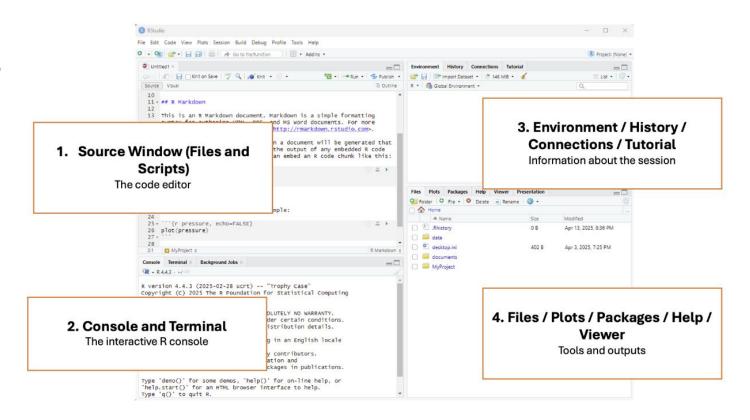
RStudio

• While R is the programming language, RStudio is an Integrated Development Environment (IDE), which is a fancy term for a software application that makes writing scripts R much more user-friendly.

RStudio

- A analogy:
 - O You are writing a document in Microsoft Word. In this case, the plain human language that you are typing would be R, and Microsoft Word, which is the software that allows you to view, format, add visuals, and save your work, is RStudio.
- In the morning of day 2 we'll do a session on using Excel vs. a coding language like R, and then we'll jump into R and RStudio in the afternoon.

RStudio

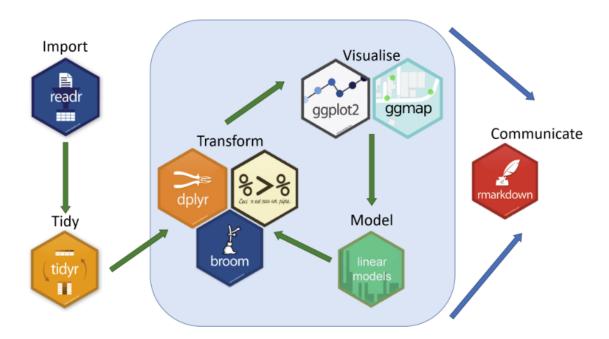


Across the Tidyverse

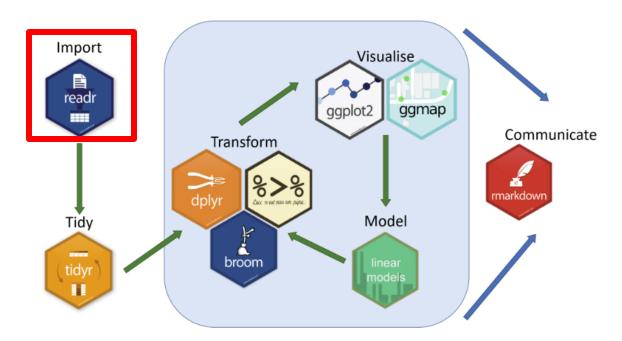
The Tidyverse

- A group of R "packages" for data manipulation, exploration, and visualization.
- Focused on the connections between the activities of reproducible workflows.
- Provides a good reference point for those new to coding to engage with coding.

The Data Science Workflow with Tidyverse



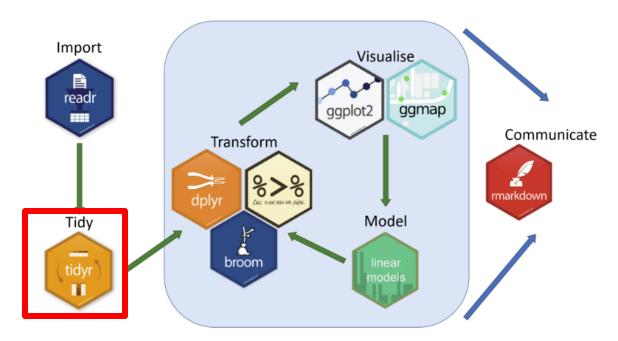
Import – Covered in day 2



Import – Covered in day 2

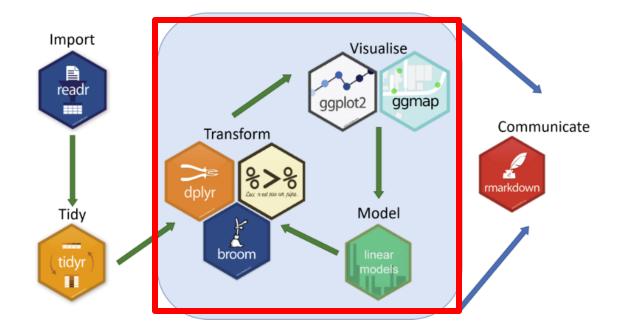
Unlike Excel, which can immediately start working with spreadsheet data, for R
to be able to work with data, the data first needs to be brought in, or imported,
to the coding environment.

Tidy - Covered in day 3



Tidy – Covered in day 3

- Most of the data we collect or retrieve is quite messy, and needs to be cleaned up, or tidied, before we can start asking questions from it.
- This can include things like:
 - Changing the names of columns
 - Accounting for missing values / missing data
 - Re-formatting data
 - 0 ..



Transform – Covered in day 3

- Once data has been tidied, we can start playing around!
- Data transformations, such as:
 - Filtering rows and columns that meet certain criteria
 - Creating new columns based on the values of existing columns
 - Subsetting specific portions of the data
- Transforming data can be a good way to set up data visualizations or modeling, but can also be conducted after visualization and modeling if you notice something interesting in the data.

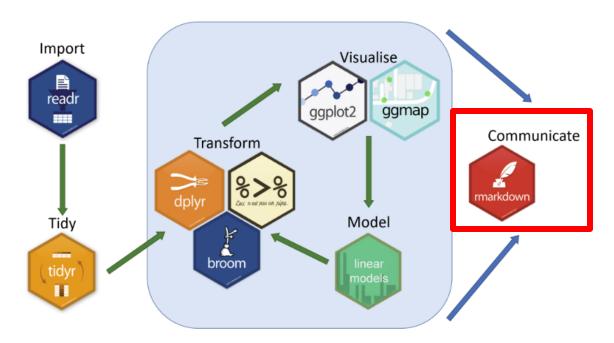
Visualize - Covered in day 4

- Excel can be pretty good at creating visualizations quickly, but it can also be frustrating due to rigid structures and it not doing what you want.
- Visualizations in R require a bit more upfront work, but have much more flexibility and power in what you can create.
- Visualizations can be used to communicate your findings, but can also be a valuable tool in exploring data and discovering interesting trends.

Model - Out of scope

- Applying statistical models to data is the real "science" of data science, but is a bit beyond the scope of this series.
- However, it's worth mentioning here as a key part of the data science workflow, and we encourage you to use the skills developed throughout the week as a foundation to start entering this area.
- It's also possible that we may build out a more advanced level of the series, and if that's of interest, please let us know!

Communicate - Covered in day 2-5



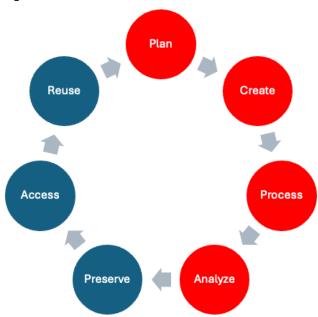
Communicate – Covered in day 2-5

 Communicating your findings is the most fulfilling part of the data science workflow (it's finally over!), but to do this in an open and transparent way, it needs to be considered throughout the entire research life cycle.

Communicate – Covered in day 2-5

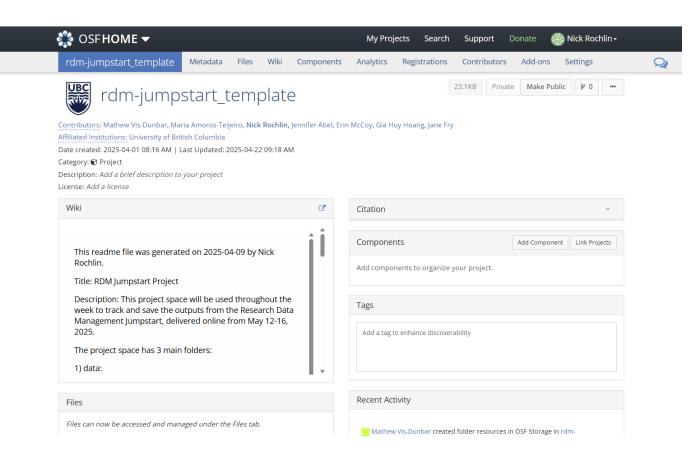
- Rmarkdown is a plain text format that combines R code, results, and explanatory text into a single, shareable, reproducible document.
- Allows you to create various file formats like PDF, Word, HTML, and presentations.

Storage & Backup

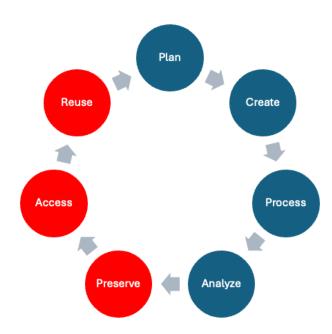


Open Science Framework (OSF) - Whole week

- Cloud-based platform designed to promote transparency, collaboration, and reproducibility in research.
- For the purpose of the Jumpstart, it will be used to:
 - Version control documents
 - Backup and share working documents
- OSF will be introduced in further detail this afternoon and used throughout the week.

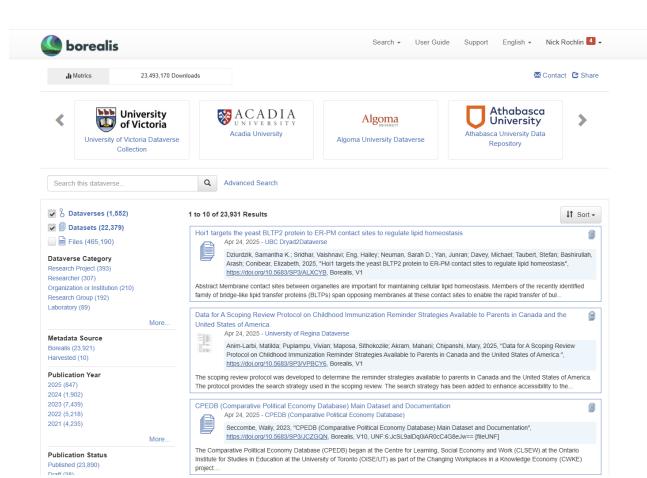


Data Preservation & Sharing



Borealis / Dataverse - Covered in day 5

- A Canadian data repository supported by academic libraries and institutions across Canada.
- Facilitates the sharing, preservation, and discovery of Canadian research data.
- Throughout the week we'll be discussing considerations for sharing and preserving data, and on day 5 we'll cover Borealis and will deposit materials from the week into the repository.



Questions?