

Git and GitHub for Reproducible Data Analysis

Data Standards and Open Data Community of Practice

Alice Byers
Data Innovation Team, Scottish Government

27 April 2023

Slides: https://github.com/alicebyers5/presentations/tree/main/2023-04-27_version-control

Aims

- What is version control?
- Introduction to Git and GitHub
- How Git and GitHub can be used for reproducible data analysis
- Note: This aim of this presentation is to introduce concepts and will be non-technical. If you would like technical support, please get in touch. (Contact details are at the end of the slides.)

What is version control?

Version control is the practice of tracking and managing changes to files.

Does this look familiar?

```
|— stats-publication
|   |— publication-analysis-code.R
|   |— publication-analysis-code-v2.R
|   |— publication-analysis-code-v2 NEW METHODOLOGY.R
|   |— publication-analysis-code-v2 NEW METHODOLOGY AB-changes.R
|   |— publication-analysis-code-final.R
|   |— publication-analysis-code-final April 2023.R
```

Git



Git is a free and open source software for version control.

To use:

- Install Git on your computer
- Initiate Git in a project folder (also called a repository)
- Record any changes you make to files (these records are called 'commits')
- Undo changes and revert to previous version of files (if required)
- Collaborate with others on the same project using branches

Git is used via a terminal-like tool called Git Bash, or via RStudio.

Using Git means you don't need to save multiple copies of the same file to retain older versions. This information is stored by Git.

Commits contain information on:

- **what** change was made,
- **when** the change was made,
- **why** the change was made, and
- **who** made the change.

Git tracks changes to the content of files, not just the file as a whole. This means the information above can be recorded for changes as small as one character on one line of code.

What is in a code repository?

A version controlled code repository will usually contain files for one project and can contain:

- Code (e.g. R scripts)
- Documentation (e.g. README)
- Configuration files
- ...but **NOT DATA!**

Data can be stored within your code repository, but it should not be tracked by Git. To ensure data files (and any other files containing sensitive information) are not tracked, a **.gitignore** file can be used to exclude them. Alternatively, data can be stored outside of your code repository.

More information on using Git safely can be found in the [Duck Book](#).

GitHub



GitHub is a web interface for hosting version controlled code. GitHub is owned by Microsoft but is mostly free to use, although some extra features are available for a fee. Rivals also exist, e.g. GitLab, Bitbucket.

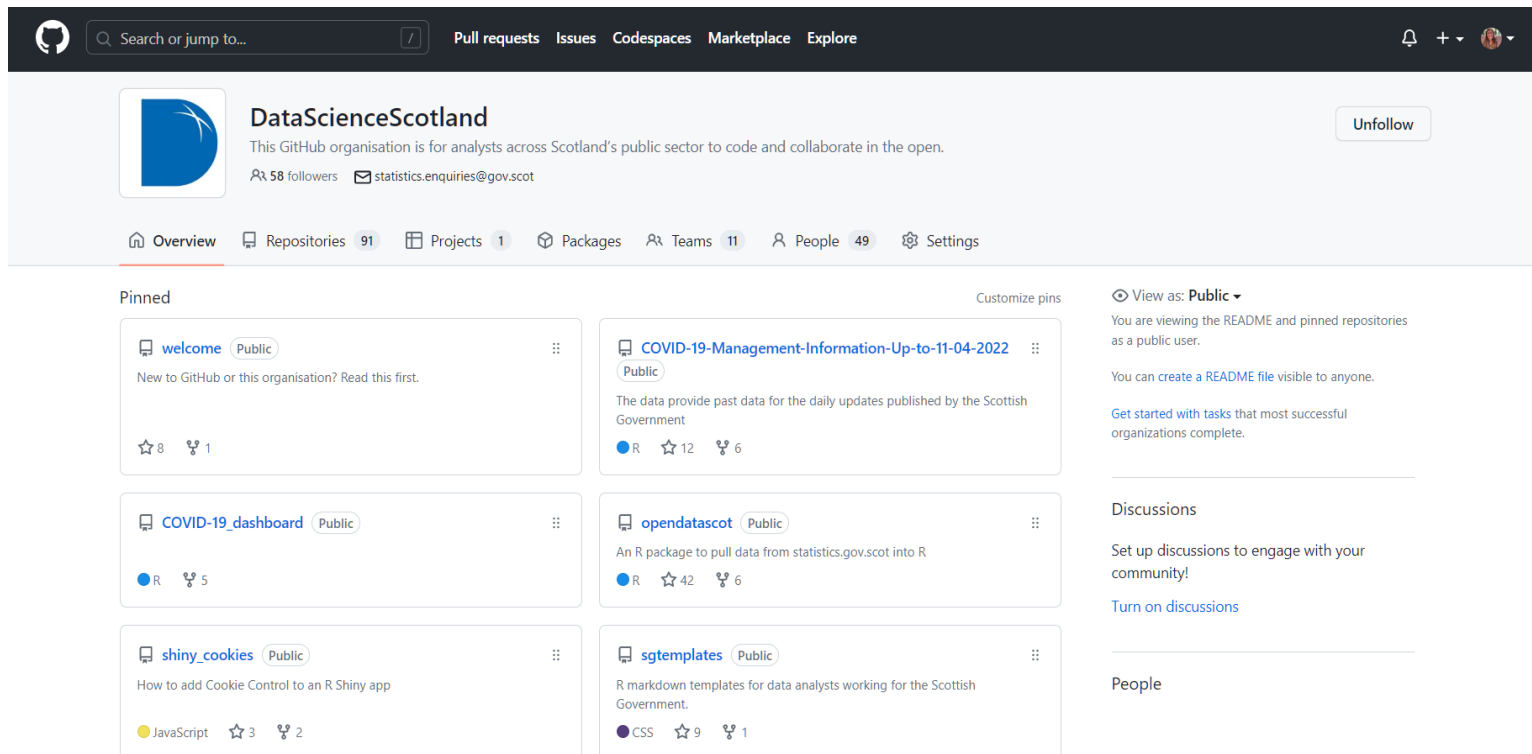
GitHub can be used to:

- Make code publicly available (although repositories can also be private)
- Facilitate code review (using 'pull requests')
- Manage projects using tools such as issue tracking
- Navigate Git history and view previous versions of files
- View other people's code and collaborate

GitHub Organisation



<https://github.com/DataScienceScotland>



The screenshot shows the GitHub organization page for DataScienceScotland. The header includes the GitHub logo, a search bar, and navigation links for Pull requests, Issues, Codespaces, Marketplace, and Explore. The organization's profile section displays the DataScienceScotland logo, a description, 58 followers, and an email address. Below this is a navigation bar with tabs for Overview, Repositories (91), Projects (1), Packages, Teams (11), People (49), and Settings. The main content area is titled 'Pinned' and shows a grid of six pinned repositories. On the right, there are sections for 'View as: Public', 'Discussions', and 'People'.

DataScienceScotland
This GitHub organisation is for analysts across Scotland's public sector to code and collaborate in the open.
58 followers | statistics.enquiries@gov.scot

Unfollow

Overview Repositories 91 Projects 1 Packages Teams 11 People 49 Settings

Pinned Customize pins

welcome (Public)
New to GitHub or this organisation? Read this first.
8 stars 1 fork

COVID-19-Management-Information-Up-to-11-04-2022 (Public)
The data provide past data for the daily updates published by the Scottish Government
R 12 stars 6 forks

COVID-19_dashboard (Public)
R 5 stars

opendatascot (Public)
An R package to pull data from statistics.gov.scot into R
R 42 stars 6 forks

shiny_cookies (Public)
How to add Cookie Control to an R Shiny app
JavaScript 3 stars 2 forks

sgtemplates (Public)
R markdown templates for data analysts working for the Scottish Government.
CSS 9 stars 1 fork

View as: Public
You are viewing the README and pinned repositories as a public user.
You can [create a README file](#) visible to anyone.
[Get started with tasks](#) that most successful organizations complete.

Discussions
Set up discussions to engage with your community!
[Turn on discussions](#)

People

How to use Git and GitHub

- Git can be used without GitHub
- GitHub is often used as the main copy of a code repository (or 'remote'). Analysts or developers can take a copy (or 'clone') of the repository from GitHub to work on locally.
- Use Git locally to track changes and regularly 'push' to GitHub
- Use GitHub to facilitate code review and merging of branches

Why use Git and GitHub?

- Preferable to lots of copies of the same file with various names!
- [Reproducible Analytical Pipelines \(RAP\)](#)
 - **Reproducible:** You can rerun your code as it was at any point in time.
 - **Auditable:** You have a record of when changes were made and why.
 - **Transparent:** Code is publicly available on GitHub and available for others to review or reuse.
 - **Good quality:** Code review is built into the GitHub workflow.

Links and Resources

- Version control [Saltire pages](#) (accessible on Scottish Government SCOTS network only)
- [Data Science Scotland](#) GitHub Organisation
- Duck Book
 - [Version Control](#)
 - [Using Git safely](#)
 - [GitHub features](#)
- Government Analysis Function guidance on [open sourcing analytical code](#)
- ONS Learning Hub (contact Data.Science.Campus.Faculty@ons.gov.uk to request an account)
 - [Command Line Basics](#)
 - [Introduction to Git](#)

Contact

Alice Byers
RAP Developer
Data Innovation Team, Scottish Government

- Email - alice.byers@gov.scot
- GitHub - [alicebyers5](https://github.com/alicebyers5)