

EL4: Tooling

Positron and version control

Reading and practicing

- [1] introduces the Unix shell and the file system, which are essential parts for using git. Read and practice according to the instructions. (You may use the [Terminal tab in RStudio](#) or similarly in Posotron). Only the first sections are required:
 - [Introducing the shell](#)
 - [Navigating Files and Directories](#)
 - [Working With Files and Directories](#)
- [2, ch. 4] introduces the basics of Git for R users and applies to all common operating systems. Read and practice by following the instructions.

Recommended references:

- A bit old but still relevant: [Happy Git with R](#)
- T. L. Staples [3] illustrates the constant evolution in technology. As a professional, you should not expect that what you learn in this course will stay relevant for ever. This article is a good illustration of how the linguistic use of R has changed in less than a decade. It is also a practical example of how you can use GitHub in a slightly different way and it touches briefly on `{data.table}`, which we will introduce later. It also illustrates that a lot of R code is not written for statistical analysis (the last and final step) but for data management. The article also mentions some statistical techniques which you will meet in later courses (ignore the details for now).

Motivation

It is widely acknowledged that the **most fundamental developments in statistics** in the past 60 years are **driven by information technology (IT)**. We should not underestimate the importance of pen and paper as a form of IT but it is since people start using **computers to do statistical analysis** that we really changed the role statistics plays in our research as well as normal life.

Although: “Let’s not kid ourselves: the most widely used piece of software for statistics is Excel.” /Brian Ripley (2002)

[Short overview](#)

Mathematics



That textbook was written thousands of years ago, and it is still as useful and relevant as ever.



Physics



Oh, that textbook is outdated. It was written before Newtonian mechanics.



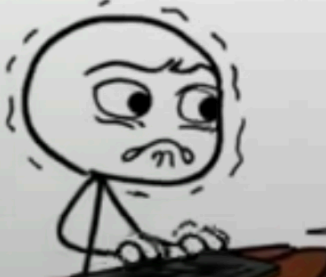
Chemistry



Oh, that textbook is obsolete. It was written before the electron was discovered.



Programing



Cries in programmer documentation becoming obsolete before they are even finished.



Panta rei!

- We teach you the present
- But your work is in the future
- We might use history to predict the future?
- At least we should learn that things changes constantly!

Terminal

- Don't rely too heavily on your computers graphical user interface (GUI)
- You know how to use the R console!
- To use a terminal/Unix shell/Bash... is similar
- To navigate the file system is basic practice
- A Unix Terminal is included if you are using Mac (and Linux)
- Emulators are common for Windows (one is included in Git, use that one after installation)

Early computer languages

- **ALGOL / ALGOL 60**
 - Algorithmic Language
 - Influential in academic/scientific computing
 - Introduced structured programming concepts that shaped later languages
- **PL/I**
 - Used in some government and industrial contexts
 - Combined scientific and business computing features (Fortran + Cobol)

General-Purpose Programming Languages

Early statistical computing relied heavily on:

- **FORTAN**
 - Dominant language for numerical and statistical computation
 - Statistical methods implemented as libraries and subroutines
 - Still used as numerical back-end (e.g., BLAS/LAPACK) in modern statistical software
- **C**
 - Emerged in the 1970s as a systems programming language
 - Widely used for implementing statistical software infrastructure
 - Core language of many modern statistical environments (e.g., R, parts of SAS, Stata)
 - Often used to interface with high-performance numerical libraries

📌 These languages required substantial programming expertise.

i FORTAN and C

FORTAN is still used in R for subroutines such as [least squares](#). Even in modern packages!

Same for C, such as for [lm](#), which is popular for high performance computing (fast execution time).

Early Statistical Packages

Several dedicated statistical systems emerged:

- **SPSS** (1968)
 - Originally batch-oriented
 - Widely used in social sciences
 - Originally “Statistical Package for the Social Sciences”
- **BMDP** (1960s)
 - Bio-Medical Data Package
 - Developed at UCLA
 - Common in medical statistics
- **GENSTAT** (1968)
 - Focused on agricultural statistics
- **MINITAB** (1972)
 - Designed for teaching and education
 - Still popular in quality control

SAS: A Transitional System

- **SAS** (early 1970s)
- Developed for agricultural and biostatistical analysis
- Script-based, but largely batch-oriented
- Became a standard in:
 - Government agencies
 - Large organizations
- Known for strong data management capabilities
- Still widely used in pharmaceutical industry

📌 SAS predates S and influenced later statistical workflows.

i SAS data

- Still common to receive data in SAS-format (`data_file.sas7bdat`) from register holders!
- Sometimes necessary to use SAS to export the data (genAI might be a good aid in those cases)

Limitations of Pre-S Systems

Common limitations included:

- Batch processing rather than interactivity
 - But we still sometimes need batch processing for large scale projects!
 - `Rscript script.R`
- Separation of data management and analysis
- Limited graphics capabilities
- High barriers to exploratory data analysis

S

- S takes form at Bell Laboratories (interactive statistical computing).
- John Chambers leads the effort.
- **1976**: first working version of S runs on GCOS
- **1979**: S2 is ported to UNIX; UNIX becomes the primary platform
- **1980**: S is first distributed outside Bell Labs
- **1981**: source versions are made available
- **1984**: key S books published (often called the “Brown Book” era)

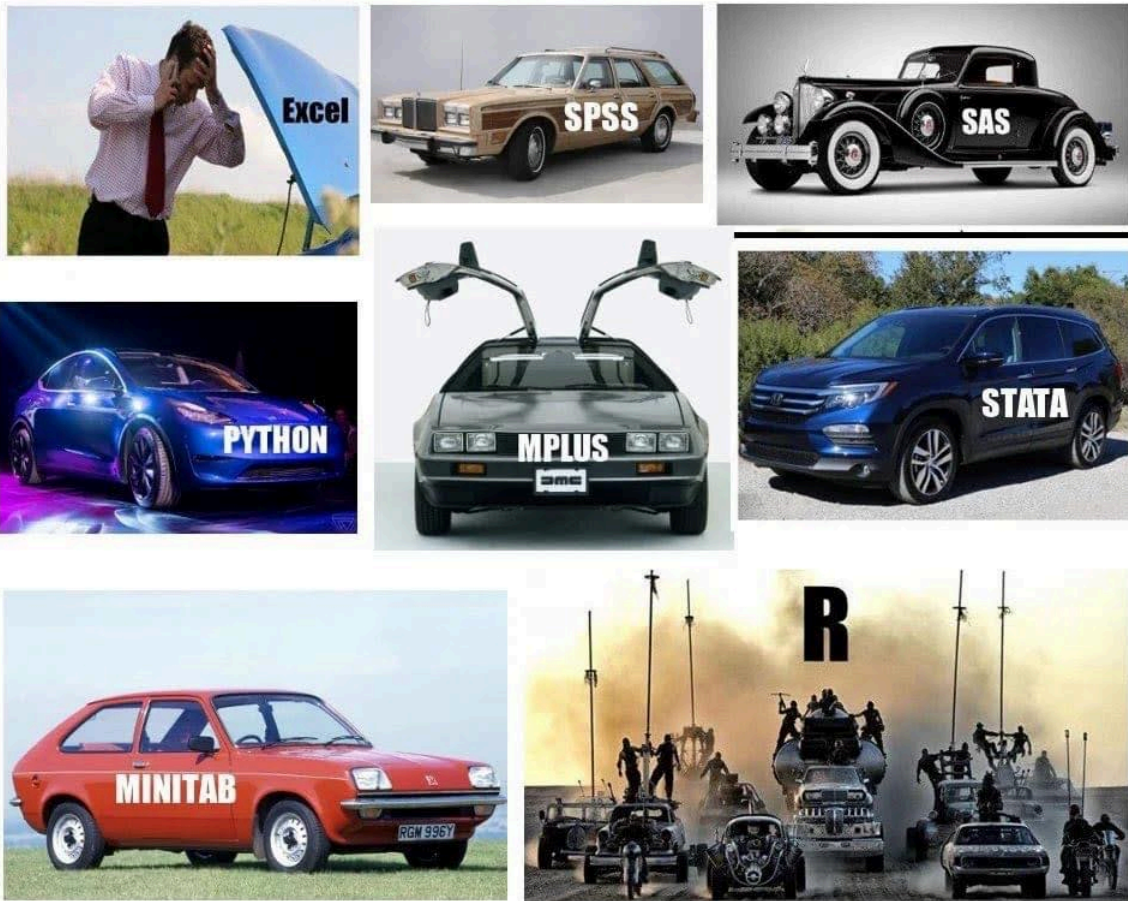
The New S Language

- **1988**: “New S” is released (major language redesign)
- **1988**: **S-PLUS** is first released as a commercial implementation of S
 - Last seen in Tibco Spotfire 2012
- **1991**: *Statistical Models in S* (“White Book”) popularizes formula notation (the ~ operator), data frames, and modeling workflows

R

- **1993**: first versions of **R** is published (Auckland; Ross Ihaka & Robert Gentleman)
 - **1995**: R becomes open source (GPL)
 - **1997**: the R Core group forms; **CRAN** is founded (Kurt Hornik)
 - **2000**: **R 1.0.0** is released 2000-02-29
-

If statistics programs/languages were cars...



RStudio brings an IDE to the R community

- **2009:** RStudio (the company) is founded
- **2011:** RStudio IDE is introduced as an open-source IDE for R (desktop + server)

Microsoft (Revolution Analytics)

- **Jan 2015:** Microsoft announces it will acquire **Revolution Analytics**
- Microsoft promotes enterprise R offerings (e.g., Microsoft R Open / R Server)
- **2016:** SQL Server 2016 introduces **R Services** (in-database R)
- **2017:** Microsoft expands the stack under “Machine Learning Server” branding
- **June 2021:** Microsoft announces retirement of **Microsoft Machine Learning Server**
- **2023:** Microsoft is still a relevant player
 - platinum member of [The R consortium](#)
 - Owner of GitHub
 - VS Code

RStudio becomes Posit

- **July 27, 2022:** RStudio rebrands as **Posit**
- **July 28, 2022:** **Quarto** is announced as a next-generation scientific and technical publishing system (multi-language, multi-engine)
- A public benefit company (not only relevant for its shareholders)
- Also a platinum member of The R consortium # Modern IDE

Positron

- New generation IDE for data science (and of course statistics!)
- From Posit PBC
- Free for individual use
- Based on Code OSS (open source version of VS Code from Microsoft)
- For both R and Python ([Julia?](#))

The screenshot displays the Positron IDE interface. The main editor shows a file named 'documentation' with the following content:

```
191 - Also a platinum member of The R consortium # secure envr
192 # Modern IDE
193
194 ## Positron
195
196 - New generation IDE for data science (and of course statistics!)
197 - From Posit PBC
198 - Free for individual use You, 4 weeks ago * EL2 - Fix #3
199 - Based on Code OSS (open source version of VS Code from Microsoft)
200 - For both R and Python (Julia?) (https://github.com/posit-dev/positron/issues/3672)
201
202 ## Quick tour
203
204 <iframe width="568" height="315" src="https://www.youtube.com/embed/4ir_0Mrihw?si=w9g-bo08748K0SUZ"
205 title="YouTube video player" frameborder="0" allow="accelerometer; autoplay; clipboard-write;
206 encrypted-media; gyroscope; picture-in-picture; web-share"
207 referrerpolicy="strict-origin-when-cross-origin" allowfullscreen>
208 </iframe>
209
210 ::: callout-important
211 # Positron assistant (mentioned in the video)
212
213
214 This feature will most likely be disabled in any secure working environment. Such environments often
215
216
```

The terminal window at the bottom shows the following output:

```
create mode 100644 images/paste-13.png
create mode 100644 images/paste-4.png
create mode 100644 images/paste-5.png
create mode 100644 images/paste-6.png
create mode 100644 images/paste-7.png
create mode 100644 images/paste-8.png
create mode 100644 images/paste-9.png
origin https://github.com/STA220/documentation.git (fetch)
origin https://github.com/STA220/documentation.git (push)
To https://github.com/STA220/documentation.git
202e44c..94f8ad9 HEAD -> gh-pages

NOTE: GitHub Pages sites use caching so you might need to click the refresh
button within your web browser to see changes after deployment.

[ ] Published to https://sta220.github.io/documentation/

NOTE: GitHub Pages deployments normally take a few minutes (your site updates
will be visible once the deploy completes)

> plot(10:10)
>
```

Quick tour

! Positron assistant (mentioned in the video)

This feature will most likely be disabled in any secure working environment. Such environments often have strict rules about data privacy and security, which may conflict with the assistant's functionality. Health data in SENSITIVE and SECURE environments must not be shared with external services, including AI assistants, to comply with data protection regulations and institutional policies.

It is recommended to not rely on such tools during the course (even if all our data is synthetic). If you start to rely on such tools, you might get difficulties the day you work with real data (might lead to prosecution for "brott mot tystnadsplikten" which is not only public, but actually civil law ("Brottsbalken") with prison sentence as a possibility). Society put an extreme emphasis on protecting health data, and rightfully so!

[Read more](#)

RStudio or Positron

- We will use Positron in this course
- This is for you to learn and see an alternative IDE ...
- ... and to get used to the fact that you should never stop learning!
- RStudio, however, is still a great tool and it is still developed and maintained.
- You may use both in the future (maybe Positron on your computer but RStudio in a secure server, which tend to be more slowly updated)

Some differences

- No package installer (yet). You need to use `pak::pkg_install()` or `install.packages()` etc.
- No inline rendering of results in Quarto documents (yet)
- You can use multiple active R sessions at once
- Great integrated tools from VS Code and extensions, such as GitHub integration

Version control

"FINAL".doc



FINAL.doc!



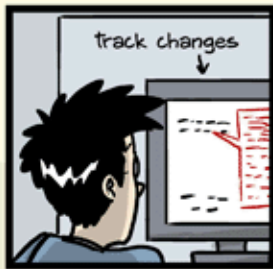
FINAL_rev.2.doc



FINAL_rev.6.COMMENTS.doc



FINAL_rev.8.comments5.
CORRECTIONS.doc



FINAL_rev.18.comments7.
corrections9.MORE.30.doc



FINAL_rev.22.comments49.
corrections.10.#@\$%WHYDID
ICOMETOGRADSCHOOL?????.doc



JORGE CHAM © 2012

Before Version Control Systems

1950s–1970s: Early software development relied on:

- Manual file naming:
 - analysis_final.f
 - analysis_final_v2.f
- Physical media:
 - [Punch cards](#)
 - Magnetic tape
- Centralised mainframes

✂ No automated tracking of changes.

Floppy discs, Mail, and Shared Directories

1970s–1980s: Common practices included:

- Copying files to:
 - Floppy disks
 - Magnetic tapes (actually still used for backups)
- Sending media by **postal mail**
- Sharing files via:
 - Network drives
 - FTP servers

✂ Version control was social, not technical.

Early Version Control Systems

1980s: First-generation tools focused on single files:

- **SCCS** (Source Code Control System, 1972)
- **RCS** (Revision Control System, 1982)

Characteristics:

- Versioning per file
- Linear history
- Central storage

Centralised Version Control

1990s: Project-level systems emerge:

- **CVS** (Concurrent Versions System)
- **Subversion (SVN)**

Key features:

- Central repository
- Multiple users

- Check-in / check-out model

✂ Still required constant access to the central server.

Limitations of Centralised Systems

Common problems:

- Single point of failure
- Poor support for branching and merging
- Difficult offline work
- Slow operations on large repositories

These limitations became critical for large projects.

Git

- **Git** was created by Linus Torvalds in 2005
- Original motivation:
 - Support Linux kernel development
 - Replace proprietary tools

Design principles:

- Distributed architecture
- Fast local operations
- Strong support for branching and merging

Terminal

- Windows: the Git installer comes with a terminal
 - Linux/Unix (incl. Mac): use your terminal/Warp
 - Terminal pane in RStudio/Positron
-

The image shows the Quarto IDE interface. The top bar displays the document name 'EL4.qmd - documentation' and the version 'R 4.5.2'. The left sidebar contains an Explorer and Documentation pane. The main editor shows a document with sections on 'The New S Language', 'R', and 'Early Version Control Systems'. The terminal window at the bottom displays system statistics for Darwin 26.3 64bit, including CPU usage, memory, and network activity.

Terminal Output:

```

n1-238911-086 Darwin 26.3 64bit
Frequency 3.70/3.78GHz  user system  idle  nice  MEM  56.3%  SWAP  0.0%  LOAD  12core
CPU0 [ 46.2%]  44.2%  18.0%  39.2%  0.0%  total  64.0k  total  0  1 min  4.67
CPU0 [ 46.0%]  35.0%  11.8%  53.2%  0.0%  used  36.0k  used  0  5 min  3.83
CPU0 [ 37.9%]  28.1%  17.8%  62.1%  0.0%  free  1.44k  free  0  15 min  6.31
CPU1 [ 53.0%]  37.2%  15.0%  57.0%  0.0%
CPU2 [ 24.3%]  29.1%  15.4%  55.5%  0.0%
MEM [ 56.2%]
LOAD [ 51.0%]

NETWORK
en0  Rx/s Tx/s
en0  0b 0b
en1  0b 0b
en2  0b 0b
en3  0b 0b
en4  0b 0b
en5  0b 0b
en6  0b 0b
en7  0b 0b
en8  14kb 29kb

DISK I/O  R/s W/s
disk0  1.90k 4.94k
disk1  1.2 0.7
disk2  0.9 0.6

FILE SYS  Used Total
2025-02-20 12:39:32 ESTW
Quarto: 1.8.26
  
```

GUI

The screenshot displays the RStudio IDE interface with the following components:

- Source Control Panel (Left):** Shows a list of changes including 'EL4.qmd', 'EL6.qmd', 'ES1.html_site', 'ES1.pdf_site', 'search.json_site', 'paste-14.png images', and 'paste-15.png images'. A commit message field is visible with the text 'Message (#Enter to commit on "ma...'. The 'Commit' button is highlighted.
- Document Editor (Top):** Displays the 'EL4.qmd' document with a preview of the rendered content. The preview includes a heading 'GUI' and a diagram with nodes for 'Open source and free' and 'Public and private repository hosting'. The rendered content is:


```

      GUI

      Open source and free
      Public and private repository hosting

      Can work as a
      
```
- Console (Bottom):** Shows system statistics for 'm1-230911-006 Darwin 26.3 64bit'. The output includes:


```

      Frequency 3.70/3.70GHz
      CPU0 [ 36.8%]
      CPU1 [ 32.1%]
      CPU2 [ 29.4%]
      CPU3 [ 27.9%]
      CPU4 [ 16.7%]
      MEM [ 56.7%]
      LOAD [ 52.8%]
      
```

 A table of tasks is also shown:

NAME	THR	NT	S
0 R Positron Helper (GPU)	--type		
0 R Spotify Helper (Renderer)	--		
0 R Positron Helper (Renderer)	--		
0 R Positron Helper (Renderer)	--		
0 R Positron Helper (Renderer)	--		
0 R Python /opt/homebrew/bin/gla			
0 R Positron Helper (Renderer)	--		
0 R screencapture -pdin -z keybo			
0 R Grammarly Desktop			
0 R zotero			
0 R Spotify Helper --type=gpu-pr			
0 R Windows App			
0 R Microsoft Teams WebView Help			
0 R M5Teams			
- Graph Panel (Bottom Left):** Shows a commit history graph with nodes for 'EL4, EL5, ECS1, ES1, Erik Bülow' and 'main'. The graph indicates a sequence of commits and branches.
- Terminal (Bottom):** Shows the current shell prompt as 'main*' and the user 'stefvanbuuren/fimdbook#14'.



The basics

[Four short videos to watch](#)

[Cheat-sheet](#)

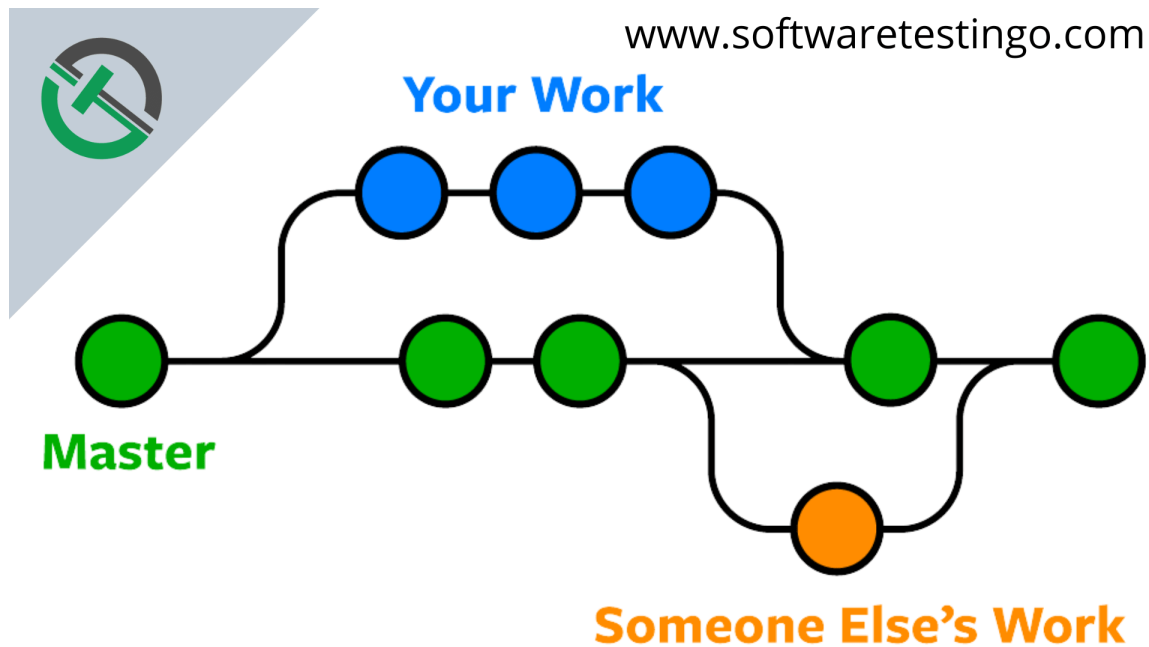
Some interactive learning tools:

- [Git practice](#)
- [Git Master](#)
- [learn git branching](#)
- [git simulator online](#)

Some basics

- You initiate a folder as a git project
- Git will track all changes made in that folder
 - New files
 - Modified files (especially text, such as programming scripts/functions etc)
 - Deleted files

Branching



Distributed Version Control with Git

Key ideas in Git:

- Every clone is a full repository (“folder”)
- Local commits without network access
- Cheap and fast branches
- Cryptographic integrity (hash-based)

🔗 Collaboration becomes more flexible and robust.

Hosting Platforms

Platforms built around Git:

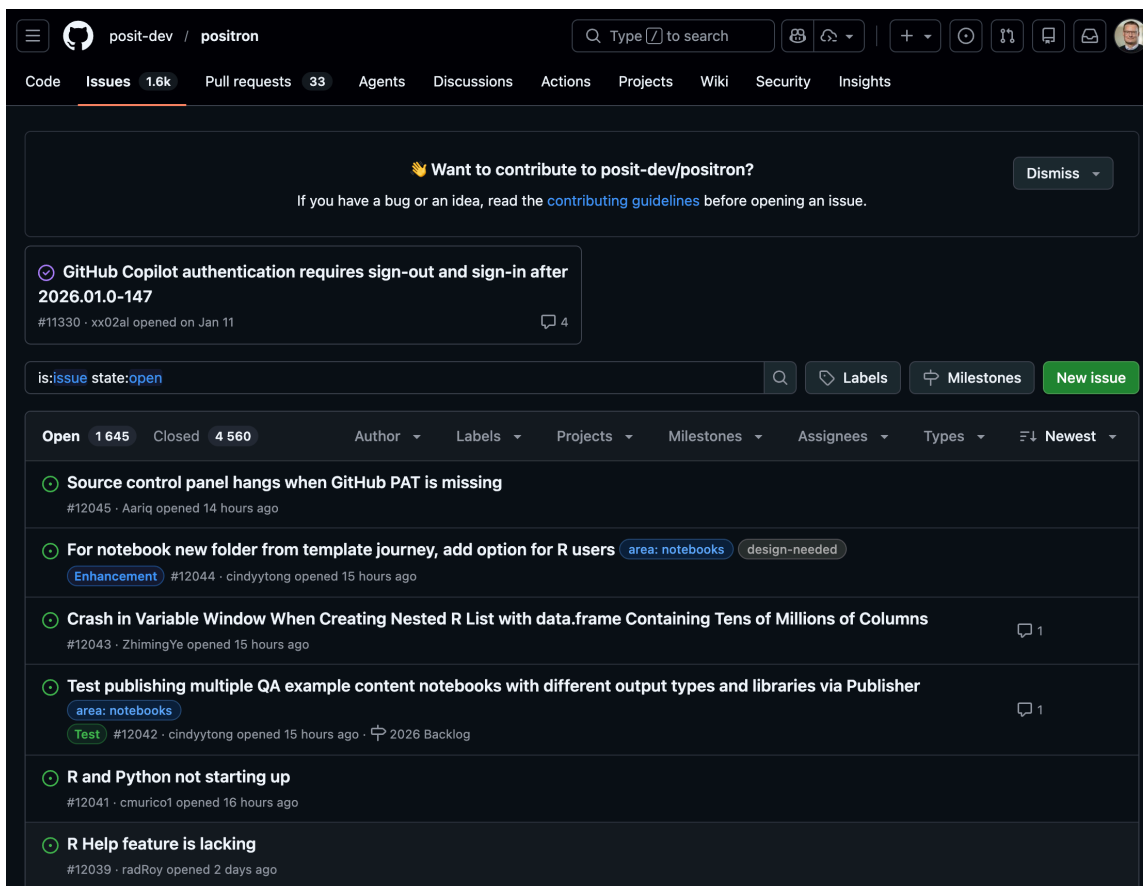
- [GitHub](#) (2008)
- [Bitbucket](#) (2008)
- [GitLab](#) (2011)
- [Gitea](#) - open source alternative to GitHub (get the same functionality locally or on a server)

They add:

- Pull requests / merge requests (collaboration tools)
- Issue tracking (discussions and bug reports etc)
- Code review
- CI/CD integration (advanced)

Issue tracking

- Issue tracking is not part of Git but is implemented in most (if not all) hosting platforms.
 - Used for bug reports and discussions between developers and users
 - Issues can be closed when fixed/adressed but are still found in the history
 - Each issue gets a number (in order) and those can be referenced in commit messages etc (ex: Fix #37, which will automatically close the issue)
-



The screenshot shows the GitHub Issues page for the repository `posit-dev/positron`. The page is in dark mode. At the top, there is a navigation bar with the repository name, a search bar, and various utility icons. Below the navigation bar, there are tabs for `Code`, `Issues` (1.6k), `Pull requests` (33), `Agents`, `Discussions`, `Actions`, `Projects`, `Wiki`, `Security`, and `Insights`. A prominent message asks if the user wants to contribute to `posit-dev/positron`, with a link to `contributing guidelines` and a `Dismiss` button. Below this, a featured issue is shown: `GitHub Copilot authentication requires sign-out and sign-in after 2026.01.0-147` (#11330), opened on Jan 11. The main content area displays a list of issues with filters for `is:issue state:open`. The list includes: `Source control panel hangs when GitHub PAT is missing` (#12045), `For notebook new folder from template journey, add option for R users` (#12044) with labels `area: notebooks` and `design-needed`, `Crash in Variable Window When Creating Nested R List with data.frame Containing Tens of Millions of Columns` (#12043), `Test publishing multiple QA example content notebooks with different output types and libraries via Publisher` (#12042) with label `area: notebooks` and `Test`, `R and Python not starting up` (#12041), and `R Help feature is lacking` (#12039). The page also shows filters for `Open` (1645) and `Closed` (4560) issues, and sorting options like `Newest`.

STA220 / crap

Q Type to search

Code Issues Pull requests Agents Actions Projects Wiki Security Insights Settings

Create new issue

Add a title *

Title

Add a description

Write Preview H B I

Type your description here...

Assignees: No one - [Assign yourself](#)

Labels: No labels

Type: No type

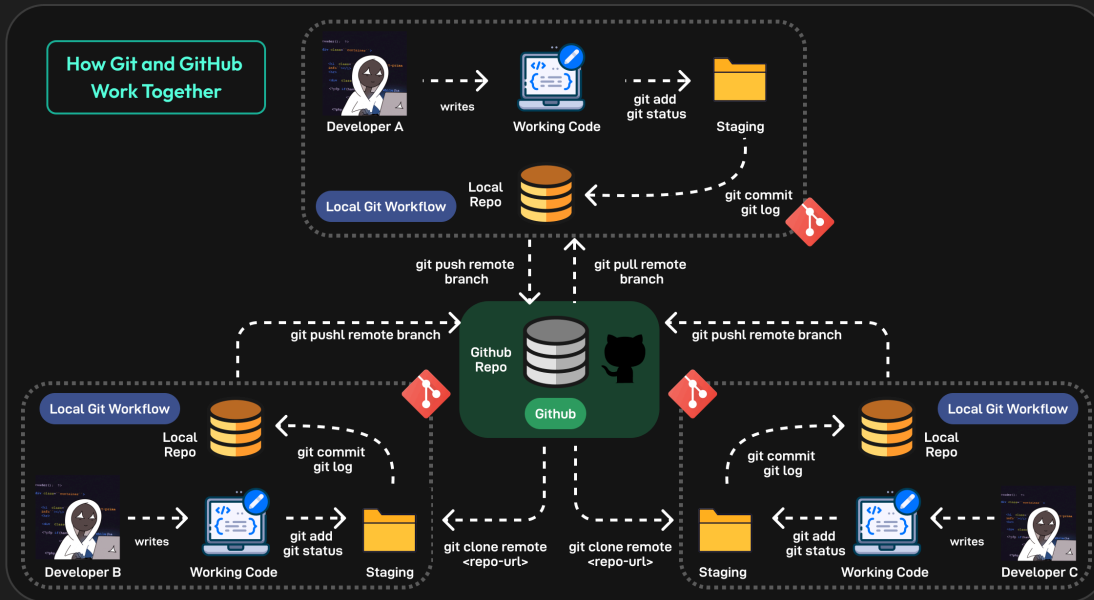
Projects: No projects

Milestone: No milestone



Paste, drop, or click to add files Write with Copilot

Create more

Git versus GitHub



Git vs GitHub Comparison

	 Git	 GitHub
Type	Git is a free, open-source version control tool	GitHub is a cloud-based, pay-for-use service that runs Git in the cloud
Installation	Git is installed locally on a developer's machine	GitHub is hosted in the cloud
Ownership	Git is maintained by the Linux Foundation	GitHub is owned by Microsoft
Use	Tool to manage different versions of edits, made to files in a git repository	It is a space to upload a copy of the Git repository
Features	Version control and source code management	Hosting code, collaboration, and project management
Tools	Minimal external tool configuration	Active marketplace for tool integration

Version Control Today

Modern usage includes:

- Code
- Documentation
- Data analysis (scripts, notebooks)

- Configuration and infrastructure

Git is integrated into:

- IDEs (RStudio, VS Code, Positron)
- CI/CD pipelines
- Cloud platforms

Version Control Beyond Code

Today, version control supports:

- Reproducible research
- Collaborative writing
- Data science workflows
- Teaching and learning

♥ Version control is now a core professional skill.

WARNING!

- Not everything should be shared!
- Scripts and documentation yes!
- But **Health data is sensitive!**
 - Do not share it!
 - Avoid unintentional sharing!
 - Private repositories are still shared with hosting provider!
- Avoid explicit file paths and sensitive info in scripts!
 - Can give information about data location and internal structure!
- No hard-coded passwords or API keys!

Git basics

(After installing the Git software)

- Collect all files related to a project in a folder
- Initialize a git repository in that folder
- Make changes to files
- Stage changes for commit
- Commit changes with a message
- Possibly push commits to remote repository

```
cd path/to/your/project
git init
git status
# make changes to files
git add filename1 filename2
git commit -m "Descriptive message about changes"
git remote add origin
```

Video tutorials

The video below is a good start to understand the basic concepts of Git and GitHub (and there are others to be found on YouTube).

Inspirational video

Watch this video even though some parts might be overwhelming. It gives a good overview of the current state (2025), even though many things will be too advanced for this course (it is not specifically aimed for statisticians or R users).

! .gitignore

The .gitignore file is very important in settings with health data! Pay close attention to [this section](#) of the video!

Git in Positron

- Remember that Positron is build on Code OSS (which shares a lot of features with VS Code).
- Branching and merging is possible but we will not cover that here.

Short official introduction from Microsoft:

More detailed introductions. Watch both! The first one is based on a Windows version of VS code and the second on Mac but the concepts are the same:

i Overwhelmed?

This video includes some parts which might be overwhelming if you are new to Git and GitHub. Don't worry! You don't need to understand everything right away. Just try to follow along with the basic concepts and steps. You will get more comfortable with practice.

Statistics projects

Not a single R script

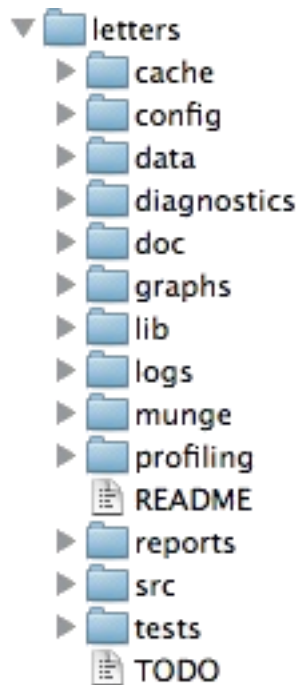
- Real projects are more complex than a single R script!
- Multiple scripts
- Multiple data files
- Documentation
- Reports
- Version control
- Reproducible workflows

Project structure

Common file structures

- Help you organize your thoughts

- Help others to collaborate
- simplifies paths used in your code



Example structure

```

/.../my_project/
├─ README.md      - project documentation
├─ TODO          - what should be done next?
├─ .git          - handled by git (hidden folder)
├─ .gitignore    - used by git but your responsibility!
├─ data/        - your data files (not under version control!)
│   └─ cancer.csv
│   └─ patients.qs
├─ R/          - your saved R functions
│   └─ function1.R
│   └─ function2.R
├─ reports/
└─ _targets.R   - targets pipeline script

```

README.md

- Document the purpose of the project
- What is it about?
- What is the aim?
- Who to contact for questions?
- In what circumstances was it created?

[Markdown format](#) (simple text with some possible formatting)

data folder

- Store your data files here as they are when you get them
- Avoid any modifications to the raw files!
- It is very easy to forget what you do if it can not be traced by code
- Do NOT include this folder in version control!
- Git is not good at handling large files
- Sensitive data should not be shared!
- Add data/* to your .gitignore file
- In realistic projects, data might come in varying formats
 - csv, txt, xlsx, sas7bdat, sav, dta, etc
 - some files might be very big (gigabytes not uncommon)

R folder

- Store your R functions here
- Document their purpose inline!
- Helps you to reuse code
- Easier to read main scripts
- Easier to test and debug code
- Easier to share code between projects

reports folder

- Store your reports here
- Quarto documents
- Document your analysis
- Include figures and tables
- Share with external collaborators

Computer practice

In [ECS1](#) we will use Positron and git/GitHub in action!

Also see the “Reading and practicing” section above for a more in-depth introduction (homework).

What you should know

- Reflect on the use of different software and how the rapid development in this field interplay with other important aspects of our field
- Be able to describe the principles of basic git commands (init, add, stage, commit, push, pull) and what they are used for (may be theoretical questions in the written exam)
- Use Git and GitHub in practice (but you can choose to do it either by commands or the GUI), this will be assessed in computer exercises and a later project.
- Similarly, you need to organize your projects according to best practice (but we will be the focus of [EL5](#)).

Bibliography

- [1] “The Unix Shell: Summary and Setup.” [Online]. Available: <https://swcarpentry.github.io/shell-novice/>
- [2] B. Rodrigues, “Building reproducible analytical pipelines with R.” [Online]. Available: <https://raps-with-r.dev/>
- [3] T. L. Staples, “Expansion and evolution of the R programming language,” *Royal Society Open Science*, vol. 10, no. 4, p. 221550, Apr. 2023, doi: [10.1098/rsos.221550](https://doi.org/10.1098/rsos.221550).