

Lecture 2: R workflow, Git, and R markdown

Data Science for Business Analytics

Thibault Vatter <thibault.vatter@unil.ch>

Department of Statistics, Columbia University and HEC Lausanne, UNIL

26.02.2018

1 R workflow

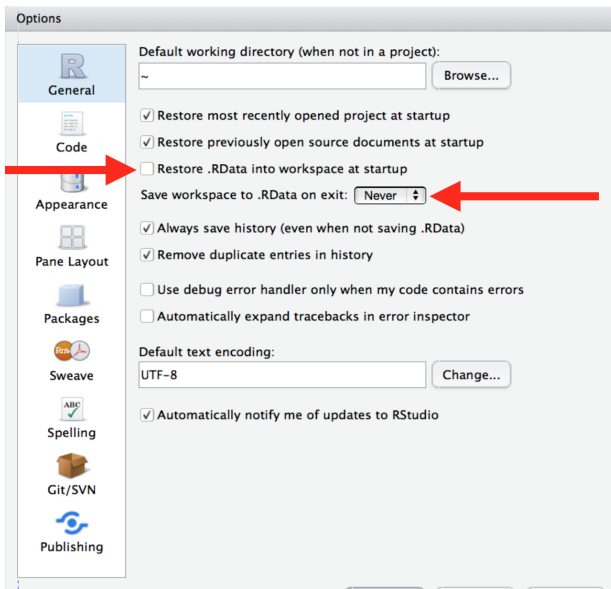
2 Git

3 R markdown

Two questions

- What about your analysis is “real”?
- Where does your analysis “live”?

What about your analysis is “real”?



- The console
- R scripts
- **RStudio projects:** *make it straightforward to divide your work into multiple contexts, each with their own working directory, workspace, history, and source documents.*

DEMO!

- Create an RStudio project for each data analysis project.
- Keep data files there.
- Keep scripts there.
- Save your outputs (plots and cleaned data) there.
- Only ever use relative paths, not absolute paths.

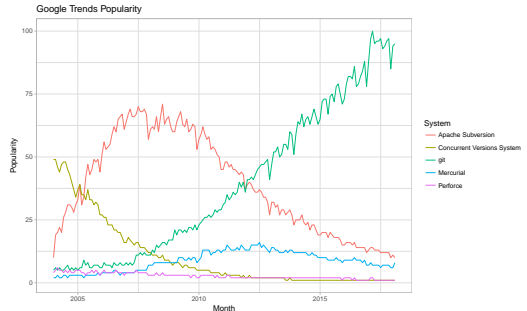
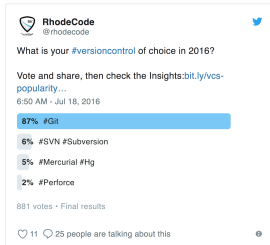
Everything you need is in one place, and cleanly separated from all the other projects that you are working on.

1 R workflow

2 Git

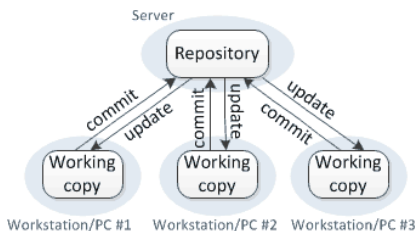
3 R markdown

- **wikipedia:** *management of changes to documents, computer programs, large web sites, and other collections of information*

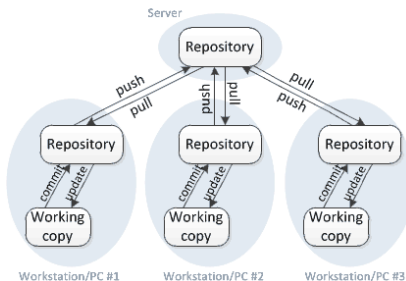


- Created by Linus Torvalds in 2005, his criteria:
 - ▶ Patching should take < 3 seconds
 - ▶ CVS as an ex. of what not to do (in doubt, do the opposite)
 - ▶ Distributed workflow
 - ▶ Strong safeguards against corruption (accidental or malicious)
- Maintained by Junio Hamano since 2005
- Part of the GNU free software project
- Source code written primarily in C, Shell, Perl, Tcl, Python
- Available for Windows, macOS, and Linux

Centralized version control



Distributed version control



Version control concepts and best practices by Michael Ernst

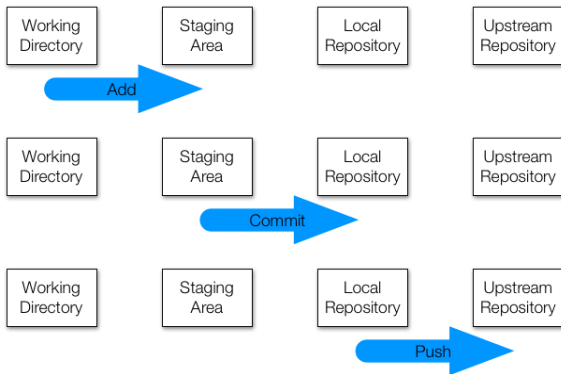


- Web-based version control service using git
- Bug tracking, feature requests, task management, and wikis for every project
- 20 million users and 57 million repositories (April 2017)
- Private and public repos
- [GitHub Student Developer Pack](#)
- Create account & email user to Natasha by **March 19**.

DEMO!

Workflow for your assignments

- Work on your assignment
- Commit changes to your local repository
- Push the changes to the github repo



source: github.com/datasciencelabs

- Use a descriptive commit message
- Make each commit a logical unit
- Avoid indiscriminate commits
- Incorporate others' changes frequently
- Share your changes frequently
- Coordinate with your co-workers
- Remember that the tools are line-based
- Don't commit generated files
- Understand your merge tool

See “Version control concepts and best practices” by Michael Ernst

- Last commit before midnight of due date as final submission
 - ▶ If there are commits after midnight, we will take the last commit up to the due date at 11:59 pm as the final version
- Check that the final commit is showing in your Github repo
 - ▶ “I forgot to push” is not an acceptable excuse
- Detailed tutorials (with lots of pictures):
 - ▶ [Setting-up Github](#)
 - ▶ [Git and RStudio](#)
 - ▶ [Github and RStudio \(alternative\)](#)
 - ▶ [Github and RStudio \(alternative 2\)](#)

1 R workflow

2 Git

3 R markdown

The two components:

- Literate programming
- Markdown

- Motivation: helps peers understand and replicate your results, find errors and suggest enhancements
- Introduced by Donald Knuth
- **wikipedia:** *a program is given as an explanation of the program logic in a natural language, such as English, interspersed with snippets of macros and traditional source code, from which a compilable source code can be generated [... It] represents a move away from writing programs in the manner and order imposed by the computer, and instead enables programmers to develop programs in the order demanded by the logic and flow of their thoughts*

What does this R code do?

```
1 data(women)
2 plot(women)
3 fit <- lm(weight ~ height, data = women)
4 abline(fit)
```

And this one?

```
1 # Analysis of the 'women' dataset in R
2 data(women) # Load the data
3 attach(women) # Attach data to path
4 plot(weight ~ height) # Make a scatter plot
5 fit <- lm(weight ~ height) # Fit linear model
6 abline(fit) # Add a line of best fit to the plot
```

Real programmers don't comment their code. If it was hard to write, it should be hard to understand. – unknown

OR

If you can't write clearly, you probably don't think nearly as well as you think you do. – Kurt Vonnegut

The *World Almanac and Book of Facts* (1975) includes a dataset of heights (in) and weights (lbs) of 15 American women aged 30–39. It is built into R:

```
1 data(women)
```

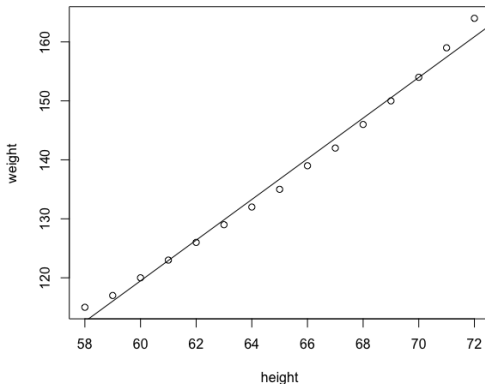
As height increases, weight appears to increase (almost) linearly: every inch in height adds approximately 3.45 lbs. This was determined by fitting a simple linear regression model of weight against height:

```
1 fit <- lm(weight ~ height, data = women)
```

Can't we do better? cont'd

The resulting least-squares regression line can be drawn on a scatter plot of height against weight, where the model seems appropriate:

```
1 plot(weight ~ height, data = women)  
2 abline(fit)
```



A lightweight markup language

■ Markup:

- ▶ A system for annotating a document in a way that is syntactically distinguishable from the text
- ▶ E.g., LaTeX and HyperText Markup Language (HTML)

■ Lightweight:

- ▶ A markup language with simple, unobtrusive syntax
- ▶ E.g., Markdown and R markdown

Here is some text:

- in *italics*,
- in **boldface**.

In LaTeX:

```
1 Here is some text:  
2 \begin{itemize}  
3 \item in \textit{italics},  
4 \item in \textbf{boldface}.  
5 \end{itemize}
```

In Markdown:

```
1 Here is some text:  
2 * in italics,  
3 * in boldface.
```

A markdown-based literate programming system

DEMO!

- Essential: [R Markdown cheat sheet](#)
- RStudio's [R markdown website](#)
 - ▶ [Tutorial](#) (to get you started)
 - ▶ [Output formats](#) (e.g., HTML, Word documents, PDFs, presentations, etc.)
- stuff written by [Yihui](#)
 - ▶ [knitr](#) and especially [its the options page](#)
 - ▶ [bookdown](#) to write technical reports
 - ▶ [blogdown](#) to even build your own website