

INTRODUCTION TO REPRODUCIBLE RESEARCH WITH SWEAVE AND ORG-MODE

YURII AULCHENKO

CONSULTANT AND INDEPENDENT RESEARCHER

YURII [DOT] AULCHENKO [AT] GMAIL [DOT] COM

REPRODUCIBLE RESEARCH

- Embeds executable code in research reports and publications, with the aim of allowing readers to re-run the analyses described

MY VIEW

- Precise way of documenting analysis
- Allows me to return to my previous analysis and
 - Reproduce / verify it
 - Build more analyses on the top
 - Find my previous analyses / results
- Allows running same analysis on different data set (“pipelining”)

LITERATE PROGRAMMING

- Enhances traditional software development by embedding code in explanatory essays and encourages treating the act of development as one of communication with future maintainers

MY VIEW

- Write the code / pipeline in a way that I can understand what is going on years after (or that someone can understand the code without me explaining it)
- This allows easy maintenance and code re-use

PROBLEM

Prepare data

Import Data

Perform QC

Analyze association

Generate report

PRIMITIVE WAY TO DO RR

```
drwxr-xr-x 17 yuryaulchenko staff 578 May 8 23:01 .
drwxr-xr-x 12 yuryaulchenko staff 408 May 9 15:06 ..
-rw-r--r-- 1 yuryaulchenko staff 390 May 8 22:49 .pversion
-rwxr-xr-x 1 yuryaulchenko staff 305 May 8 22:42 00_prepare_data.sh
-rw-r--r-- 1 yuryaulchenko staff 178 May 8 21:47 01_import2R.R
-rw-r--r-- 1 yuryaulchenko staff 206 May 8 22:46 02_doQC.R
-rw-r--r-- 1 yuryaulchenko staff 110 May 8 22:47 03_analysis.R
-rw-r--r-- 1 yuryaulchenko staff 7581 May 8 22:50 Rplots.pdf
-rw-r--r-- 1 yuryaulchenko staff 14639 May 8 22:50 dta.RData
-rw-r--r-- 1 yuryaulchenko staff 13213 May 8 22:50 dtaQCed.RData
-rw-r--r-- 1 yuryaulchenko staff 2003 May 8 22:49 phenotypes.txt
-rw-r--r-- 1 yuryaulchenko staff 2470 May 8 22:49 phenotypes_goodEOL.csv
-rw-r--r-- 1 yuryaulchenko staff 1844 May 8 22:49 plink.log
-rw-r--r-- 1 yuryaulchenko staff 59220 May 8 22:50 plink.raw
-rw-r--r-- 1 yuryaulchenko staff 1824 May 8 22:49 plink.tfam
-rw-r--r-- 1 yuryaulchenko staff 251312 May 8 22:49 plink.tped
-rw-r--r-- 1 yuryaulchenko staff 105 May 8 22:41 reformat.pl
```

R SWEAVE

- A module allowing documenting R analyses
- Our example:
 - Import files to R
 - Perform QC
 - Run association analysis

SUMMARY R SWEAVE

- **Great tool** for pedagogic; for RR in R (and may be some shell commands...)
- Lack of instant feedback (-RR): with larger analysis it may be very annoying that you need to re-run large chunks
- Not so good to combine multiple tools and scripts and pieces of own code

ORG-MODE

- OrgMode paper
 - <http://www.jstatsoft.org/v46/i03>
- OrgMode web-site
 - <http://orgmode.org>

LANGUAGES SUPPORTED

Language	Identifier	Language	Identifier
Asymptote	asymptote	Awk	awk
Emacs Calc	calc	C	C
C++	C++	Clojure	clojure
CSS	css	ditaa	ditaa
Graphviz	dot	Emacs Lisp	emacs-lisp
gnuplot	gnuplot	Haskell	haskell
Java	java		
Javascript	js	LaTeX	latex
Ledger	ledger	Lisp	lisp
Lilypond	lilypond	MATLAB	matlab
Mscgen	mscgen	Objective Caml	ocaml
Octave	octave	Org mode	org
Oz	oz	Perl	perl
Plantuml	plantuml	Python	python
R	R	Ruby	ruby
Sass	sass	Scheme	scheme
GNU Screen	screen	shell	sh
SQL	sql	SQLite	sqlite

LANGUAGES SUPPORTED

Language	Identifier	Language	Identifier
Asymptote	asymptote	Awk	awk
Emacs Calc	calc	C	C
C++	C++	Clojure	clojure
CSS	css	ditaa	ditaa
Graphviz	dot	Emacs Lisp	emacs-lisp
gnuplot	gnuplot	Haskell	haskell
Java	java	LaTeX	latex
Javascript	js	Lisp	lisp
Ledger	ledger	MATLAB	matlab
Lilypond	lilypond	Objective Caml	ocaml
Mscgen	mscgen	Org mode	org
Octave	octave	Perl	perl
Oz	oz	Python	python
Plantuml	plantuml	Ruby	ruby
R	R	Scheme	scheme
Sass	sass	shell	sh
GNU Screen	screen	SQLite	sqlite
SQL	sql		

ORG-MODE

- Platform for RR & Literate Programming
- Supports >30 languages
- Instant feedback
- Flexible controls
- Reports in many formats - PDF, HTML, ...