

6: GNU Make

`bit.ly/2018rr`

Here, we'll talk about automating the entire workflow for a project using the tool GNU Make.

GNU Make

- ▶ Automation the full project
- ▶ Document dependencies
- ▶ Only re-run things that need to be re-run

2

GNU Make is an old tool that was originally for automating the compilation of large, complex programs. But it's useful much more generally, though it does have some quirks.

In addition to automating, you'll be documenting the dependencies among the steps, and since the dependencies are defined, only stuff that needs to be re-run will be.

Automate the process (GNU Make)

```
R/analysis.html: R/analysis.Rmd Data/cleandata.csv
    cd R;R -e "rmarkdown::render('analysis.Rmd')"

Data/cleandata.csv: R/prepData.R RawData/rawdata.csv
    cd R;R CMD BATCH prepData.R

RawData/rawdata.csv: Python/xls2csv.py RawData/rawdata.xls
    Python/xls2csv.py RawData/rawdata.xls > RawData/rawdata.csv
```

3

GNU Make is an old (and rather quirky) tool for automating the process of building computer programs. But it's useful much more broadly, and I find it valuable for automating the full process of data file manipulation, data cleaning, and analysis.

In addition to **automating** a complex process, it also **documents** the process, including the dependencies among data files and scripts.

Automation with GNU Make

- ▶ Make is for more than just compiling software
- ▶ The **essence** of what we're trying to do
- ▶ Automates a workflow
- ▶ Documents the workflow
- ▶ Documents the dependencies among data files, code
- ▶ Re-runs only the necessary code, based on what has changed

4

People usually think of Make as a tool for automating the compilation of software, but it can be used much more generally.

To me, Make is the essential tool for reproducible research: automation plus the documentation of dependencies and workflows.

Installing Make

- ▶ On Macs, Make should be installed. Type “make --version” to check.
- ▶ On Windows, probably the easiest is to install **Rtools**, which includes Make.

cran.r-project.org/bin/windows/Rtools

Installation of these sorts of command-line tools on Windows can be a bit difficult.

How do you use Make?

- ▶ If you name your make file `Makefile`, then just go into the directory containing that file and type `make`
- ▶ If you name your make file `something.else`, then type `make -f something.else`
- ▶ Actually, the commands above will build the **first** target listed in the make file. So I'll often include something like the following.

```
all: target1 target2 target3
```

Then typing `make all` (or just `make`, if `all` is listed first in the file) will build all of those things.

- ▶ To be build a specific target, type `make target`. For example, `make Figs/fig1.pdf`

Details on the use.

Make with R Markdown

To use Make with R Markdown, you'll use a command like:

```
R -e "rmarkdown::render('my_report.Rmd')"
```

You'll need to tell your operating system where it can find **pandoc**. **RStudio** includes pandoc, but you need to add the relevant directory to your PATH.

Mac:

```
/Applications/RStudio.app/Contents/MacOS/pandoc
```

Windows:

```
"c:\Program Files\RStudio\bin\pandoc"
```

7

To use GNU Make with R Markdown, you'll need to be able to turn the thing into HTML (or PDF or Word) from the command line. The key trick is getting the operating system to know where pandoc can be located. You may also need to tell it about the location of R.

PATH is an environment/system variable that indicates which directories the operating system should look through to find executable programs.

So if you type **R** or **pandoc** at the command line, the operating system will look through all the folders in your **PATH** to find the program to run.

It can be tricky to get the **PATH** variable set properly.

~/.bash_profile

```
export PATH=$PATH:/Applications/RStudio.app/Contents/MacOS/pandoc

noclobber=1      # prevent overwriting of files
IGNOREEOF=1      # disable Ctrl-D as a way to exit
HISTCONTROL=ignoredups

alias cl='clear;cd'
alias rm='rm -i'
alias mv='mv -i'
alias cp='cp -i'
alias ls='ls -GF'
alias 'l.='='ls -d .[a-zA-Z]*'
alias ll='ls -lh'
alias md='mkdir'
alias rd='rmdir'
alias rmb='rm .*~ *~ *.bak *.bk!'

alias Rb='R CMD build --force --resave-data'
alias Ri='R CMD INSTALL --library=/Users/kbroman/Rlibs'
alias Rc='R CMD check --library=/Users/kbroman/Rlibs'
alias Rcc='R CMD check --as-cran --library=/Users/kbroman/Rlibs'
```

8

Use the `.bash_profile` file to define various variables and aliases to make your life easier.

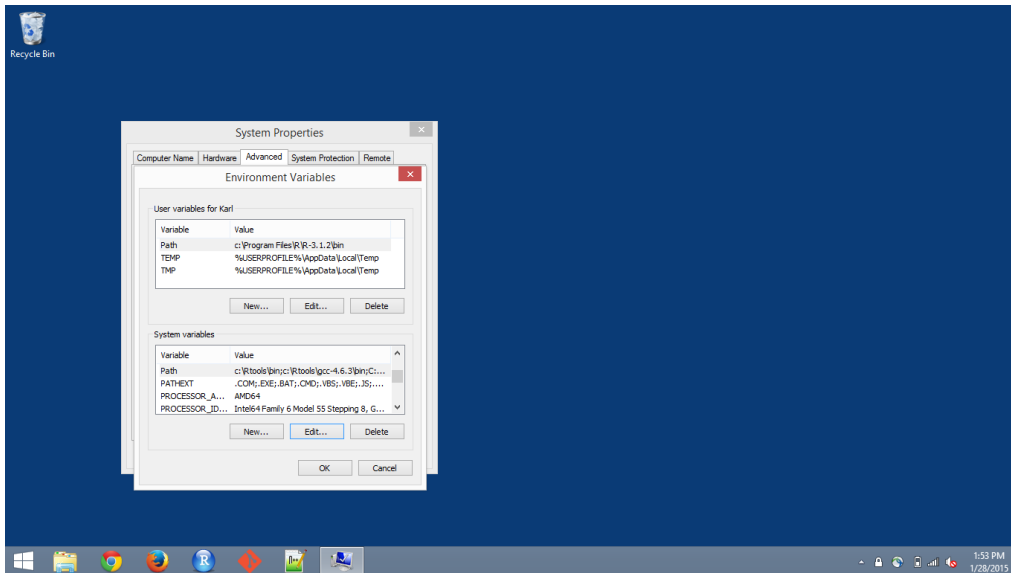
The most important variable is `PATH`: it defines the set of directories where the shell will look for executable programs. If “.” isn’t part of your `PATH`, you’ll need to type something like `./myscript.py` to execute a script in your working directory. So put “.” in your `PATH`.

My `.bash_profile` file sources a `.bashrc` file; I don’t quite understand when one is used versus the other. Google “`.bashrc` vs `.bash_profile`.” There are links to my `.bash_profile` and `.bashrc` files on the resources page at the course web site; some of it might just be total crap.

If you’re using Windows and Git Bash, the `.bash_profile` file will be in your Documents folder (I think).

Important note: use of aliases within your code will create reproducibility issues; another user will need those same aliases. Consider testing your code on a more basic account.

PATH in Windows



9

With Git Bash, you can have a `~/.bash_profile` file that adds stuff to your `PATH`, just as in Mac OS X and Linux.

But things will also be added to the `PATH` variable via the Path system variable and/or a Path user variable. You can get to these via the “Control panel,” but it’s a bit cumbersome.

The simplest way to get to the relevant dialog box seems to be to click Win-w (the little windows key and the w key) and searching for “path”.

Variables

- ▶ Define a **variable** like
`R_OPTS=--vanilla`
- ▶ Use it with a \$ and () or {}, for example:
`R CMD BATCH $(R_OPTS) fig1.R`

10

Variables are useful shorthand, for bits of code that you want to use repeatedly.

Automatic variables

There are a bunch of **automatic variables** that you can use to save yourself a lot of typing.

Here are the ones I use most:

<code>\$@</code>	the file name of the target
<code>\$<</code>	the name of the first dependency
<code>\$^</code>	the names of all dependencies
<code>\${@D}</code>	the directory part of the target
<code>\${@F}</code>	the file part of the target
<code>\${<D}</code>	the directory part of the first dependency
<code>\${<F}</code>	the file part of the first dependency

You'll see on the next slide how automatic variables get used.

Pattern rules

Pattern rules are like wildcards for file names: if a bunch of files are to be built the same way, you can use the symbol `%` as a wildcard.

For example, if you have two figures `fig1.pdf` and `fig2.pdf` that are to be built by `fig1.R` and `fig2.R`, respectively, you might do:

```
Figs/%.pdf: R/%.R  
    cd $(<D);R CMD BATCH $(<F)
```

The two figures' file names will need to be spelled out somewhere, for example as dependencies.

12

Pattern rules can greatly reduce the length of your **Makefile**, but they can also be rather frustrating. I recommend keeping things really simple initially and then move to pattern rules later, after you've been working with Make for a while.

Fancier example

```
FIG_DIR = Figs

mypaper.pdf: mypaper.tex $(FIG_DIR)/fig1.pdf $(FIG_DIR)/fig2.pdf
    pdflatex mypaper

# One line for both figures
$(FIG_DIR)/%.pdf: R/%.R
    cd R;R CMD BATCH $(<F)

# Use "make clean" to remove the PDFs
clean:
    rm *.pdf Figs/*.pdf
```

13

As I said, you can get really fancy with GNU Make.

Use variables for directory names or compiler flags. (This example is not a good one.)

Use pattern rules and automatic variables to avoid repeating yourself. With %, we have one line covering both `fig1.pdf` and `fig2.pdf`. The `$(<F)` is the file part of the first dependency. More on this later.

Look at the manual for Make and the many online tutorials, such as the one from Software Carpentry, or http://kbroman.org/minimal_make.

Resources

- ▶ `kbroman.org/minimal_make`
- ▶ `bost.ocks.org/mike/make`
- ▶ `robjhyndman.com/hyndsight/makefiles`
- ▶ Search github with `filename:Makefile`
 - `R CMD BATCH filename:Makefile`
 - `filename:Makefile user:yihui`

A lot of people are using Make; you can search for their files on GitHub. And there are some good introductory tutorials out there.

Activity

Go back to your R Markdown documents from this morning.

- ▶ Write a `Makefile` to produce different types of outputs from your various Rmd files.
- ▶ Add `make all` and `make clean` as targets
- ▶ What happens if you run `make all` twice?

15

An activity on Make.