

Workshop 1 - Reproducibility and R

Andrew Stewart

Andrew.Stewart@manchester.ac.uk



@ajstewart_lang



<https://github.com/ajstewartlang>

Workshop	Topic
1	Reproducibility and R
2	General Linear Model (Regression)
3	General Linear Model (ANOVA)
4	Mixed Models
5	Data Simulation and Advanced Data Visualisation
6	Reproducible Computational Environments and Presentations

Assignment

Assignment to be completed by Semester 1 exam period.

This Unit

- Everything we cover in this Unit will be taught following the principles of Open Science.
- All of the statistical analyses you do will be conducted using the R open source data science language.
- We will cover core topics in Statistics with an emphasis on reproducibility and transparency.
- You will learn how to produce reports in R Markdown - these reports contain your analysis code, your output and narrative describing what it all means.
- You will learn how to use GitHub and Binder for full computational reproducibility.
- You will learn how to use the `xaringan` package for reproducible presentations.

- The workshops will be a blend of seminars and hands-on labs.
- If you have a laptop, I recommend you use that rather than the cluster PCs.
- Another option is to use the RStudio Cloud server space I have set up: <https://bit.ly/3301Zd5>

Assessment

- The assessment for this Unit is via one coursework assignments, worth 100%.
- The assessment will cover data wrangling, data visualisation, and statistical modelling.
- You'll need to write your assessment using R Markdown (which we'll cover).

Replication and Reproducibility in Science

- Ioannidis (2005), *PLOS Medicine*, most published research findings are false.
- Prinz et al. (2011), *Nature Reviews Drug Discovery*, around 65% of cancer biology studies do not replicate.
- Button et al. (2013), *Nature Reviews Neuroscience*, small sample size undermines the reliability of neuroscience.
- MacLeod et al. (2014), *Lancet*, 85% of biomedical research resources are wasted.
- Baker (2015), *Nature*, 90% of scientists recognise a ‘reproducibility crisis’.
- Nosek & Errington (2017), *eLife*, out of first 5 replication attempts of preclinical cancer biology work, only 2 have replicated.



Power Posing: Brief Nonverbal Displays Affect Neuroendocrine Levels and Risk Tolerance

Psychological Science
XX(X) 1–6
© The Author(s) 2010
Reprints and permission:
sagepub.com/journalsPermissions.nav
DOI: 10.1177/0956797610383437
<http://pss.sagepub.com>


Dana R. Carney¹, Amy J.C. Cuddy², and Andy J. Yap¹

¹Columbia University and ²Harvard University

Abstract

Humans and other animals express power through open, expansive postures, and they express powerlessness through closed, contractive postures. But can these postures actually cause power? The results of this study confirmed our prediction that posing in high-power nonverbal displays (as opposed to low-power nonverbal displays) would cause neuroendocrine and behavioral changes for both male and female participants: High-power posers experienced elevations in testosterone, decreases in cortisol, and increased feelings of power and tolerance for risk; low-power posers exhibited the opposite pattern. In short, posing in displays of power caused advantaged and adaptive psychological, physiological, and behavioral changes, and these findings suggest that embodiment extends beyond mere thinking and feeling, to physiology and subsequent behavioral choices. That a person can, by assuming two simple 1-min poses, embody power and instantly become more powerful has real-world, actionable implications.

(In)famous studies...

- Power posing
- Ego depletion
- Social priming
- Marshmallow test performance predicts future achievement
- Stanford prison experiment
- Growth mindset
- Learning styles
- Any others you know of?

- Button et al. (2013), *Nature Reviews Neuroscience*, small sample size undermines the reliability of neuroscience. Nord et al., (2017), *Journal of Neuroscience*, highlight wide heterogeneity in power in neuroscience studies.

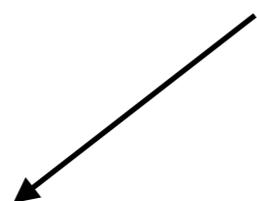


Table 2. Median, maximum, and minimum power subdivided by study type

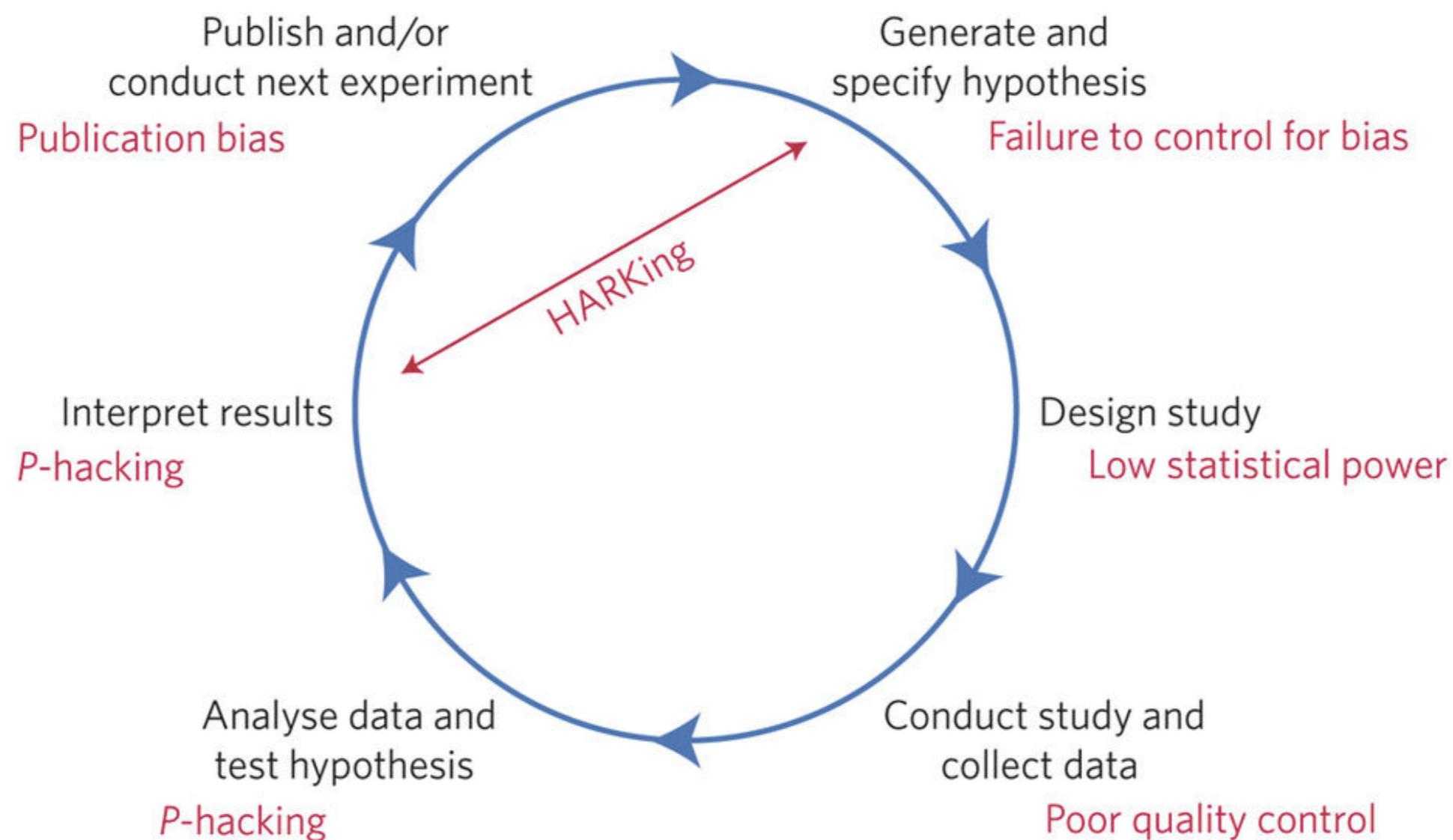
Group	Median power (%)	Minimum power (%)	Maximum power (%)	2.5 th and 97.5 th percentile (based on raw data)	95% HDI (based on GMMs)	Total N
All studies	23	0.05	1	0.05–1.00	0.00–0.72, 0.80–1.00	730
All studies excluding null	30	0.05	1	0.05–1.00	0.01–0.73, 0.79–1.00	638
Genetic	11	0.05	1	0.05–0.94	0.00–0.44, 0.63–0.93	234
Treatment	20	0.05	1	0.05–1.00	0.00–0.65, 0.91–1.00	145
Psychology	50	0.07	1	0.07–1.00	0.02–0.24, 0.28–1.00	198
Imaging	32	0.11	1	0.11–1.00	0.03–0.54, 0.71–1.00	65
Neurochemistry	47	0.07	1	0.07–1.00	0.02–0.79, 0.92–1.00	50
Miscellaneous	57	0.11	1	0.11–1.00	0.09–1.00	38

Is there not just “good science” and “bad science”?

Without realising it, good scientists have been
engaging in questionable research practices (QRPs)...

Problems include *p*-hacking, lack of power, HARKing, failing (refusal) to share data and code, too many researcher degrees of freedom...

From: [A manifesto for reproducible science](#)



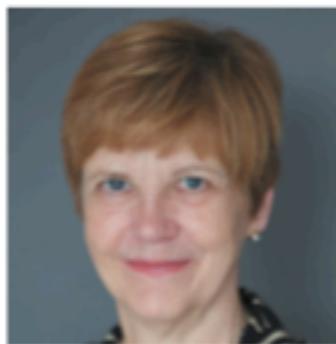
Munafo et al. (2017), *Nature Human Behaviour*

Replicable Science ≠ Reproducible Science

Replicable Science is when someone else can run a study the same as or conceptually equivalent to your one, and find a similar pattern of effects.

Reproducible Science is when someone else can take your data and your analysis code, run it and then find the same effects that you have reported.

ROBERT TAYLOR



Rein in the four horsemen of irreproducibility

Dorothy Bishop describes how threats to reproducibility, recognized but unaddressed for decades, might finally be brought under control.

More than four decades into my scientific career, I find myself an outlier among academics of similar age and seniority: I strongly identify with the movement to make the practice of science more robust. It's not that my contemporaries are unconcerned about doing science well; it's just that many of them don't seem to recognize that there are serious problems with current practices. By contrast, I think that, in two decades, we will look back on the past 60 years — particularly in biomedical science — and marvel at how much time and money has been wasted on flawed research.

How can that be? We know how to formulate and test hypotheses in controlled experiments. We can account for unwanted variation with statistical techniques. We appreciate the need to replicate observations.

Yet many researchers persist in working in a way almost guaranteed not to deliver meaningful results. They ride with what I refer to as the four horsemen of the reproducibility apocalypse: publication bias, low statistical power, *P*-value hacking and HARKing (hypothesizing after results are known). My generation and the one before us have done little to rein these in.

In 1975, psychologist Anthony Greenwald noted that science is prejudiced against null hypotheses; we even refer to sound work supporting such conclusions as 'failed experiments'. This prejudice leads to publication bias: researchers are less likely to write up studies that show no effect, and journal editors are less likely to accept them. Consequently, no one can learn from them, and researchers waste time and resources

be adequately powered. Other disciplines have yet to catch up.

I stumbled on the issue of *P*-hacking before the term existed. In the 1980s, I reviewed the literature on brain lateralization (how sides of the brain take on different functions) and developmental disorders, and I noticed that, although many studies described links between handedness and dyslexia, the definition of 'atypical handedness' changed from study to study — even within the same research group. I published a sarcastic note, including a simulation to show how easy it was to find an effect if you explored the data after collecting results (D. V. M. Bishop *J. Clin. Exp. Neuropsychol.* **12**, 812–816; 1990). I subsequently noticed similar phenomena in other fields: researchers try out many analyses but report only the ones that are 'statistically significant'.

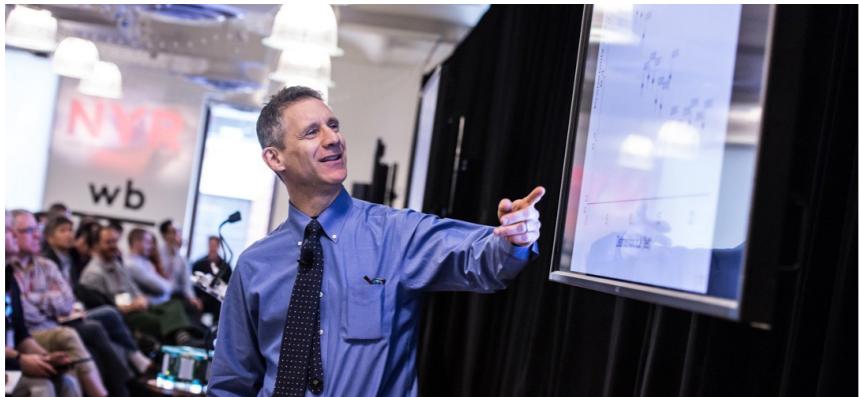
This practice, now known as *P*-hacking, was once endemic to most branches of science that rely on *P* values to test significance of results, yet few people realized how seriously it could distort findings. That started to change in 2011, with an elegant, comic paper in which the authors crafted analyses to prove that listening to the Beatles could make undergraduates younger (J. P. Simmons *et al. Psychol. Sci.* **22**, 1359–1366; 2011). "Undisclosed flexibility," they wrote, "allows presenting anything as significant."

The term HARKing was coined in 1998 (N. L. Kerr *Pers. Soc. Psychol. Rev.* **2**, 196–217; 1998). Like *P*-hacking, it is so widespread that researchers assume it is good practice. They look at the data, pluck out a finding that looks exciting and write a paper to tell a story around this result. Of course, researchers should be free to explore their

MANY RESEARCHERS
PERSIST IN WORKING
IN A WAY ALMOST
GUARANTEED
NOT
TO DELIVER
MEANINGFUL
RESULTS.

How do we make our science more replicable?

How do we make our science more reproducible?



<http://www.stat.columbia.edu/~gelman/>

Andrew Gelman gives the following recommendations to researchers:

- Analyze all your data.
- Present all your comparisons.
- Make your data public.
- Put in the effort to take accurate measurements (low bias, low variance, and a large enough sample size).
- Do repeated-measures comparisons where possible.

Open Science practices include...

- Pre-registering experiments.
- Registered reports.
- Using preprint servers (e.g., bioRxiv, PsyArXiv).
- Making data and analysis code freely available (e.g., via GitHub, OSF).
- Open access to journal articles.
- ...and more.

Open Science recently recognised by G7 Science Ministers...

Focus: Incentives and the researcher ecosystem

Ambition: Foster a research environment in which career advancement takes into account Open Science activities, through incentives and rewards for researchers, and valuing the skills and capabilities in the Open Science workforce.

Recommendations:

At national levels: G7 nations should each engage with research stakeholders to identify and implement enhancements to research evaluation and reward systems that take into consideration the Open Science activities carried out by researchers and research institutions. Topics that could be discussed include:

- Recognizing Open Science practices during evaluation of research funding proposals, and research outcomes;
- Recognizing and rewarding research productivity and impact that reflect open science activities by researchers during career advancement reviews;
- Including credit for service activities such as reviewing, evaluating, and curation and management of research data; and,
- Developing metrics of Open Science practices.

Panel criteria and working methods

200. The sub-panels welcome research practice that supports reproducible science and the application of best practice. Examples include registered reports, pre-registration, publication of data sets, experimental materials, analytic code, and use of reporting checklists for publication purposes and those relating to the use of animals in research. These contribute to the evaluation of rigour for submitted outputs. Replication studies may be submitted as outputs and will be evaluated on the extent to which they contribute significant new knowledge, improved methods, or advance theory or practice¹.

346...

Within the context of the institution's strategy, how the submitting unit is progressing towards an open research environment, including where this goes above and beyond the REF open access policy requirements, and wider activity to encourage the effective sharing and management of research data, as appropriate to the discipline. Consideration of reproducibility should also be included where relevant to the discipline.

is beginning to appear in tenure-track job adverts...

Our Department embraces the values of open and reproducible science, and candidates are encouraged to address (in their statements and/or cover letter) how they have pursued and/or plan to pursue these goals in their work.

and is forming part of Universities' teaching manifestos.

Teaching with Open Science commitment:

To teach the practices and skills of open research and science in our undergraduate and postgraduate degree programmes

- a. Promote open science in our teaching.
- b. Design a Research Methods curriculum that teaches skills for open science and uses open science to enhance teaching (for example: teach R and use open data to practice analysis skills).
- c. Learn about and adopt open educational practices in our teaching.
- d. Produce and promote tools for helping student researchers adopt open practices, including training and guidance suitable to their level of study.
- e. Author, share and use open educational resources to promote teaching with open science beyond our School and Institution.
- f. Support our colleagues to learn the skills of teaching Open Science.

Data Sharing - FAIR data



Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.

FINDABLE



Metadata and data are understandable to humans and machines. Data is deposited in a trusted repository.

ACCESSIBLE



Metadata use a formal, accessible, shared, and broadly applicable language for knowledge representation.

INTEROPERABLE



Data and collections have a clear usage licenses and provide accurate information on provenance.

REUSABLE

**How do you do Open and
Reproducible Science?**

Before Data Collection

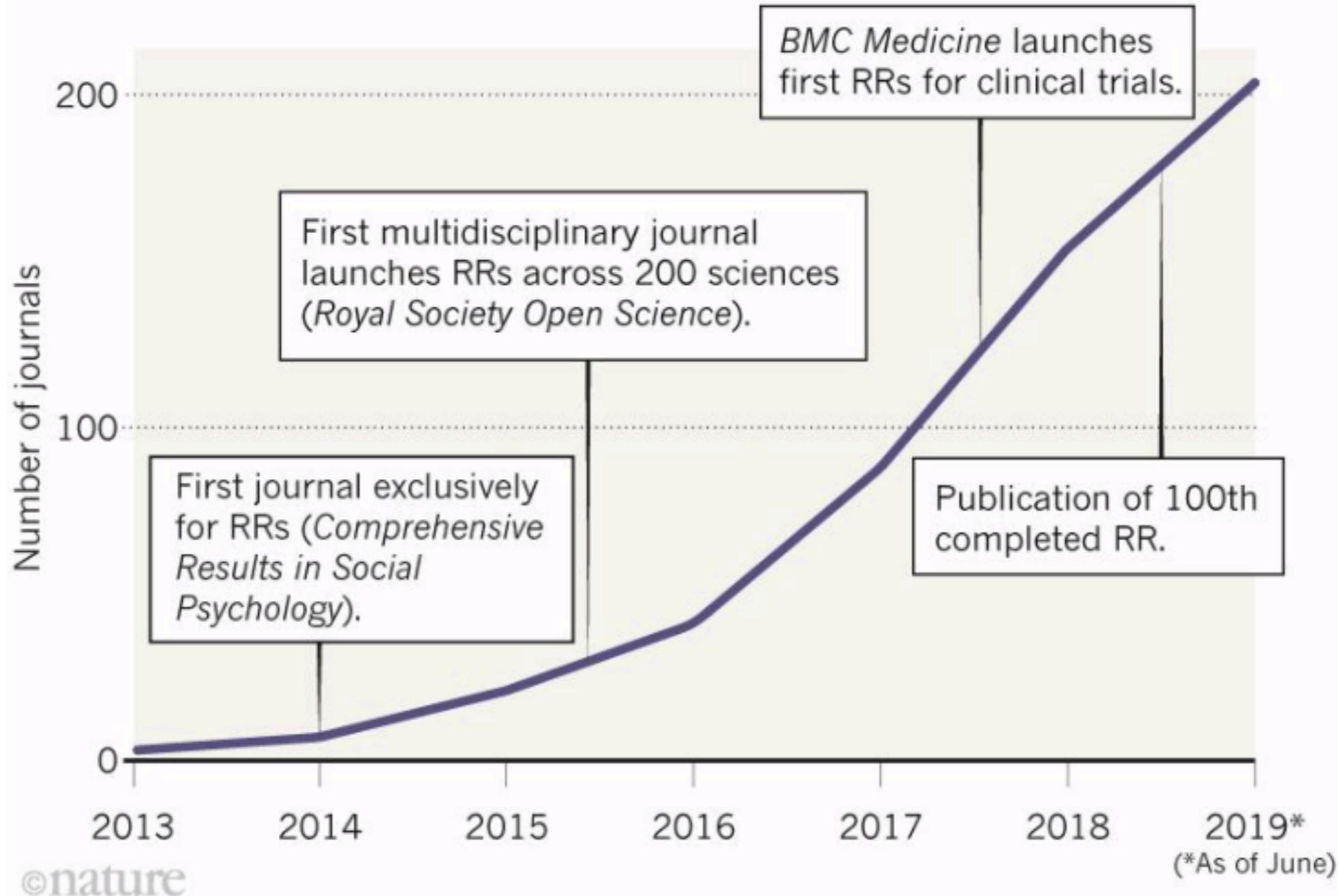
- Specify your hypotheses and analysis plan.
- **Pre-register** your hypotheses and analysis plan at osf.io
- Consider data simulation so that you can write your analysis script before you have your real data.
- Consider submitting as a **registered report** - currently more than **200** journals now support this route. This involves acceptance in principle before you have even started collecting your data.

Registered Reports



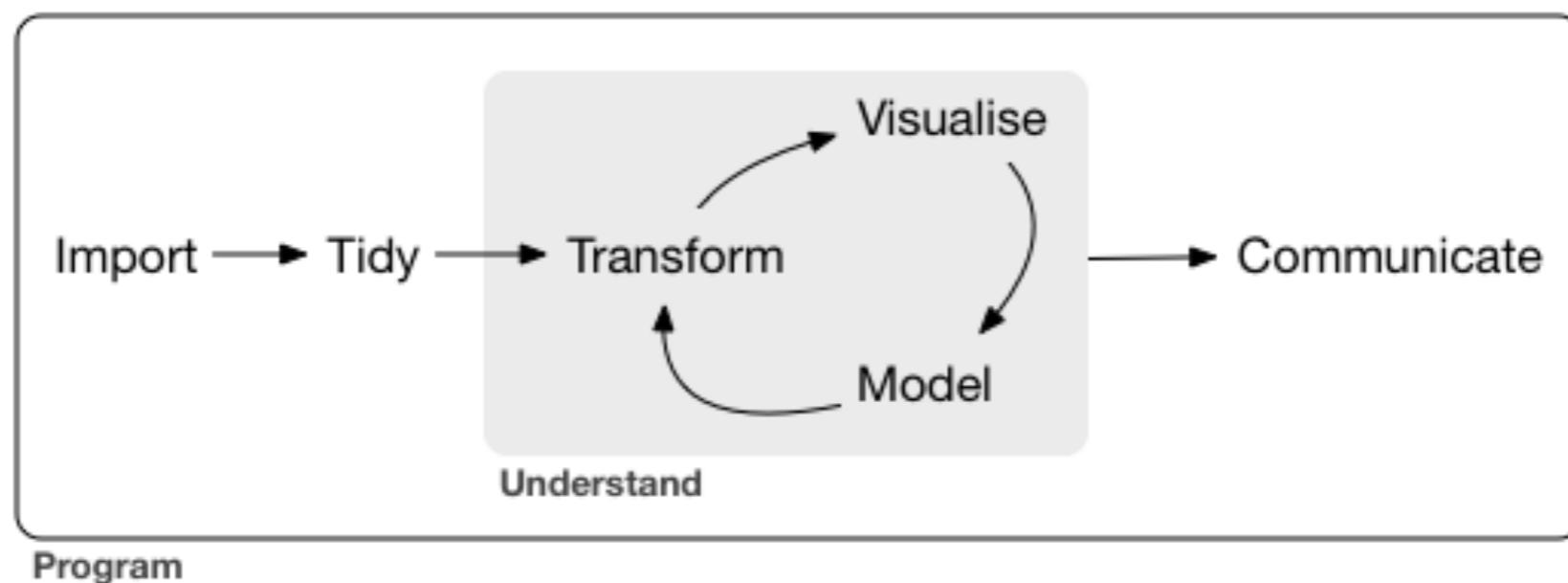
RAPID RISE

Since 2013, the number of journals offering Registered Reports (RRs) has risen to more than 200 titles.



After Data Collection

- You need to use analysis software that allows for open sharing and reproducibility of the entire data wrangling/analysis/write-up workflow.



Hadley Wickham and Garrett Grolemund

- You can share your data via [osf.io](#), on GitHub etc.:

P.s	Item	Condition	Probmanip	Speaker	statement	response	final	Meaning	Imposition
1	1	1	1708	302	1399	1867	1206	Indirect	High
1	2	2	1466	296	1377	1674	828	Indirect	Low
1	3	3	1393		1494	1950	1812	Direct	High
1	4	4	2463	530	1691	1866	965	Direct	Low
1	5	1	1552	267	1332	1477	1345	Indirect	High
1	6	2	1445	444	1004	1067	797	Indirect	Low
1	7	3	2159	501	739	1231	2240	Direct	High
1	8	4	1459		1086	946	978	Direct	Low
1	9	1	3302		1503	900	1736	Indirect	High

- alongside your analysis code

```
--  
26 FPs$Meaning <- as.factor(FPs$Meaning)  
27 FPs$Imposition <- as.factor(FPs$Imposition)  
28  
29 #this sets up the contrasts so that the intercept in the mixed LMM is the grand mean (i.e., the mean of all conditions)  
30 my.coding <- matrix (c(.5, -.5))  
31  
32 contrasts (FPs$Meaning) <- my.coding  
33 contrasts (FPs$Imposition) <- my.coding  
34  
35 #construct the models with crossed random effects for subjects and items for the pre-critical, critical and post-crtical region  
36 fpmodelprec <- lmer (Probmanip ~ Meaning*Imposition + (1+Meaning*Imposition |P.s) + (1+Meaning+Imposition |Item), data=FPs, REML=F)  
37 summary (fpmodelprec)  
38 lsmeans (fpmodelprec, pairwise~Meaning*Imposition, adjust="none")  
39  
40 fpmodelc <- lmer (statement ~ Meaning*Imposition + (1+Meaning*Imposition |P.s) + (1+Meaning*Imposition |Item), data=FPs, REML=T)  
41 summary (fpmodelc)  
42 lsmeans (fpmodelc, pairwise~Meaning*Imposition, adjust="none")  
43  
44 fpmodelpostc <- lmer (response ~ Meaning*Imposition + (1+Meaning*Imposition |P.s) + (1+Meaning+Imposition |Item), data=FPs, REML=F)  
45 summary (fpmodelpostc)  
46 lsmeans (fpmodelpostc, pairwise~Meaning*Imposition, adjust="none")  
47  
48 #Regression Path Analysis  
49 #Read in Regression Path data  
50 RPs <- read.csv("~/RPs.csv")  
51  
52 RPs$Meaning <- as.factor(RPs$Meaning)  
53 RPs$Imposition <- as.factor(RPs$Imposition)  
54  
55 contrasts (RPs$Meaning) <- my.coding  
56 contrasts (RPs$Imposition) <- my.coding  
57  
58 #construct the models with crossed random effects for subjects and items for the pre-critical, critical and post-crtical region  
59 rpmodelprec <- lmer (Probmanip ~ Meaning*Imposition + (1+Meaning*Imposition |P.s) + (1+Meaning*Imposition |Item), data=RPs, REML=F)
```

And preserve it with a DOI via Zenodo

The screenshot shows a web browser window with the URL zenodo.org/account/settings/github/. The browser's address bar also lists other sites like Google Scholar, Scopus, BBC News, etc. The Zenodo interface has a blue header with the logo, a search bar, and navigation links for Upload and Communities. A user profile is shown on the right. The main content area is titled "GitHub Repositories" and includes a "Get started" section with three steps: 1. Flip the switch (with an "ON" button), 2. Create a release, and 3. Get the badge (with a DOI example: 10.5281/zenodo.8475). Below this is a "Repositories" section showing a single repository: [ajstewartlang/Affective-Theory-of-Mind-Inferences](#), which has its switch set to "OFF". On the left, a sidebar under "Settings" shows options: Profile, Change password, Security, Linked accounts, Applications, Shared links, and GitHub (which is currently selected).

Why R?

- R allows you to engage in reproducible research.
- Statistical packages in R reflect the latest advances in the fields of Statistics and Data Science.
- Advanced models such as Linear Mixed Models (LMMs) are easy to build in R - many of the best journals now require LMMs to be used as they are more powerful and flexible than classical techniques such as ANOVA.
- Lots of amazing data visualisation packages available.
- In many institutions, the next generation of academics (M-level and PhD students, post-docs etc.) are learning R skills as part of their training.
- Plays an important role in Open Science.

What role can R play in Open Science?

- R scripts are easy to share allowing for reproducibility and easy public sharing of data and code.
- R is free, open source software that is much more flexible and powerful than SPSS.
- There is an active R community continuously updating statistical tests and packages that run in R.
- As R is a programming language, it forces you to know your data.

R vs. SPSS

“SPSS is like a bus - easy to use for the standard things, but very frustrating if you want to do something that is not already pre-programmed.

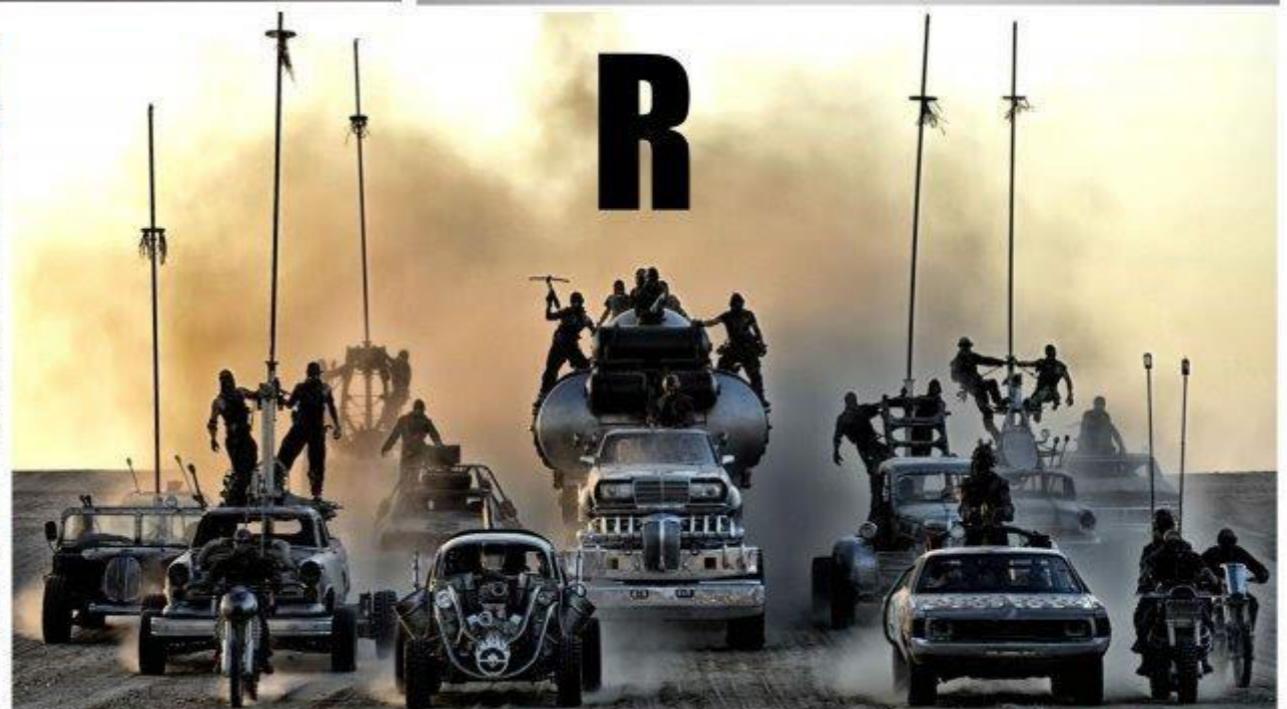
R is a 4-wheel drive off-roader, with a bike on the back, a kayak on top, good walking and running shoes in the passenger seat, and mountain climbing and spelunking gear in the back.

R can take you anywhere you want to go if you take time to learn how to use the equipment, but that is going to take longer than learning where the bus stops are in SPSS.” (Greg Snow, 2010, stackoverflow.com).

In meme form...

If statistics programs/languages were cars...

Image credit Darren Dahly @statsepi



O'REILLY®



R for Data Science

VISUALIZE, MODEL, TRANSFORM, TIDY, AND IMPORT DATA

Hadley Wickham &
Garrett Grolemund

Available electronically
for free at:

<http://r4ds.had.co.nz>

The R Series

Advanced R



Hadley Wickham

 CRC Press
Taylor & Francis Group
A CHAPMAN & HALL BOOK

Available electronically for
free at:

<https://adv-r.hadley.nz>

“Hadley Wickham, the Man Who Revolutionized R”



Chief Scientist at RStudio, author of key R packages incl. `ggplot2`, `tidyverse`, `dplyr` - all components of the tidyverse.

R User Groups...

R user groups exist all over the world with many in the UK (incl. Manchester)...

You can even buy the t-shirt...

UK

England

- Birmingham: [BRUM](#)
- Canterbury: [CanterburyR](#)
- Cambridge: [CambR](#)
- Coventry: [Warwick R User Group](#); [@WarwickRUG](#)
- Exeter: [Exeter R Users Group](#)
- London: [LondonR](#)
- Manchester: [ManchesterR](#)
- Newcastle: [R North East](#); [@RstatsNE](#)
- Nottingham: [Nottingham R User group](#)
- Oxford: [R user group Oxford](#); [@rusersoxford](#)
- Sheffield: [SheffieldR](#); [@Sheffield_R_](#)



RStudio
\$17



magrittr
\$17



Shiny
\$17



I like the way you R
\$17

...and conferences

WHO PLAN VENUE SPONSORS

R

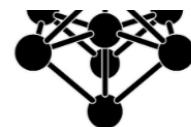
rOpenSci Unconf

May 25 - 26 2017 • Los Angeles • CA



useR!2017
BRUSSELS

04.07.2017 - 07.07.2017



NEWS ABOUT REGISTRATION PROGRAM VENUE FAQ EXTRA



FOLLOW US ON

CONFERENCE BROCHURE

NEWS

LAST CALL

Only 15 spots left for the useR!2017 Conference in Brussels next week! Be quick! Go to <https://user2017.brussels/> to register! ...

NEWS

Conference Brochure online!

The useR!2017 Conference Brochure is online now! Take a look and find an answer to all your practical questions. A printed conference brochure will be ...

NEWS

useR!2017 App live now!

We are very proud to announce that our useR!2017 app is live! Discover the full schedule, speakers and the venue on the go! Download the app here: iOS ...

NEW YORK R CONFERENCE SPEAKERS AGENDA SPONSORS VIDEOS ▾

NYR

NEW YORK R CONFERENCE

BROUGHT TO YOU BY [LANDER ANALYTICS](#) AND [WORK-BENCH](#)

useR!

useR! 2018

THE CONFERENCE FOR USERS OF R

JULY 10-13, 2018.

BRISBANE, AUSTRALIA.



HOME PROGRAMME ▾ REGISTRATION ▾ CODE OF CONDUCT NEWS TRAVEL ▾ FAQ CONTACT

R skills are in high demand...

February 11, 2014

R skills attract the highest salaries

Two recent salary surveys have shown that [R language](#) skills attract median salaries in excess of \$110,000 in the United States.

In the [2014 Dice Tech Salary Survey](#) of over 17,000 technology professionals, the highest-paid IT skill was R programming. While [big-data skills in general featured strongly](#) in the top tier, having R at the top of the list reflects the strong demand for skills to make sense of, and extract value from big data.

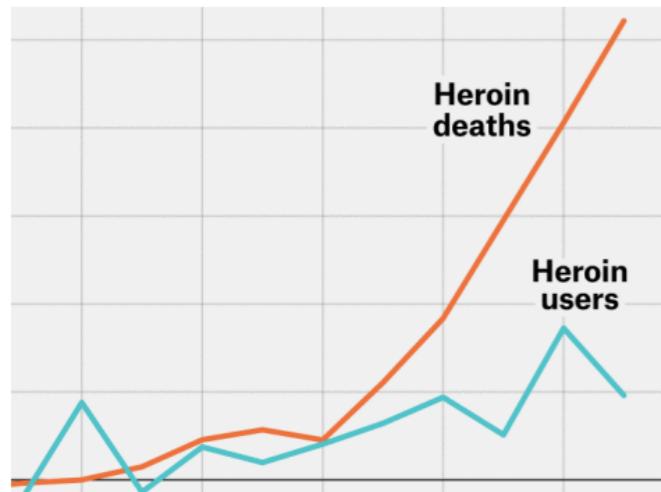
AVERAGE SALARY FOR High Paying Skills and Experience		
SKILL	2013	YR/YR CHANGE
R	\$ 115,531	n/a
NoSQL	\$ 114,796	1.6%
MapReduce	\$ 114,396	n/a
PMBok	\$ 112,382	1.3%
Cassandra	\$ 112,382	n/a
Omnigraffle	\$ 111,039	0.3%
Pig	\$ 109,561	n/a
SOA (Service Oriented Architecture)	\$ 108,997	-0.5%
Hadoop	\$ 108,669	-5.6%
Mongo DB	\$ 107,825	-0.4%

Similarly, the recent [O'Reilly Data Scientist Survey](#) also found R skills amongst those that pay in the \$110,000-\$125,000 range (albeit amongst a much smaller and specialized sample of respondents).

and R is widely used by a number of organisations (incl. the ONS)...

FiveThirtyEight

Politics Sports Science & Health Economics Culture



OPIOIDS
Data On Drug Use Is Disappearing Just When We Need It Most

By Kathryn Casteel

FEATURES

THE LATEST

JUN. 30 Why Republicans Might Be Forced To Oppose Tax Cuts

JUN. 29 Data On Drug Use Is Disappearing Just When We Need It Most

JUN. 28 The Health Care System Is Leaving The Southern Black Belt Behind

JUN. 26 The New CBO Report On Health Insurance Didn't Do Republicans Any Favors

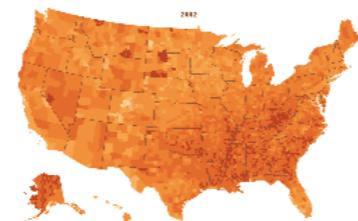
JUN. 26



Most recent: Politics podcast

INTERACTIVES

35 Years Of American Death



YouGov UK

TAKE PART

SEE RESULTS

SOLUTIONS

Login

+ Join

How the YouGov model for the 2017 General Election works



By Douglas Rivers is a professor of political science at Stanford University and Chief Scientist at YouGov PLC.

In General Election 2017, Politics

On May 31, 2017, 6 a.m.

Share

RELATED ARTICLES...

How did 2015 voters cast their ballot at the 2017 general election?

It wasn't just the kids that boosted Labour

John Humphrys - A Government Adrift?

Theresa May is now almost as unpopular as pre-campaign Corbyn



YouGov ElectionCentre

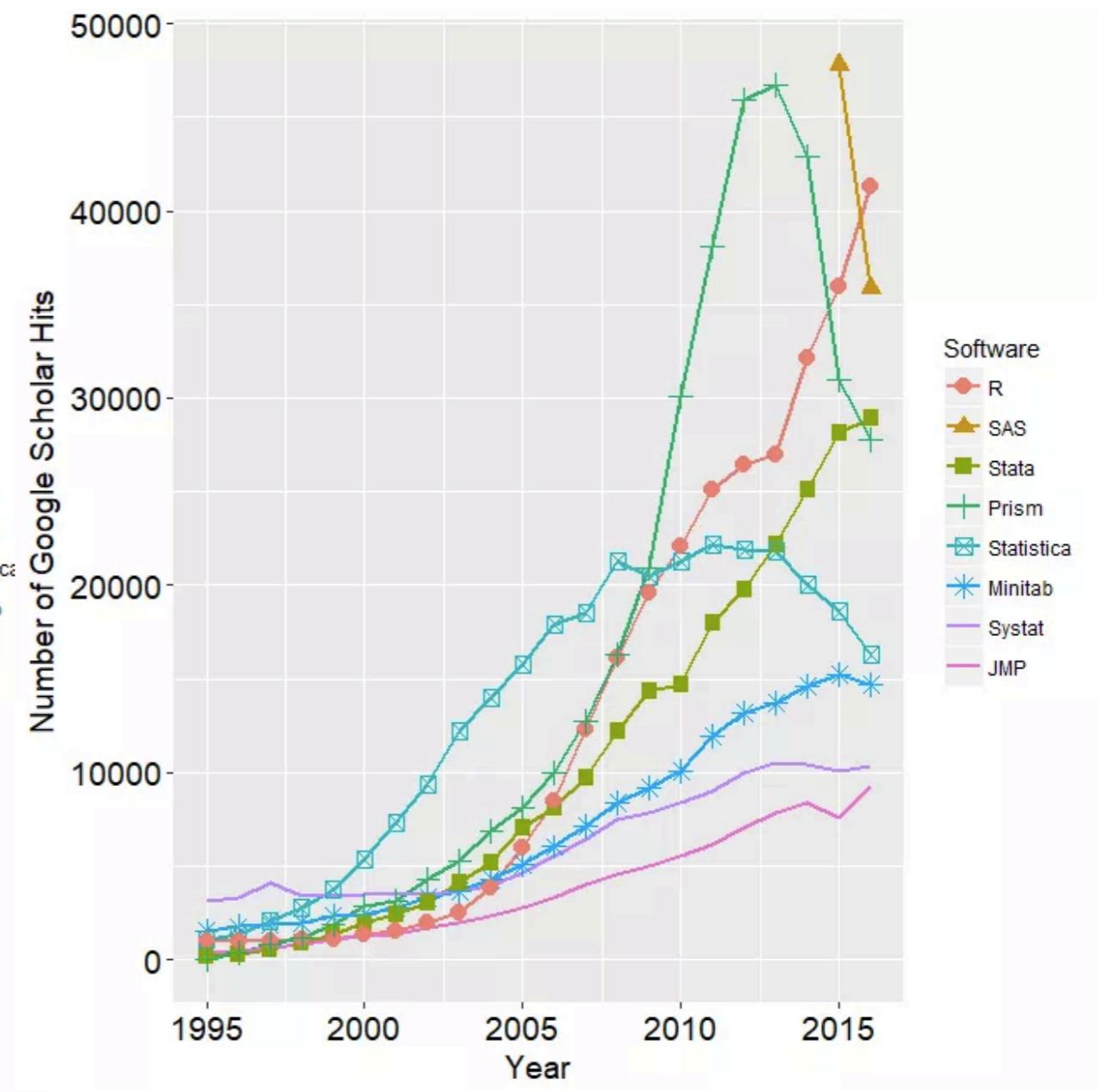
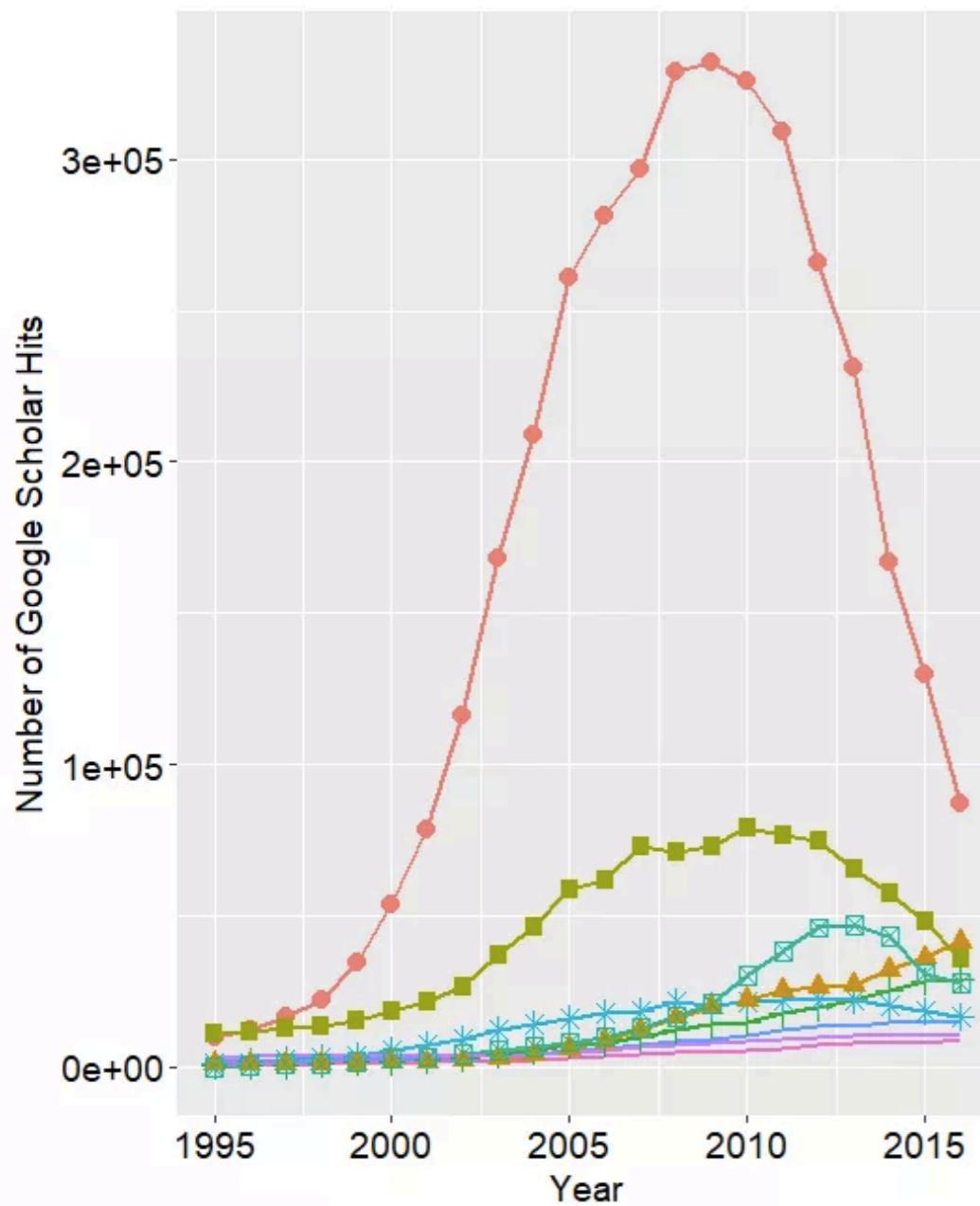
The 2017 UK General Election

Doug Rivers, YouGov's chief scientist, sets out how YouGov's 2017 General Election model works

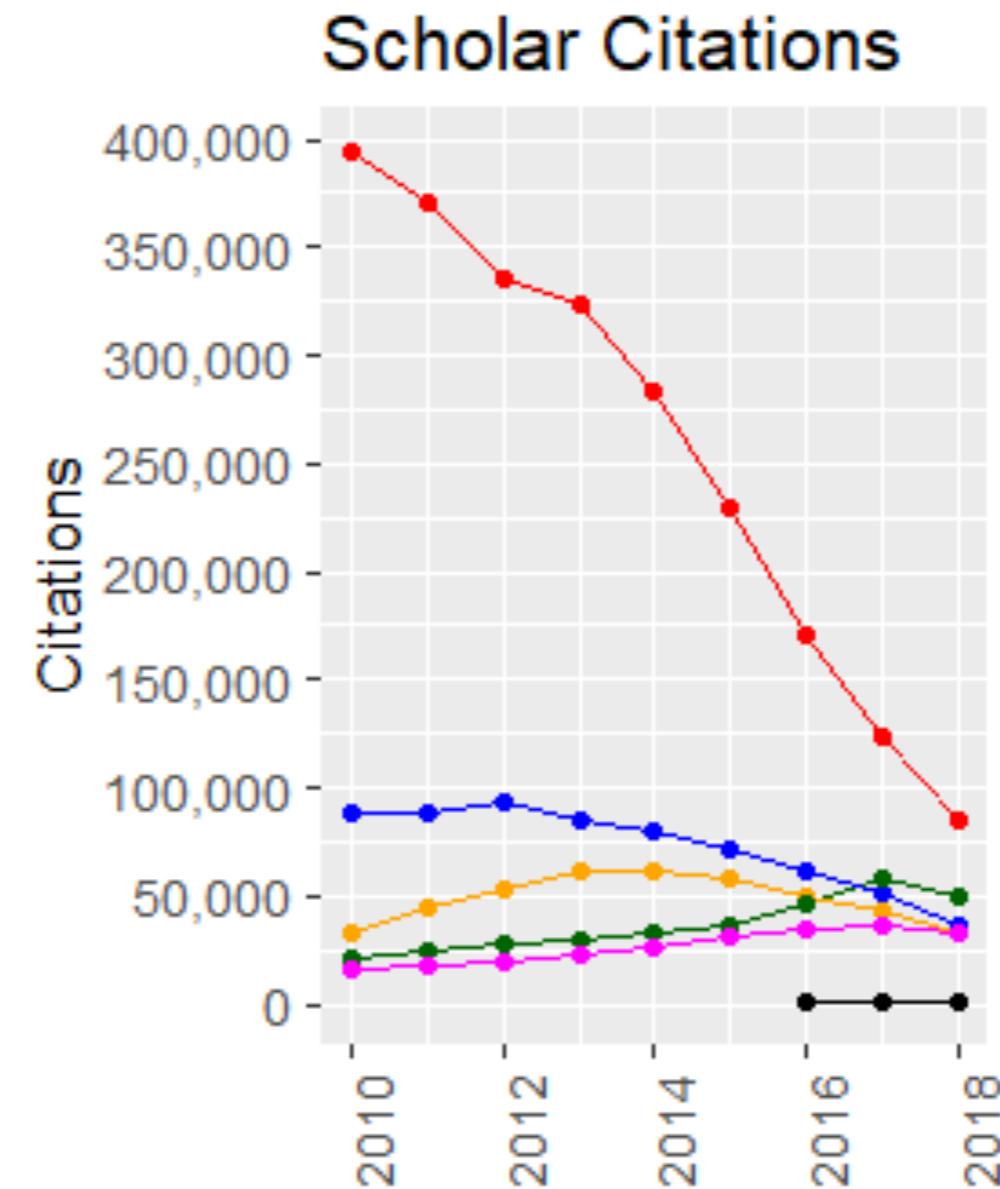
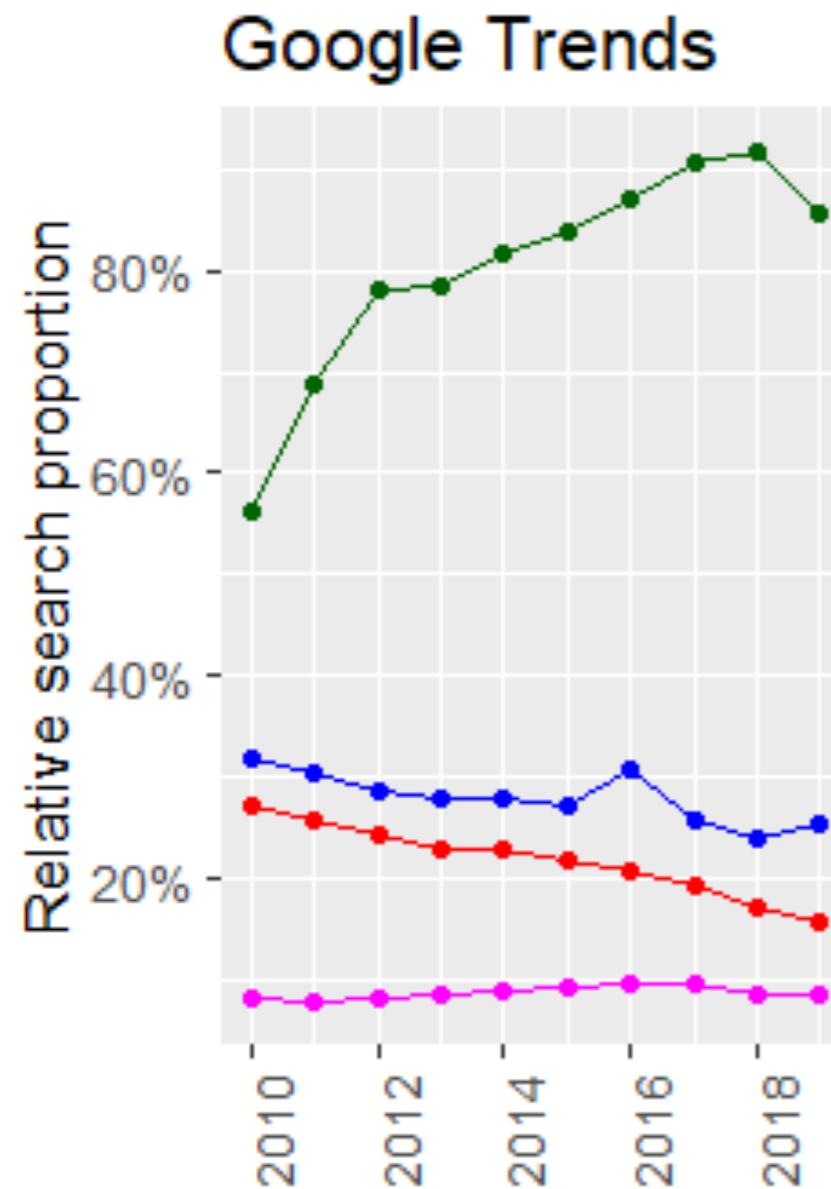
Follow @YouGov on twitter and stay up to date with the latest news and results



and across academia...



Google Search Trends and citations in the academic literature.



software

- SPSS
- R
- SAS
- STATA
- JASP
- Prism

with Data Science a growing employment destination for Psychologists and others.



AMERICAN PSYCHOLOGICAL ASSOCIATION

ABOUT APA

TOPICS

PUBLICATIONS & DATABASES

PSYCHOLOGY HELP CENTER

NEWS & EVENTS

SCIENCE

[Home](#) // [gradPSYCH Magazine](#) // [January 2013 gradPSYCH](#) // Hot jobs: Big-data psychologists

CAREER CENTER

Hot jobs: Big-data psychologists

Wanted: Scientists who can help companies make sense of consumer and employee data.



By Rebecca Voelker

Print version: page 18

Every time you use an Internet search engine, sign up for a company "rewards" program or swipe your credit card, that information is saved and stored. The industries that amass these billions of bytes of data are increasingly hiring psychologists to help make sense of it, says Douglas Reynolds, PhD, president of APA's Div. 14 ([Society for Industrial and Organizational Psychology](#)).

"Big data, and what it means for business, is a hot topic right now," says Reynolds, who also serves as vice president of assessment technology at [Development Dimensions International Inc.](#), a human resources consulting firm. "We're in high demand."

Psychologists' ability to interpret numbers and human behavior makes them key members of many industry analytics teams, adds Suzie Weaver, PhD, a psychologist and senior analytic consultant with [Epsilon](#), a global marketing and analytics company. "Analytics is at the core of everything we do, whether it's in research, the academic sphere or the business world."

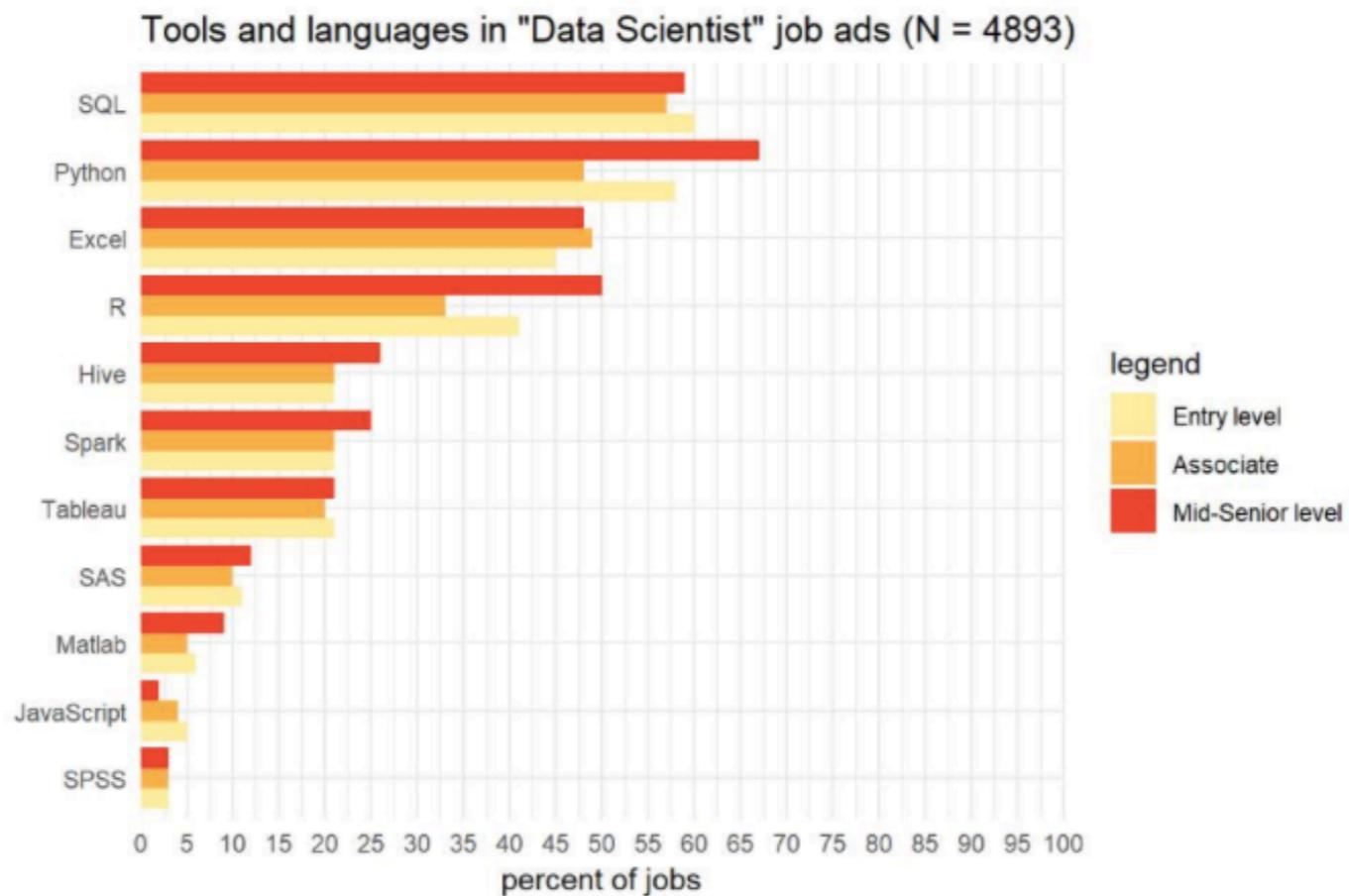
Tyler Burleigh
@tylerburleigh

▼

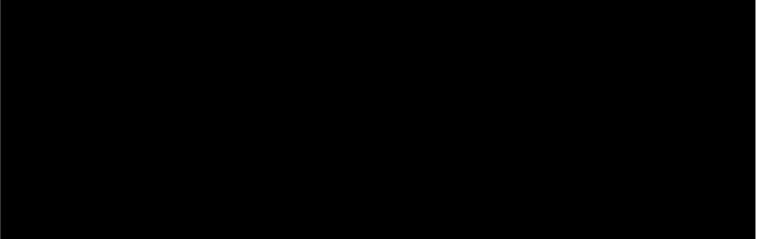
If your department teaches quantitative skills in SPSS, SAS, or Matlab, is it really preparing students for quantitative jobs outside of academia? [#datascientist](#) [#academia](#) [#rstats](#) [#python](#)

Data scientist jobs

4893 data scientist job ads were included in this analysis.



Jobs were scraped from LinkedIn in Sept 2019
Locations included NYC, SF, Seattle, Boston, and Toronto
Positions were Entry, Associate, and Senior levels

 
#brainhackschool instructors be like: familiar with docker?

Me: 😬

Familiar with jupyter?

Me: 😬

Familiar with github?

Me: 😬

Familiar with binder?

Me: 😬

Familiar with python?

Me: 😬

It feels like I've spent these past 4 years of PhD on mars



7:05 PM · Aug 7, 2019 · Twitter for iPhone

How to create BBC style graphics

Load all the libraries you need

Install the bbplot package

How does the bbplot package work?

Save out your finished chart

Make a line chart

Make a multiple line chart

Make a bar chart

Make a stacked bar chart

Make a grouped bar chart

Make a dumbbell chart

Make a histogram

Make changes to the legend

Make changes to the axes

Add annotations

Work with small multiples

Do something else entirely

BBC Visual and Data Journalism cookbook for R graphics

Last updated: 2019-01-24

How to create BBC style graphics

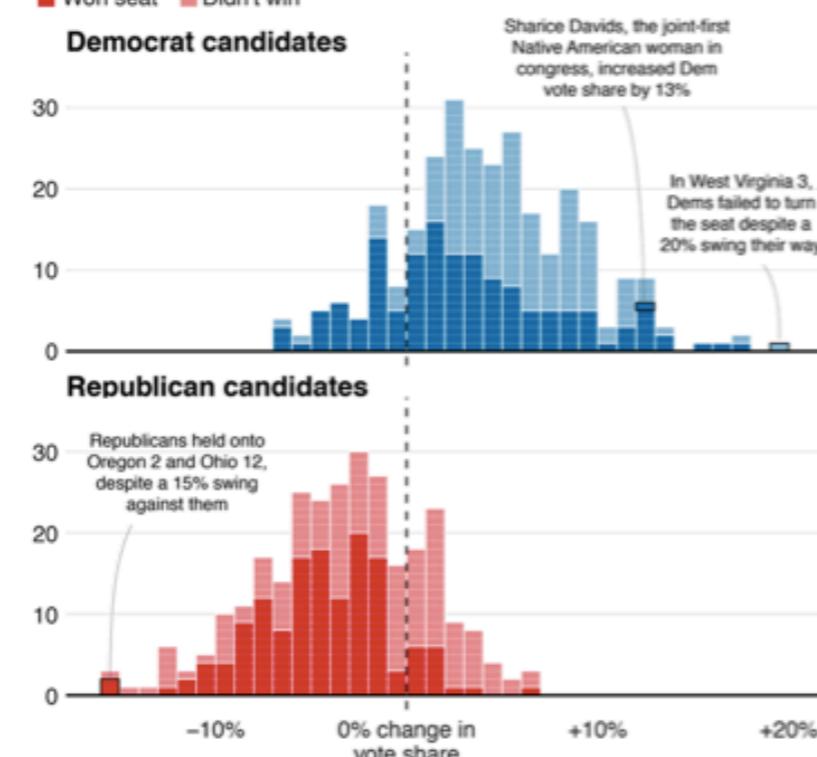
At the BBC data team, we have developed an R package and an R cookbook to make the process of creating publication-ready graphics in our in-house style using R's ggplot2 library a more reproducible process, as well as making it easier for people new to R to create graphics.

The cookbook below should hopefully help anyone who wants to make graphics like these:

Blue wave

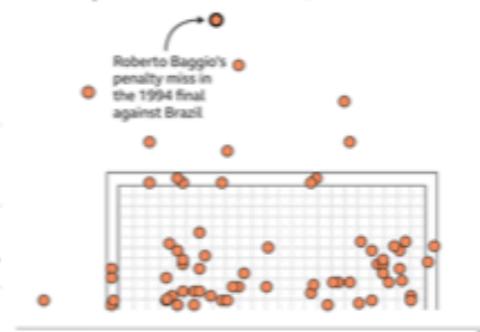
■ Won seat ■ Didn't win

Democrat candidates



Where penalties are saved

World Cup shootout misses and saves, 1982-2014

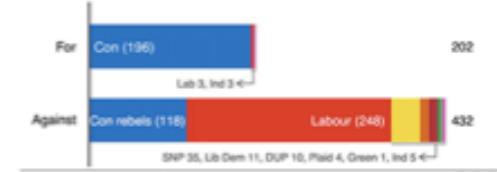


Fast-growing cities face worse climate risks

Population growth 2018-2035 over climate change vulnerability

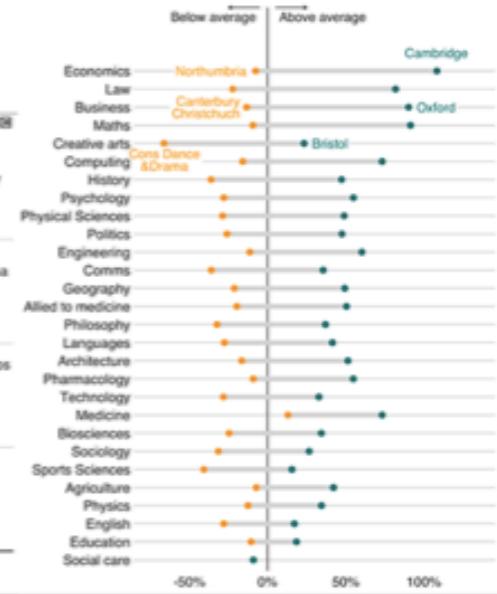


MPs rejected Theresa May's deal by 230 votes



Earnings vary across unis even within subjects

Impact on men's earnings relative to the average degree



<https://medium.com/bbc-visual-and-data-journalism/how-the-bbc-visual-and-data-journalism-team-works-with-graphics-in-r-ed0b35693535>

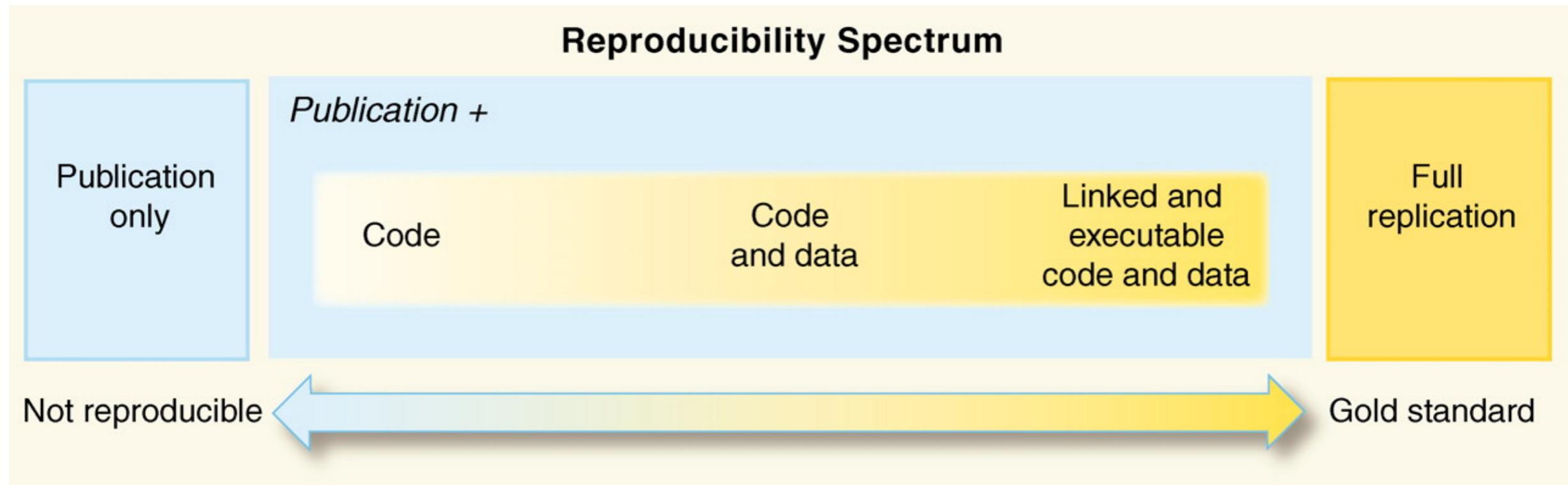
PERSPECTIVE

Reproducible Research in Computational Science

Roger D. Peng

[+ See all authors and affiliations](#)

Science 02 Dec 2011;
Vol. 334, Issue 6060, pp. 1226-1227
DOI: 10.1126/science.1213847



"You can't do data science in a GUI", Hadley Wickham

There will be coding....



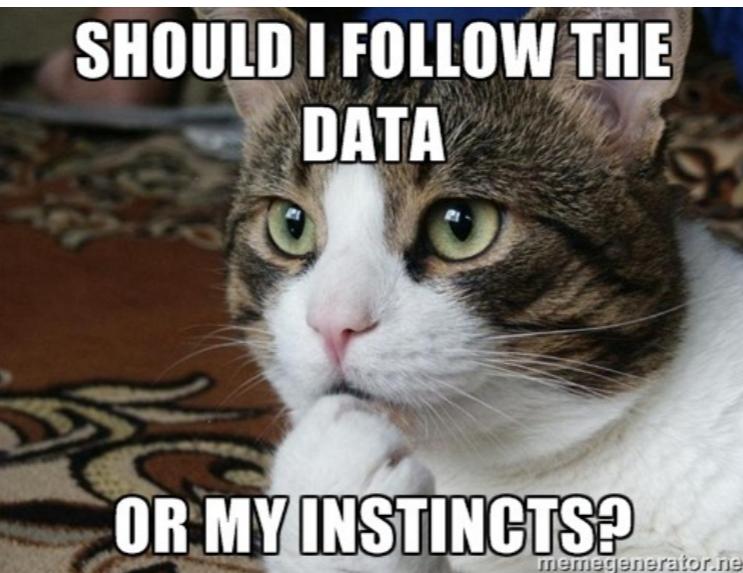
Kelly Vaughn
@kvly

How I code:

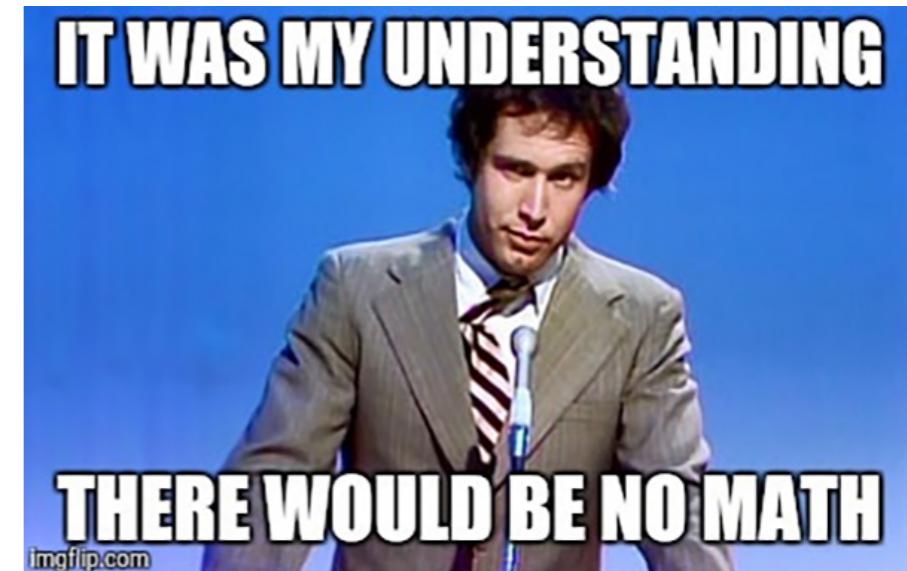
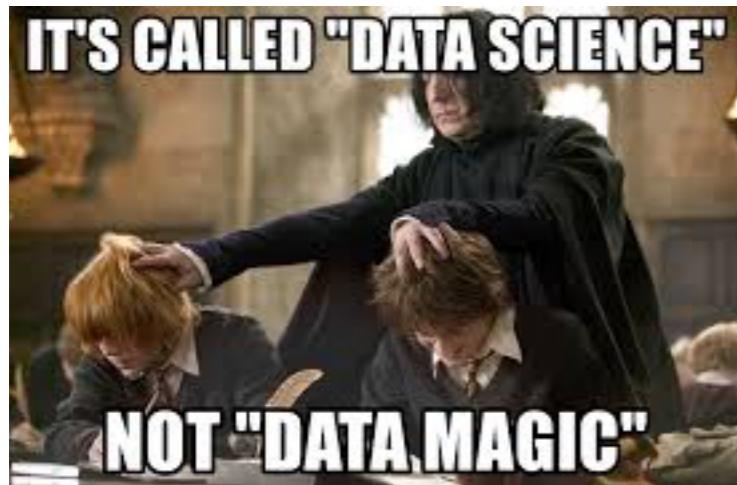
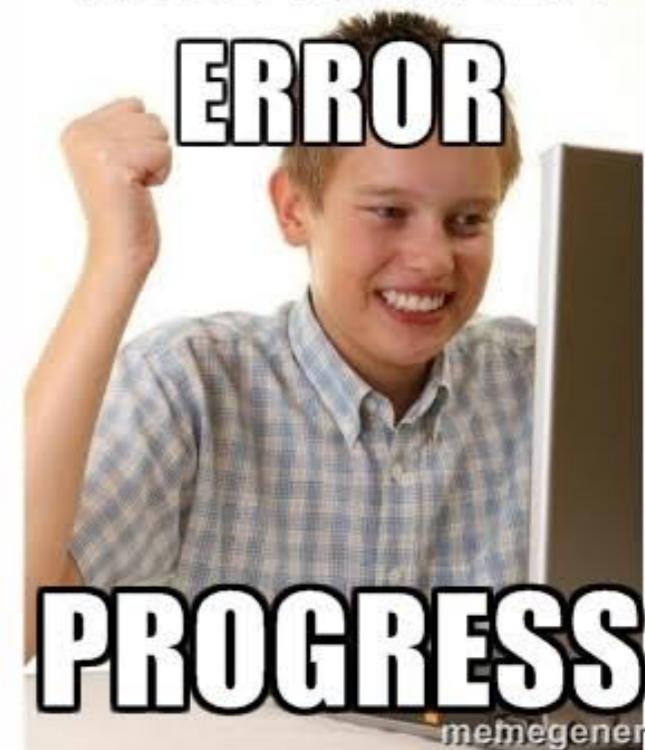
1. Make it work
2. Make it pretty
3. Make it work again because I probably broke it while trying to make it pretty

2:27 PM · Sep 13, 2019 · Twitter Web App

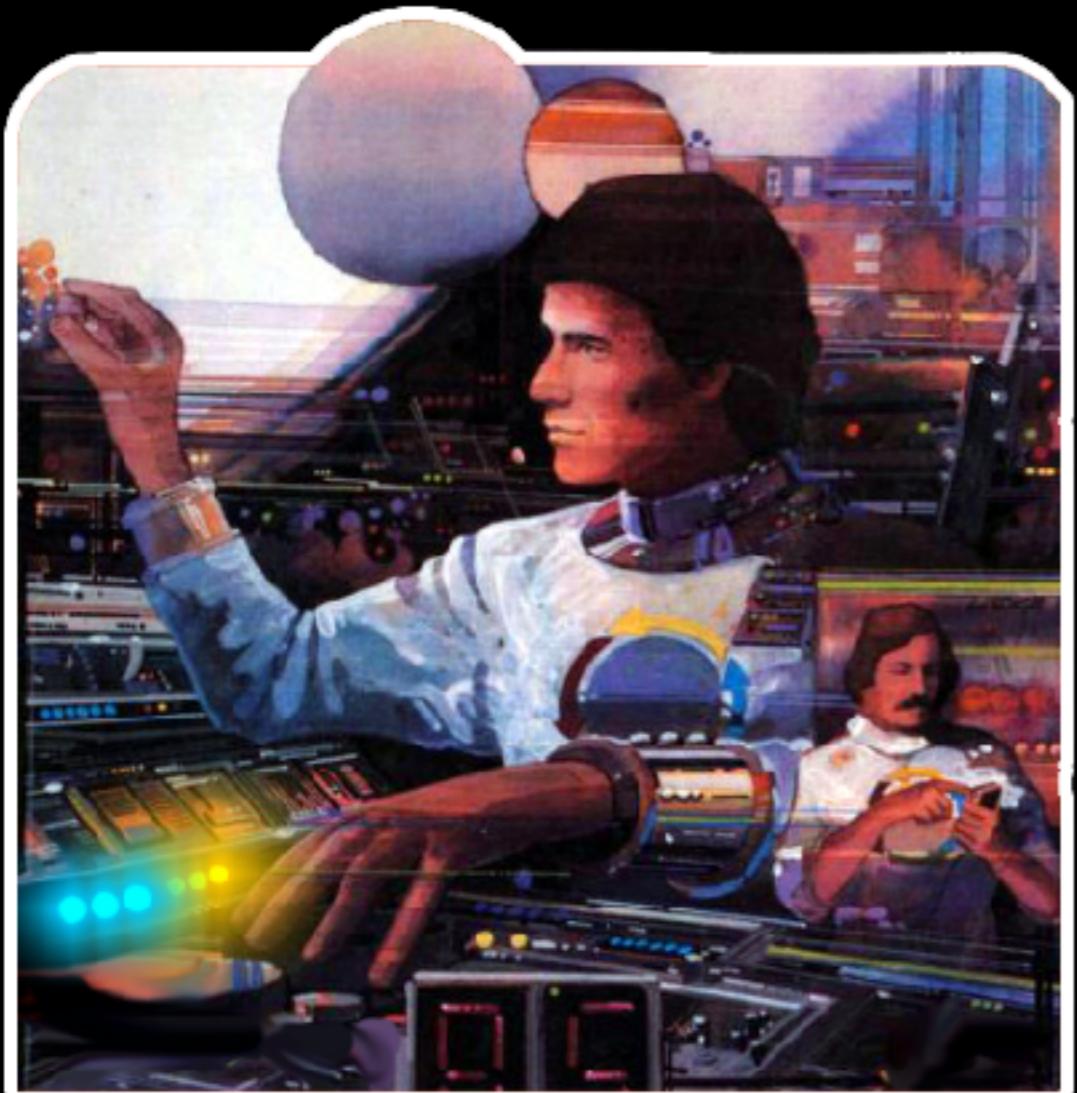
5.3K Retweets 32.4K Likes



GOT A NEW
ERROR



THE TWO STATES OF EVERY PROGRAMMER



I AM A GOD.



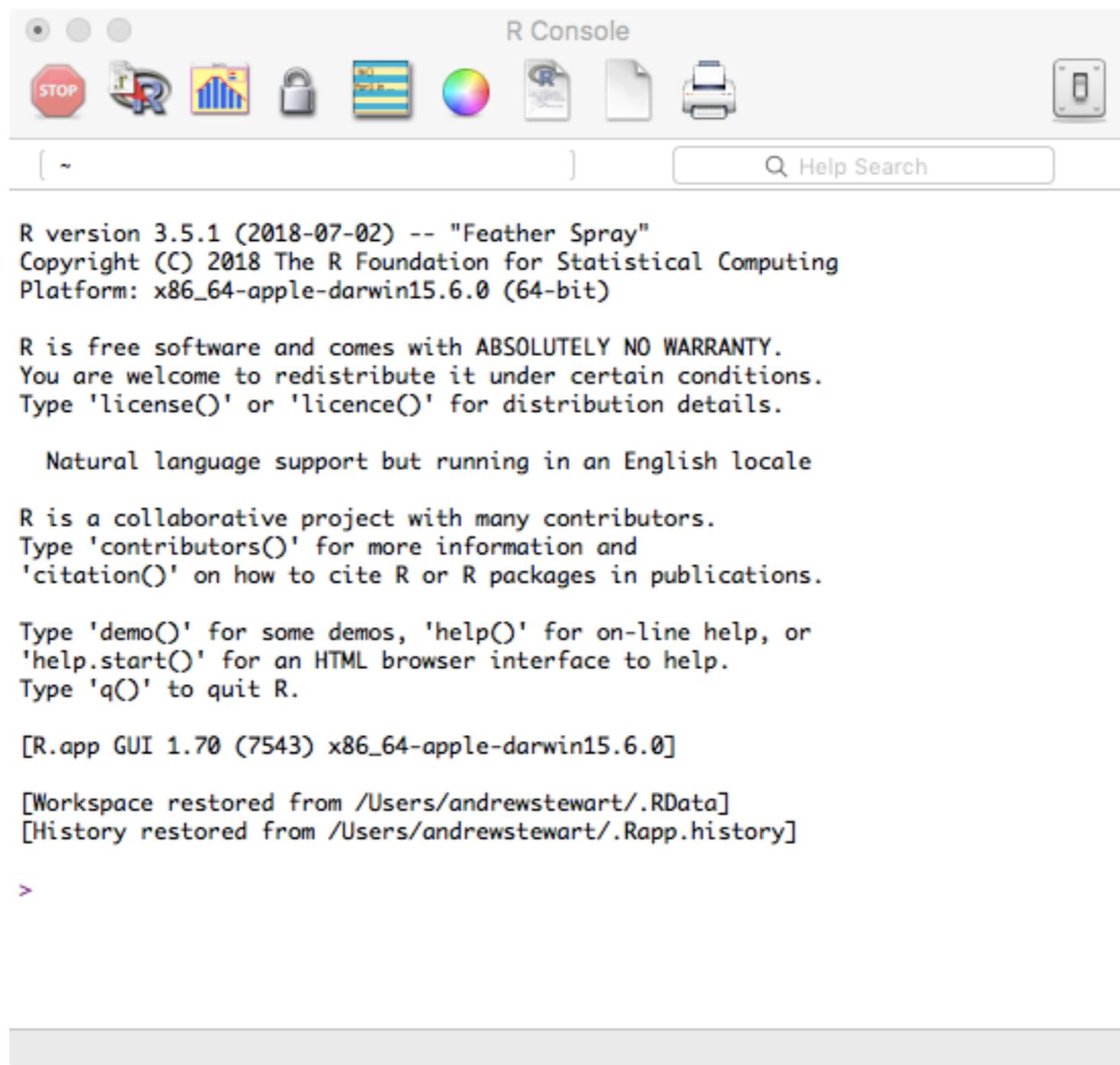
I HAVE NO IDEA
WHAT I'M DOING.

Starting R

- R is open source (i.e., free to download and add to). Main R site is:
- *www.r-project.org*
- From here you can download R (from one of the CRAN¹ mirrors for Windows, Mac and UNIX).
- R updates regularly (you need to update manually).

¹CRAN = *The Comprehensive R Archive Network*

- When you first load R it looks like this:



R version 3.5.1 (2018-07-02) -- "Feather Spray"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-apple-darwin15.6.0 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[R.app GUI 1.70 (7543) x86_64-apple-darwin15.6.0]

[Workspace restored from /Users/andrewstewart/.RData]
[History restored from /Users/andrewstewart/.Rapp.history]

>



- Rather than using the R interface, you should use the RStudio graphical interface.
- Download it from www.rstudio.com
- When you use RStudio, it looks like this:

RStudio File Edit Code View Plots Session Build Debug Tools Window Help

~/Desktop/Air Work/MRes 2016:17/R workshop MRes/R workshop directory/Workshop - RStudio

Addins

Workshop script2.R * DV

```

1 library(lmerTest)
2 library(lsmeans)
3 library(pbkrtest)
4
5 #Read in First Pass Data
6 FPs <- read.csv("~/Desktop/Air Work/R analyses/Indirect Request Expt/Experiment 1 - probability of success - Libby's data/FPs")
7
8 #this sets up the contrasts so that the intercept in the mixed LMM is the grand mean (i.e., the mean of all conditions)
9 my.coding <- matrix(c(.5, -.5))
10
11 contrasts(DV$Context)<-matrix(c(.5, -.5))
12 contrasts(DV$Sentence)<-matrix(c(.5, -.5))
13
14 #construct the models with crossed random effects for subjects and items for the pre-critical, critical and post-critical regions
15 model.full <- lmer(RT ~ Context*Sentence + (1+Context*Sentence | Subject) + (1+Context*Sentence | Item), data=DV, REML=TRUE)
16 model.null <- lmer(RT ~ (1+Context*Sentence | Subject) + (1+Context*Sentence | Item), data=DV, REML=TRUE)
17
18 summary(model.full)
19 lsmeans(model.full, pairwise~Context*Sentence, adjust="none")
20
21 model.full <- lmer(RT ~ Statement*Meaning*Probability + (1:Meaning*Probability | P_c) + (1:Meaning*Probability | Item), data=FPs)
22
23 (Top Level) ▾
  
```

Console ~ /Desktop/Air Work/MRes 2016:17/R workshop MRes/R workshop directory/Workshop/

The following object is masked from 'package:stats':

```

step
  
```

```

> library("lsmeans", lib.loc="/Library/Frameworks/R.framework/Versions/3.3/Resources/library")
Loading required package: estimability

Attaching package: 'lsmeans'

The following object is masked from 'package:lmerTest':
  
```

```

lsmeans
  
```

```

> model.full <- lmer(RT ~ Context*Sentence + (1+Context*Sentence | Subject) + (1+Context*Sentence | Item), data=DV, REML=TRUE)
  
```

Error in strsplit(keys, " * ") : non-character argument

Environment History

Import Dataset

Global Environment

Data

	DV	1680 obs. of 5 variables
my.coding	num [1:2, 1]	0.5 -0.5

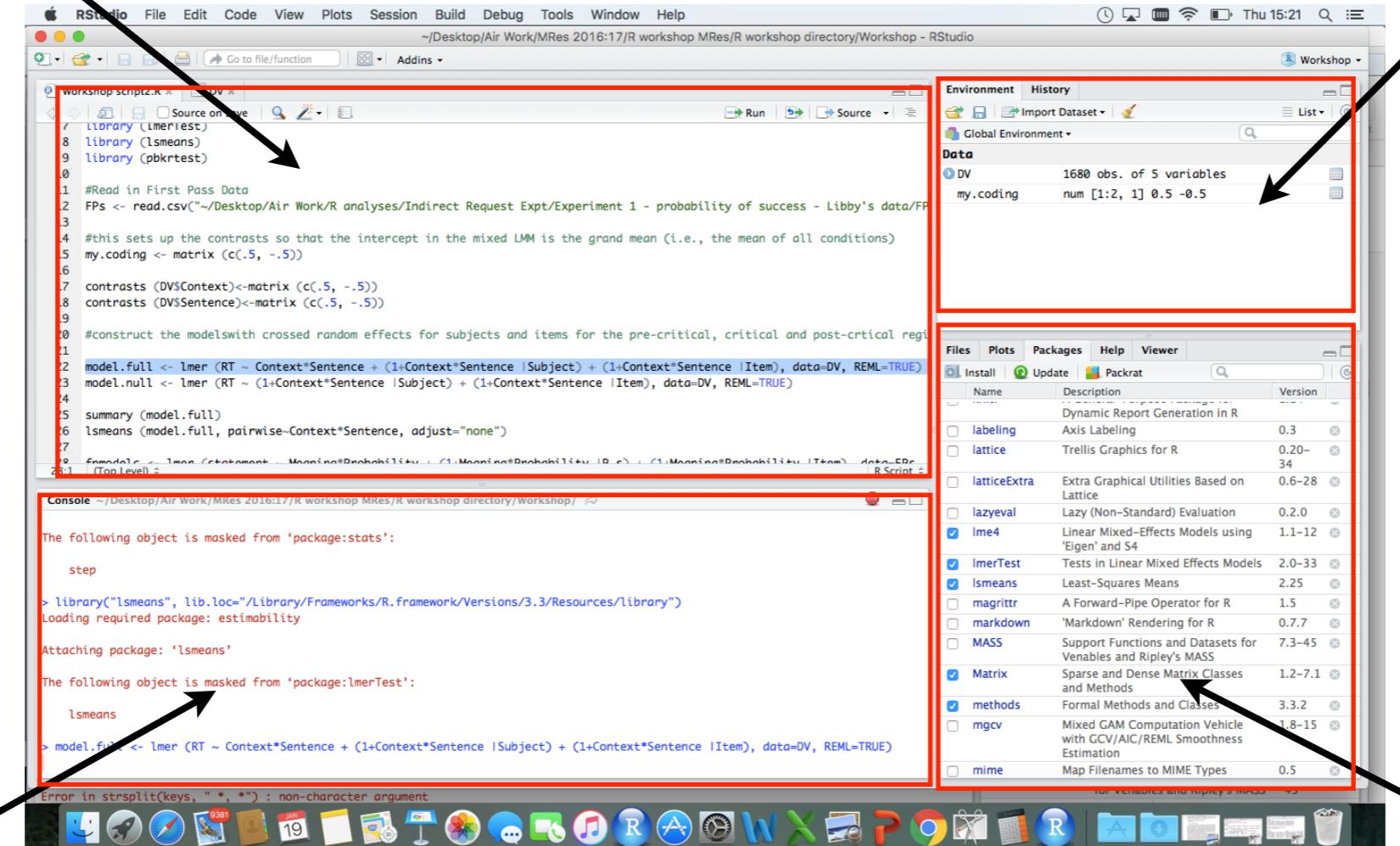
Files Plots Packages Help Viewer

Install Update Packrat

Name	Description	Version
lme4	Linear Mixed-Effects Models using 'Eigen' and S4	1.1-12
lmerTest	Tests in Linear Mixed Effects Models	2.0-33
lsmeans	Least-Squares Means	2.25
MASS	Support Functions and Datasets for Venables and Ripley's MASS	7.3-45
Matrix	Sparse and Dense Matrix Classes and Methods	1.2-7.1
methods	Formal Methods and Classes	3.3.2
mgcv	Mixed GAM Computation Vehicle with GCV/AIC/REML Smoothness Estimation	1.8-15
mime	Map Filenames to MIME Types	0.5

for Venables and Ripley's MASS 45

This is where you build your script and where data can be seen.

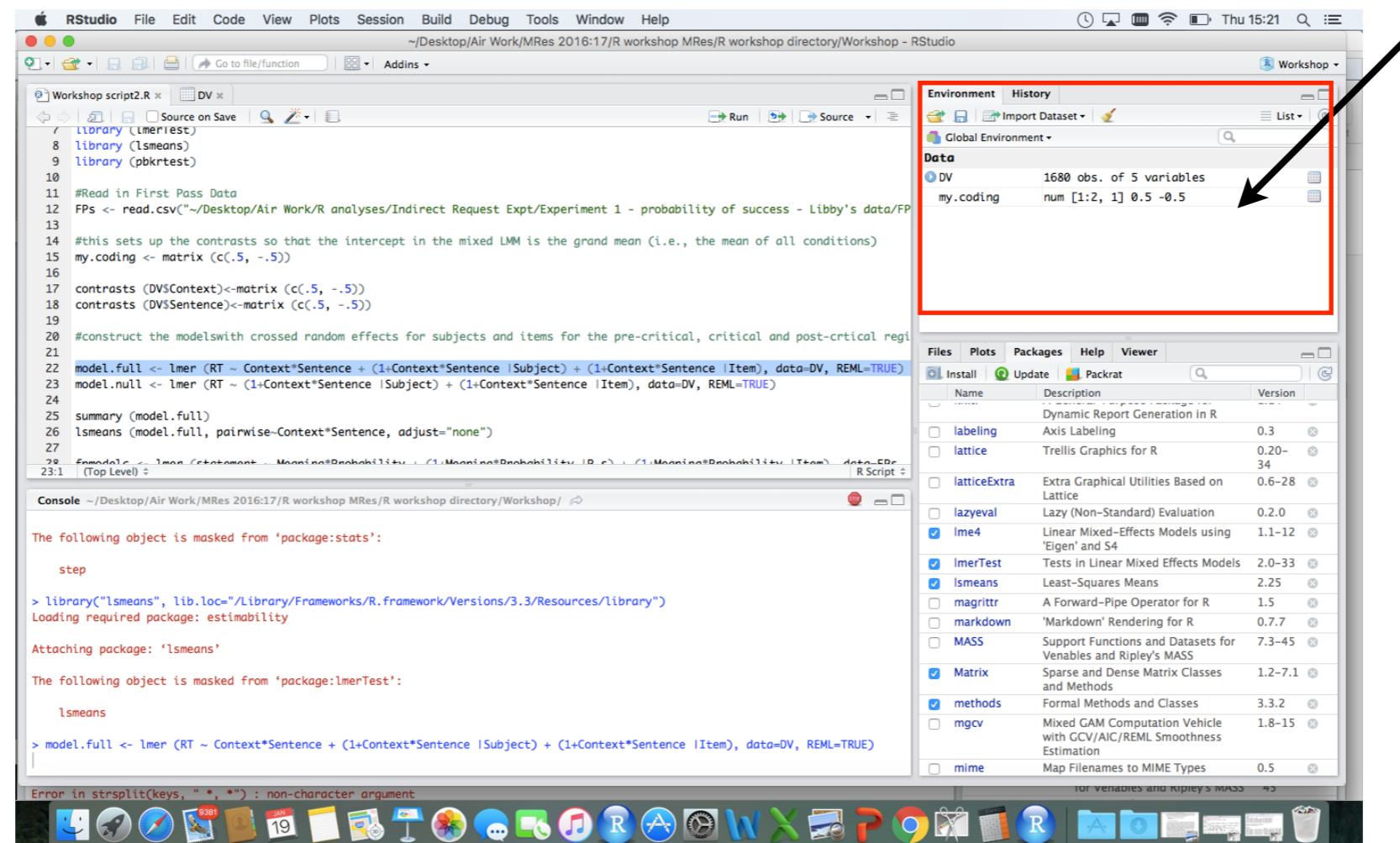


This is where you type commands.

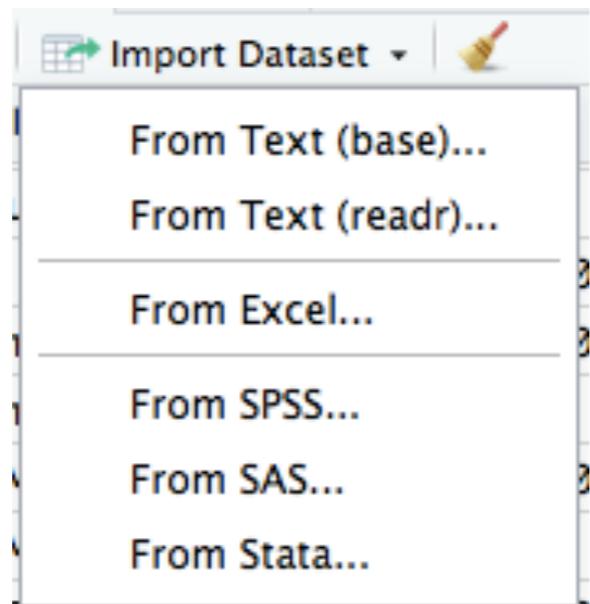
This contains information about your variables and open data sets.

This lists the packages you have loaded, and has tabs for help, graphs etc.

Here is where you import your data.



- Importing data (CSV (Text), Excel, SPSS, SAS and Stata format). CSV can be delimited by commas, tabs etc. Most of the time you'll use the `readr` import function...



Import Text Data

File/Url:

~/Desktop/Air Work/R analyses/R courses/R course/R Work Folder/priming data.txt

Data Preview:

Subject (integer) ▾	Priming Condition (character) ▾	RT (integer) ▾
1	LONG	500
2	LONG	551
3	LONG	479
4	LONG	561
5	LONG	522
6	SHORT	399
7	SHORT	423
8	SHORT	444
9	SHORT	410
10	SHORT	398

You can change the type of each variable by clicking on the down arrow.

You should change this to type Factor (LONG vs. SHORT).

Previewing first 50 entries.

Import Options:

Name: priming_data First Row as Names Delimiter: Tab Escape: None
Skip: 0 Trim Spaces Quotes: Default Comment: Default
 Open Data Viewer Locale: Configure... NA: Default

Code Preview:

```
library(readr)
priming_data <- read_delim("~/Desktop/Air Work/R analyses/R courses/R course/R Work Folder/priming data.txt",
                           "\t", escape_double = FALSE, trim_ws = TRUE)
View(priming_data)
```

Can change here how the columns are delimited.

This is the code that corresponds to what you're getting RStudio to do.

- RStudio now looks like this:

The screenshot shows the RStudio interface with the following components:

- Data View:** Displays a data frame named "priming_data" with 10 observations and 3 variables: Subject, Priming Condition, and RT.
- Environment View:** Shows the global environment with "priming_data" listed.
- Console View:** Displays the R session history, including the command to read the data and its successful execution.
- Packages View:** Shows the installed packages and their details.

```

> 
> 
> 
> 
> 
> 
> 
> RTdata <- priming_data [c("Priming Condition", "RT")]
Error: object 'priming_data' not found
> library(readr)
> priming_data <- read_delim("~/Desktop/Air Work/R analyses/R courses/R course/R Work Folder/priming data.txt",
+   "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Subject = col_integer(),
  `Priming Condition` = col_character(),
  RT = col_integer()
)
> View(priming_data)
> 
  
```

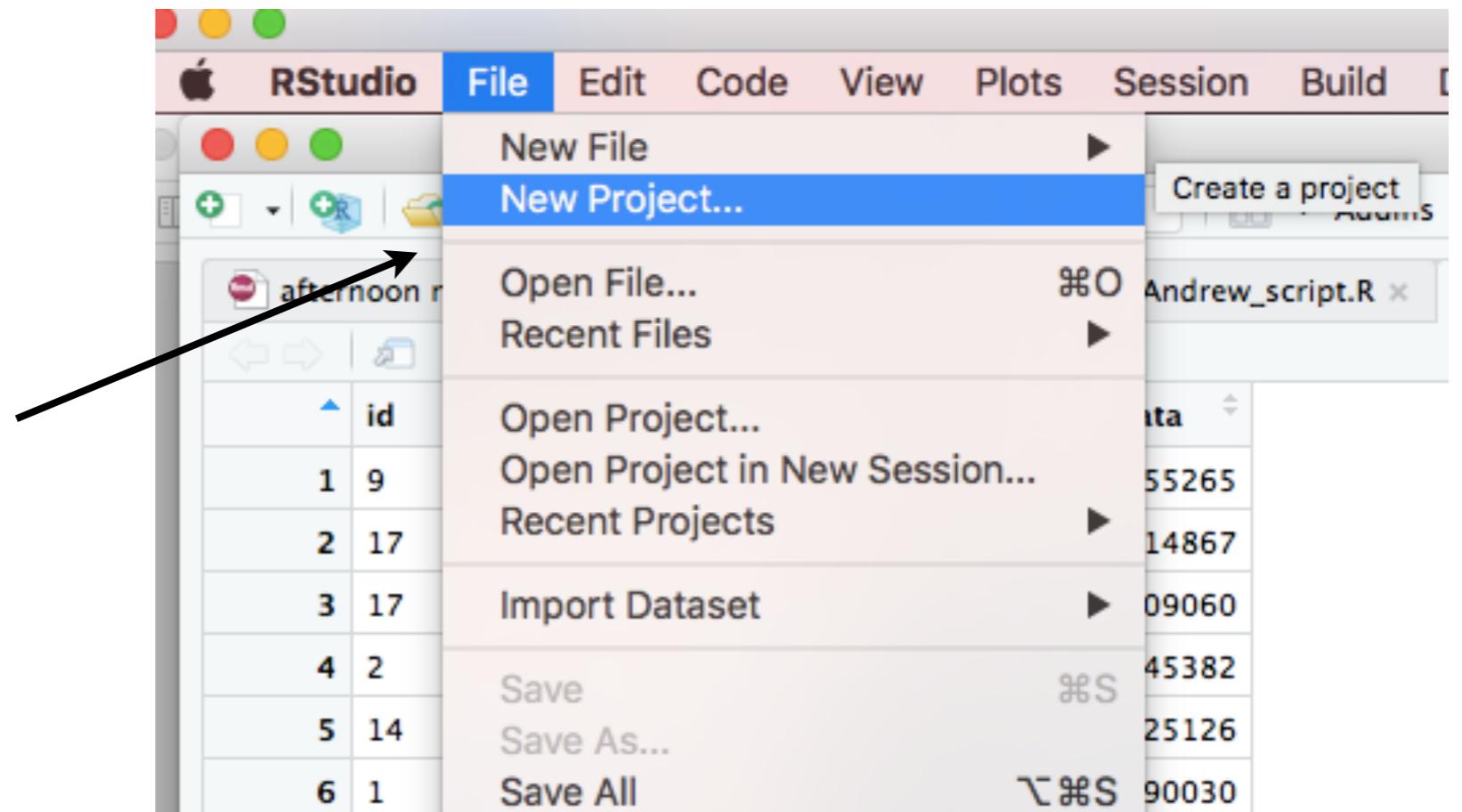
Subject	Priming Condition	RT
1	LONG	500
2	LONG	551
3	LONG	479
4	LONG	563
5	LONG	522
6	SHORT	399
7	SHORT	423
8	SHORT	444
9	SHORT	410
10	SHORT	398

Showing 1 to 10 of 10 entries

Name	Description	Version
labeling	Axis Labeling	0.3
lattice	Trellis Graphics for R	0.20-34
latticeExtra	Extra Graphical Utilities Based on Lattice	0.6-28
lazyeval	Lazy (Non-Standard) Evaluation	0.2.0
lme4	Linear Mixed-Effects Models using 'Eigen' and S4	1.1-12
lmerTest	Tests in Linear Mixed Effects Models	2.0-33
lsmeans	Least-Squares Means	2.25
magrittr	A Forward-Pipe Operator for R	1.5
markdown	'Markdown' Rendering for R	0.7.7
MASS	Support Functions and Datasets for Venables and Ripley's MASS	7.3-45
Matrix	Sparse and Dense Matrix Classes and Methods	1.2-7.1
methods	Formal Methods and Classes	3.3.2
mgcv	Mixed GAM Computation Vehicle with GCV/AIC/REML Smoothness Estimation	1.8-15
mime	Map Filenames to MIME Types	0.5
minqa	Derivative-free optimization algorithms by quadratic approximation	1.2.4
multcomp	Simultaneous Inference in General Parametric Models	1.4-6
munsell	Utilities for Using Munsell Colours	0.4.3

When Starting R

Always create a new project - this will keep all your files nice and organised.

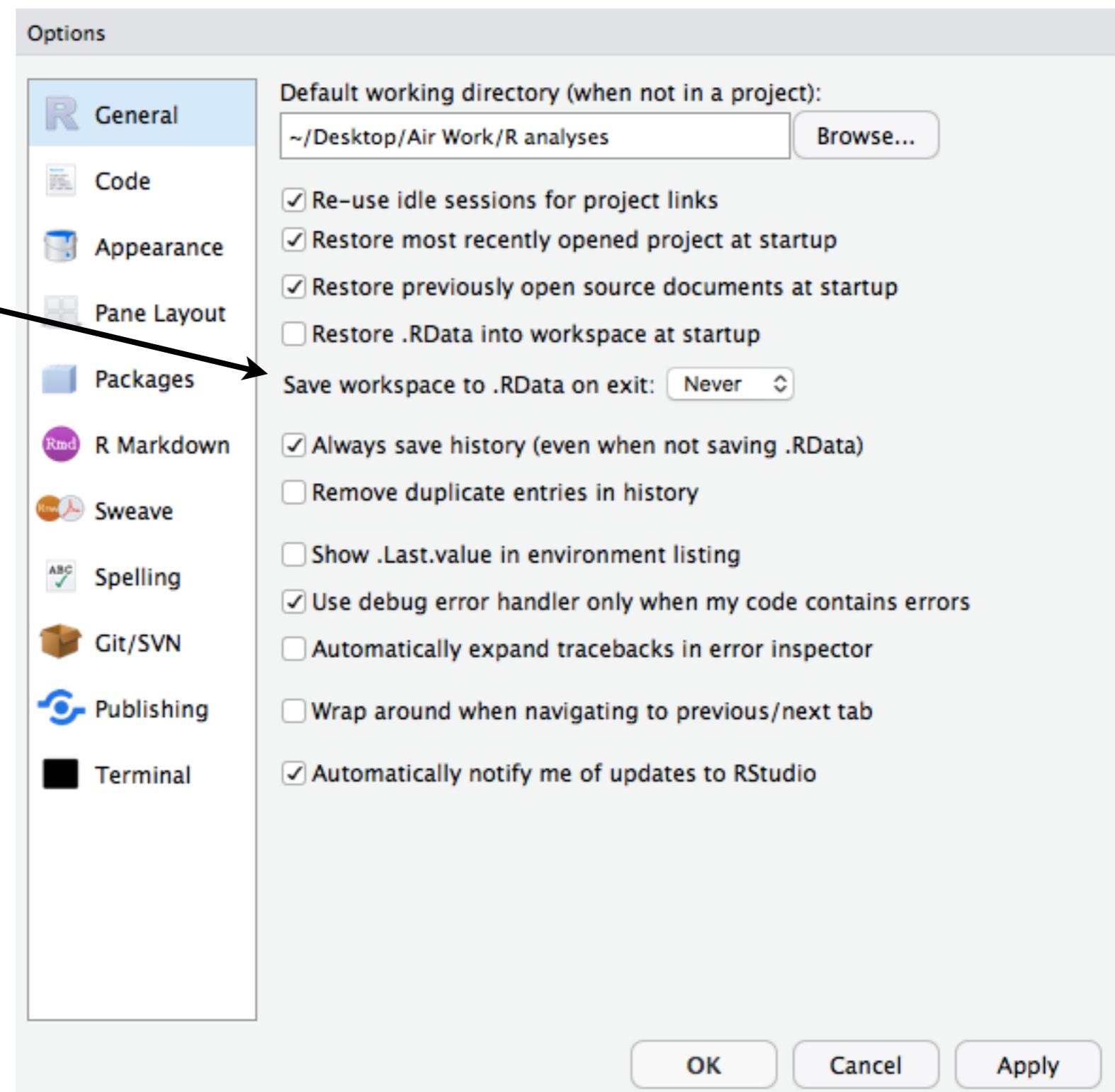


When you create an R Project in a folder, it's a good idea to keep that folder itself nicely organised:

```
name_of_project
| --name_of_project.Rproj
| --raw_data
|   |--data.csv
| --tidied_data
|   |--tidy_data.csv
| --R_scripts
|   |--data_tidy.R
|   |--data_vis.R
|   |--analysis.R
| --markdown_scripts
|   |--analysis.Rmd
| --output_reports
|   |--analysis.html
|   |--analysis.pdf
```

When Starting R

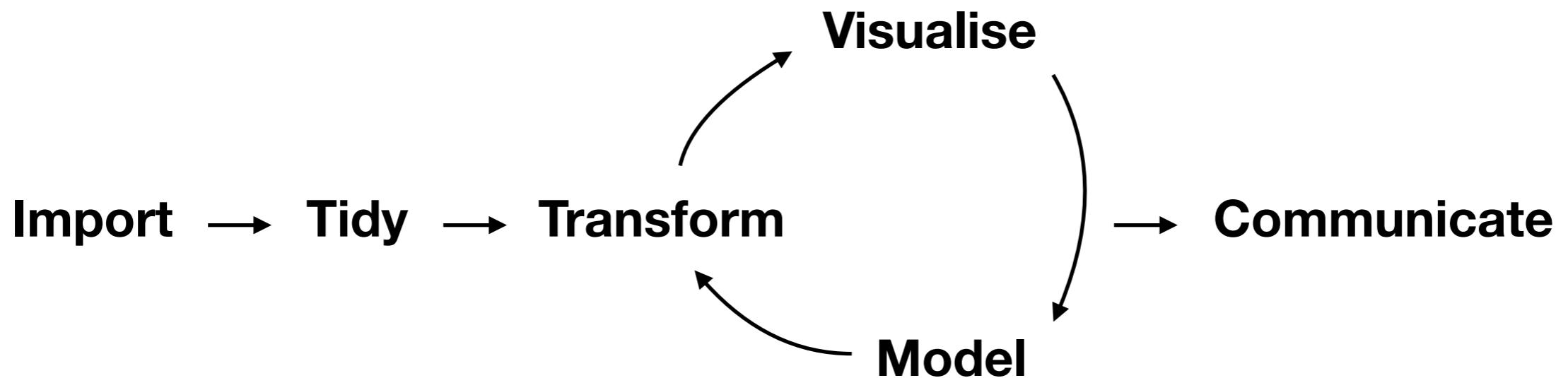
Under preferences,
make sure you
uncheck the saving
workspace
option...



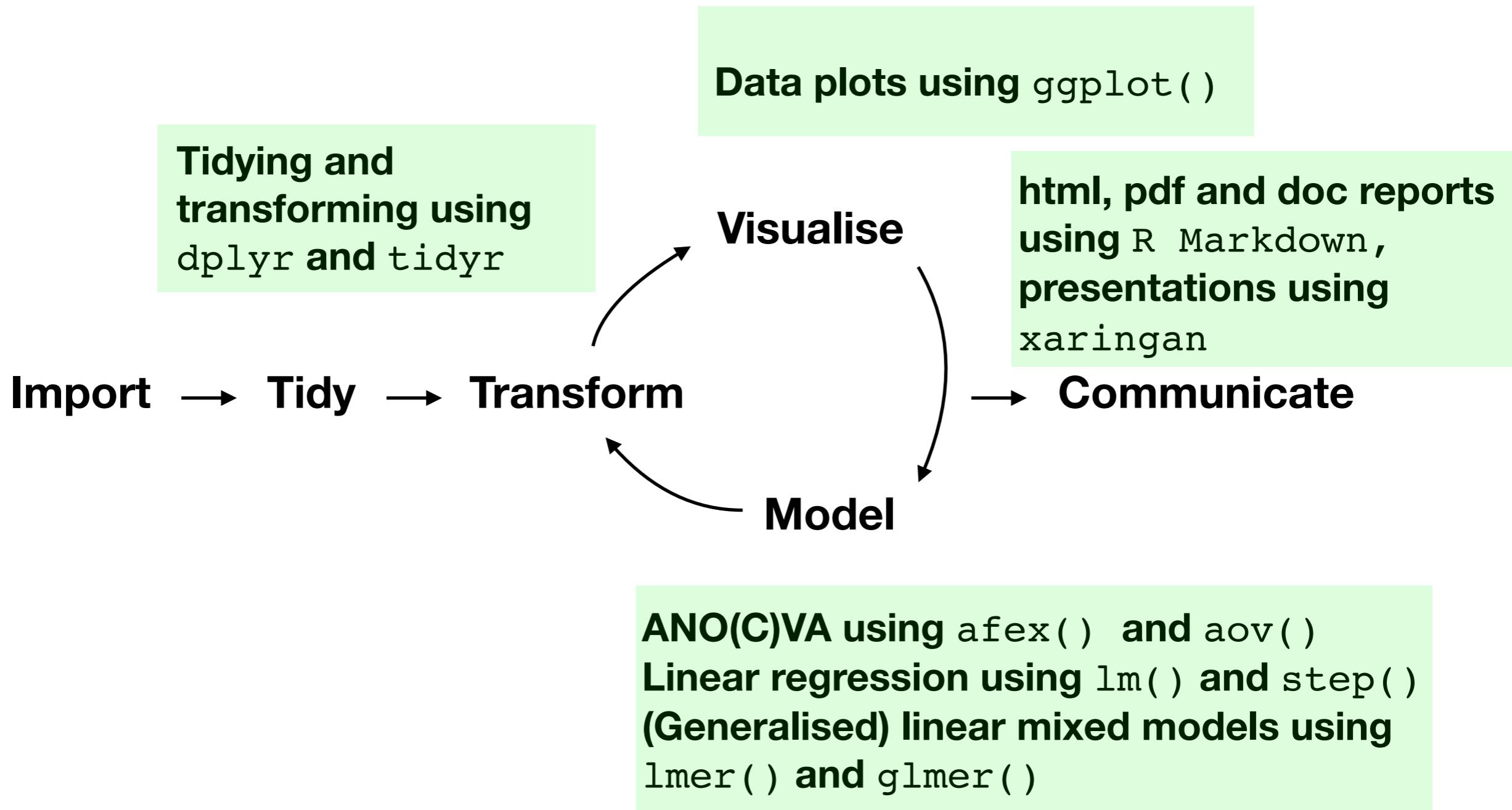
Setting up R

Analysis Workflow in the Tidyverse

(Garrett Grolemund and Hadley Wickham) - from Data to Write-up



Workflow



Packages in R

R has a number of functions already built in. These are part of “base R”. For much of what we do we need to load particular packages to let us use additional functions. We do this by:

```
> install.packages ("packagename")  
  
> library(packagename)
```

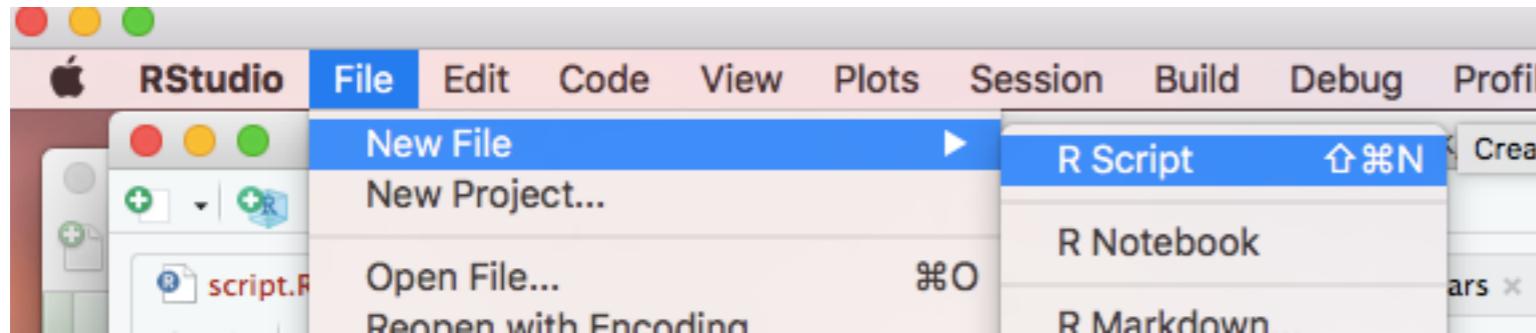
Every time you start R afresh, you need to load the packages you will use by using the `library` function. If you see notation like the following, it means we’re using a function (e.g., `summarise`) from within a package (e.g., `dplyr`)

```
dplyr::summarise
```

You don't need to re-invent the wheel...

- For any problem you want to solve, chances are others have had the same problem, solved it and have created an R package to do exactly what you want...
- One of the many great things about R is that you can add freely available packages to your library.
- ~15,000 R packages currently available
- The sites r-bloggers.com and rweekly.org are good places to find out about new packages.

Writing Scripts in R



Ultimately you will be writing scripts that will allow yourself and others to recreate your analysis just by running the script - the code at the start of the script will load the packages your analysis needs, then load your data, tidy, transform and visualise your data, and then code for your model...

You should write your code as a script, and use the Console window for single command instructions relatively rarely...

Writing Scripts in R

Scripts should have lots of comments, indicated by a # symbol - this helps others read your code, and also yourself when you look back at it later.

Scripts can be saved and uploaded to (e.g.) OSF, GitHub, or submitted as an electronic supplement alongside your journal submission. Even when journals don't require you to adhere to Open Science practices, it's a good idea to make your code and data available during the reviewing process - and publicly available once your paper is published.

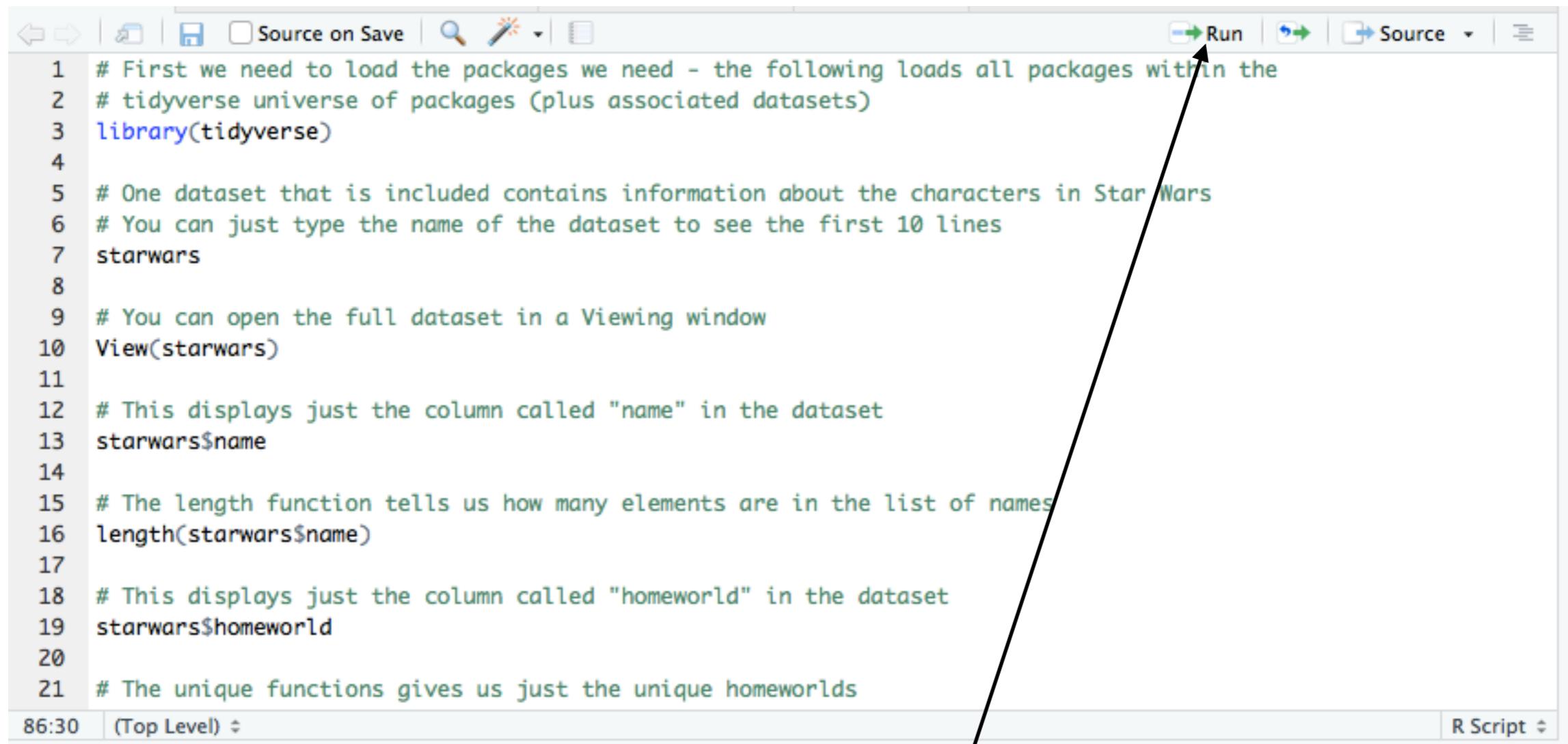
If you aren't sharing your data and code, you aren't engaged in generating reproducible or open research...

The screenshot shows the RStudio interface with an R script editor. The script contains 21 numbered lines of R code, each preceded by a line number and a descriptive comment. The code demonstrates basic data manipulation and exploration using the tidyverse package and the starwars dataset.

```
1 # First we need to load the packages we need - the following loads all packages within the
2 # tidyverse universe of packages (plus associated datasets)
3 library(tidyverse)
4
5 # One dataset that is included contains information about the characters in Star Wars
6 # You can just type the name of the dataset to see the first 10 lines
7 starwars
8
9 # You can open the full dataset in a Viewing window
10 View(starwars)
11
12 # This displays just the column called "name" in the dataset
13 starwars$name
14
15 # The length function tells us how many elements are in the list of names
16 length(starwars$name)
17
18 # This displays just the column called "homeworld" in the dataset
19 starwars$homeworld
20
21 # The unique functions gives us just the unique homeworlds
```

86:30 (Top Level) R Script

Lots of comments explaining what each section of code does, and lots of white space.



```
1 # First we need to load the packages we need - the following loads all packages within the
2 # tidyverse universe of packages (plus associated datasets)
3 library(tidyverse)
4
5 # One dataset that is included contains information about the characters in Star Wars
6 # You can just type the name of the dataset to see the first 10 lines
7 starwars
8
9 # You can open the full dataset in a Viewing window
10 View(starwars)
11
12 # This displays just the column called "name" in the dataset
13 starwars$name
14
15 # The length function tells us how many elements are in the list of names
16 length(starwars$name)
17
18 # This displays just the column called "homeworld" in the dataset
19 starwars$homeworld
20
21 # The unique functions gives us just the unique homeworlds
```

Once a script is written, you can run the entire script by highlighting it all (CMD-A in OSX) and clicking on ‘Run’, or you can highlight a section and run only that. CMD-Return will also Run the highlighted code.

Style Guide

File names (both scripts and data files) should be meaningful:

regression_model.R - Good

somemodel.R - Bad

Variable names should be meaningful and in lowercase:

age - Good

A1 - Bad

expt1_data - Good

Experiment1datawithoutliersremoved - Bad

Style Guide

Place spaces around operators like + , - , = , <-

Whitespace used sensibly makes your code easier to read. Separate discrete sections of code in your script with a blank line.

When writing scripts, continue on a subsequent line rather than writing one very long line.

Comment, comment, comment...

```
# First we load our datafile.
```

<https://style.tidyverse.org>

The screenshot shows the 'R Style Guide' sidebar on the left and the 'The tidyverse style guide' content area on the right.

Left Sidebar:

- R Style Guide
- ☰
- 🔍
- A
- ⤓
- ⤒
- The tidyverse style guide
- Welcome
- I Analyses
- 1 Files
 - 1.1 Names
 - 1.2 Organisation
 - 1.3 Internal structure
- 2 Syntax
 - 2.1 Object names
 - 2.2 Spacing
 - 2.3 Argument names
 - 2.4 Code blocks
 - 2.5 Long lines
 - 2.6 Assignment
 - 2.7 Semicolons
 - 2.8 Quotes
 - 2.9 Comments
- 3 Functions
 - 3.1 Naming

Right Content Area:

The tidyverse style guide

Hadley Wickham

Welcome

Good coding style is like correct punctuation: you can manage without it, but it sure makes things easier to read. This site describes the style used throughout the [tidyverse](#). It was derived from Google's original R Style Guide - but Google's [current guide](#) is derived from the tidyverse style guide.

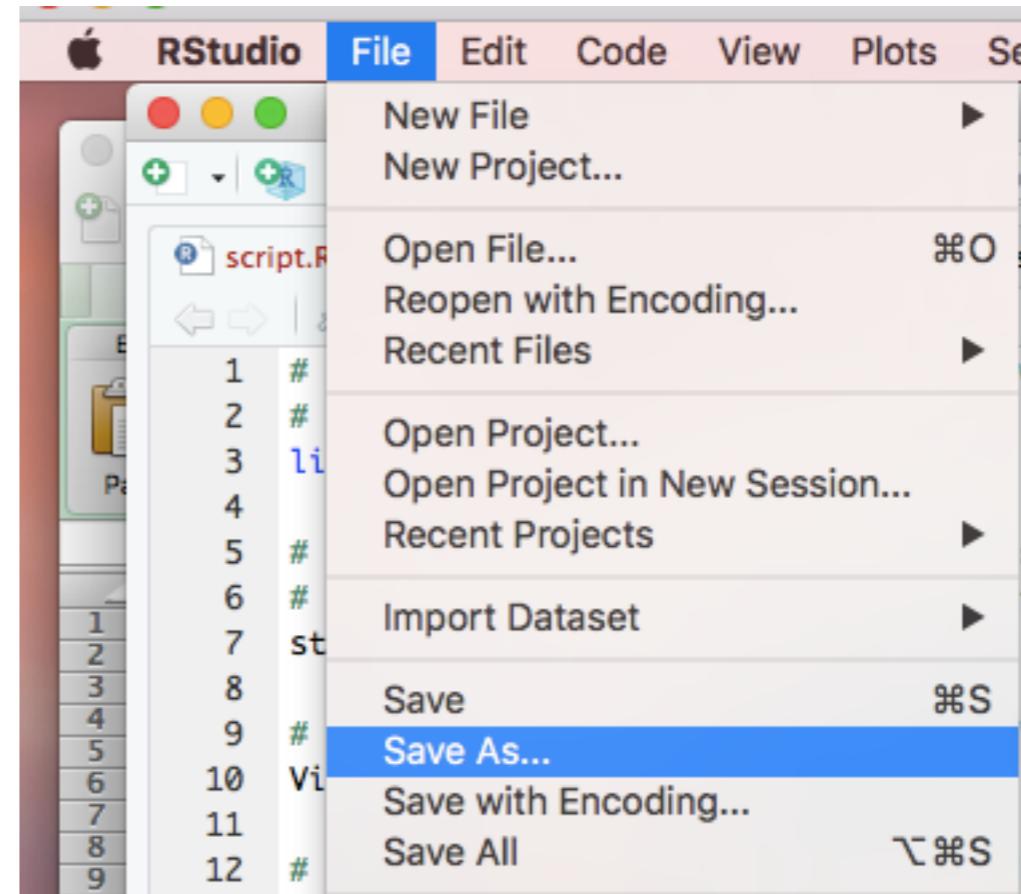
All style guides are fundamentally opinionated. Some decisions genuinely do make code easier to use (especially matching indenting to programming structure), but many decisions are arbitrary. The most important thing about a style guide is that it provides consistency, making code easier to write because you need to make fewer decisions.

Two R packages support this style guide:

- [styler](#) allows you to interactively restyle selected text, files, or entire projects. It includes an RStudio add-in, the easiest way to re-style existing code.

The screenshot shows the RStudio interface with the 'Addins' menu open. The 'STYLER' option is highlighted in the dropdown menu.

Once your
script is
finished, don't
forget to save
it...



Sharing your analysis

Journals are increasingly requesting analysis code and data to be published alongside the published paper (and sometimes at the review stage).

You can share your analysis and data even at the submission stage using something like GitHub or OSF.

Git is a powerful tool that allows for the version control in collaborative projects, and the sharing of code and projects.

You can set up a GitHub account, and install GitHub Desktop on your own computer.

Google Scholar Scopus jobs.ac.uk BBC News Chester Weather The Grauniad The Independent Google Maps Chester Weather Station GitHub

Search or jump to... / Pull requests Issues Marketplace Explore + 

 Overview Repositories 34 Projects 0 Stars 9 Followers 5 Following 1

Pinned Customize your pins

 Psychology_MRes_Stats_R_Course Slides for my MRes Stats Course   

 Affective-Theory-of-Mind-Inferences R code and data for Stewart, S.L.K., Schepman, A., Haigh, M., McHugh, R., & Stewart, A.J. (2019). Affective theory of mind inferences contextually influence the recognition of emotional facial expr...  

 Comprehension-of-indirect-requests-is-influenced-by-their-degree-of-imposition R code and data for Stewart, A.J., Le-luan, E., Yao, B., Wood, J., & Haigh, M., (2018). Comprehension of indirect requests is influenced by their degree of imposition. Discourse Processes, 55, 187-... 

 It's-hard-to-write-a-good-article R code and data for Stewart, A.J., Wood, J.S., Le-luan, E., Yao, B., & Haigh, M. (2018). "It's hard to write a good article." The online comprehension of excuses as indirect replies. Quarterly Jour... 

 Psychophysiology2017 Forked from RonanMcG/Psychophysiology2017 

 ajstewartlang.github.io



Google Scholar Scopus jobs.ac.uk BBC News Chester Weather The Grauniad The Independent Google Maps Chester Weather Station GitHub

Search or jump to... Pull requests Issues Marketplace Explore

ajstewartlang / Affective-Theory-of-Mind-Inferences Watch 0 Star 1 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights Settings

R code and data for Stewart, S.L.K., Schepman, A., Haigh, M., McHugh, R., & Stewart, A.J. (2019). Affective theory of mind inferences contextually influence the recognition of emotional facial expressions. Cognition & Emotion, 33, 272-287. Edit

Manage topics

15 commits 1 branch 0 releases 1 contributor

Branch: master New pull request Create new file Upload files Find File Clone or download

 ajstewartlang	Update README.txt	Latest commit e4b4ff4 28 days ago
	Analysis with emotional word items excluded	typos corrected
	Analysis with emotional word items excluded_old	final changes
	Both Expts comparison	final changes
	Expt 1 Script and data	final changes
	Expt 2 Script and data	final changes
	README.txt	Update README.txt

Google Scholar Scopus jobs.ac.uk BBC News Chester Weather The Grauniad The Independent Google Maps Chester Weather Station GitHub

Search or jump to... Pull requests Issues Marketplace Explore

ajstewartlang / Affective-Theory-of-Mind-Inferences Watch 0 Star 1 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Insights Settings

Branch: master Affective-Theory-of-Mind-Inferences / Expt 1 Script and data / Create new file Upload files Find file History

		Latest commit dc6b924 on 4 Apr 2018
..		
AngerFearAcc.csv	Tidied up code with consistent labelling	2 years ago
AngerFearRT.csv	Tidied up code with consistent labelling	2 years ago
AngerFeargraphacc.csv	R code and data	2 years ago
AngerFeargraphdata.csv	R code and data	2 years ago
Expt1_AngerFear_script.R	final changes	a year ago
Expt1_AngerFear_script_final.R	final changes	a year ago
Expt1_AngerFear_script_with_exclusions_final.R	final changes	a year ago
Expt1_AngerFear_script_with_item_exclusions_fin...	final changes	a year ago

This repository Search Pull requests Issues Gist + ⚙

ajstewartlang / Affective-Theory-of-Mind-Inferences Watch 0 Star 0 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Wiki Pulse Graphs Settings

Branch: master Find file Copy path

Affective-Theory-of-Mind-Inferences / Expt 1 Script and data / Expt1_AngerFear_script.R

ajstewartlang Tidied up code with consistent labelling de7b584 3 hours ago

0 contributors

59 lines (46 sloc) | 3.17 KB Raw Blame History

```
1 library(lme4)
2 library(lmerTest)
3 library(lsmeans)
4 library(pbkrtest)
5 library(readr)
6 library(ggplot2)
7
8 #script for AngerFear RT and accuracy data analysis with arousal
9
10 #this is the analysis of the RT data
11 AngerFearRT <- read_csv("~/AngerFearRT.csv")
12
13 AngerFearRT$StoryEmotion <- as.factor(AngerFearRT$StoryEmotion)
14 AngerFearRT$FaceExpression <- as.factor(AngerFearRT$FaceExpression)
15
16 contrasts(AngerFearRT$StoryEmotion) <- matrix(c(.5, -.5))
17 contrasts(AngerFearRT$FaceExpression) <- matrix(c(.5, -.5))
18
19 #with Subject, Vignette, and Face as crossed random effects with arousal
20 #full model does not converge so need to drop interaction term from the random effects - in addition, Face random effect has only random in
21 modelRTAr1 <- lmer(RT ~ StoryEmotion*FaceExpression*Arousal + (1+StoryEmotion+FaceExpression|Subject) + (1+StoryEmotion+FaceExpression|Vig
22 summary(modelRTAr1)
23 modelRT <- lmer(RT ~ StoryEmotion*FaceExpression + (1+StoryEmotion+FaceExpression|Subject) + (1+StoryEmotion+FaceExpression|Vignette) + (1
24 anova(modelRTAr1, modelRT)
25
26 #difference between models not signif - arousal does not interact with effect so drop arousal from subsequent analysis
27
28 modelRTnull <- lmer(RT ~ (1+StoryEmotion+FaceExpression|Subject) + (1+StoryEmotion+FaceExpression|Vignette) + (1+FaceExpression|Face),
29 anova(modelRT, modelRTnull)
30 summary(modelRT)
```

This is GitHub Desktop

The screenshot shows the GitHub Desktop application interface. At the top, there's a dark header bar with three colored window control buttons (red, yellow, green) on the left. To the right of these are three status indicators: "Current Repository" (Keele_Sept_2019), "Current Branch" (master), and "Fetch origin" (Last fetched just now). Below the header is a sidebar on the left containing a "Recent" section with a list of repositories. The "Keele_Sept_2019" repository is currently selected, highlighted with a light gray background. Other recent repositories listed include MRes_Advanced_Data_Skills, ajstewartlang.github.io, website, ajstewartlang, Binder_demo, Binder_RUM, CarpentryConnect2019, Change_detection_Danish, Chester_text_mining_R, Dockerfile_demo, holepunch, and SIPS_visualisation_1, SIPS_visualisation_2. To the right of the sidebar, the main area displays the message "No local changes". It includes a sub-message: "There are no uncommitted changes in this repository. Here are some friendly suggestions for what to do next." Below this are two call-to-action boxes: one for viewing files in Finder and another for viewing the repository page on GitHub. A small blue icon of a person working at a desk is positioned to the right of the main message area.

Current Repository
Keele_Sept_2019

Filter Add ▾

Recent

- MRes_Advanced_Data_Skills
- ajstewartlang.github.io
- website

ajstewartlang

- ajstewartlang.github.io
- Binder_demo
- Binder_RUM
- CarpentryConnect2019
- Change_detection_Danish
- Chester_text_mining_R
- Dockerfile_demo
- holepunch

Keele_Sept_2019

- MRes_Advanced_Data_Skills
- Psychology_MRes_Stats_R_...
- SIPS_presentation
- SIPS_visualisation_1
- SIPS_visualisation_2

No local changes

There are no uncommitted changes in this repository. Here are some friendly suggestions for what to do next.

View the files of your repository in Finder
Repository menu or ⌘ ⌘ F Show in Finder

Open the repository page on GitHub in your browser
Repository menu or ⌘ ⌘ G View on GitHub

**Let's run
through a
script...**



TWENTIETH CENTURY-FOX Presents A LUCASFILM LTD. PRODUCTION **STAR WARS**
Starring **MARK HAMILL HARRISON FORD CARRIE FISHER**
PETER CUSHING
and
ALEC GUINNESS

Written and Directed by
GEORGE LUCAS

Produced by
GARY KURTZ JOHN WILLIAMS

PANAVISION® PRINTS BY DE LUXE® TECHNICOLOR®

Original Motion Picture Soundtrack on 20th Century Records and Tapes

Making Films Sound Better
DOLBY SYSTEM
Noise Reduction - High Fidelity

**STAR
WARS**





R packages for data science

The tidyverse is an opinionated **collection of R packages** designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

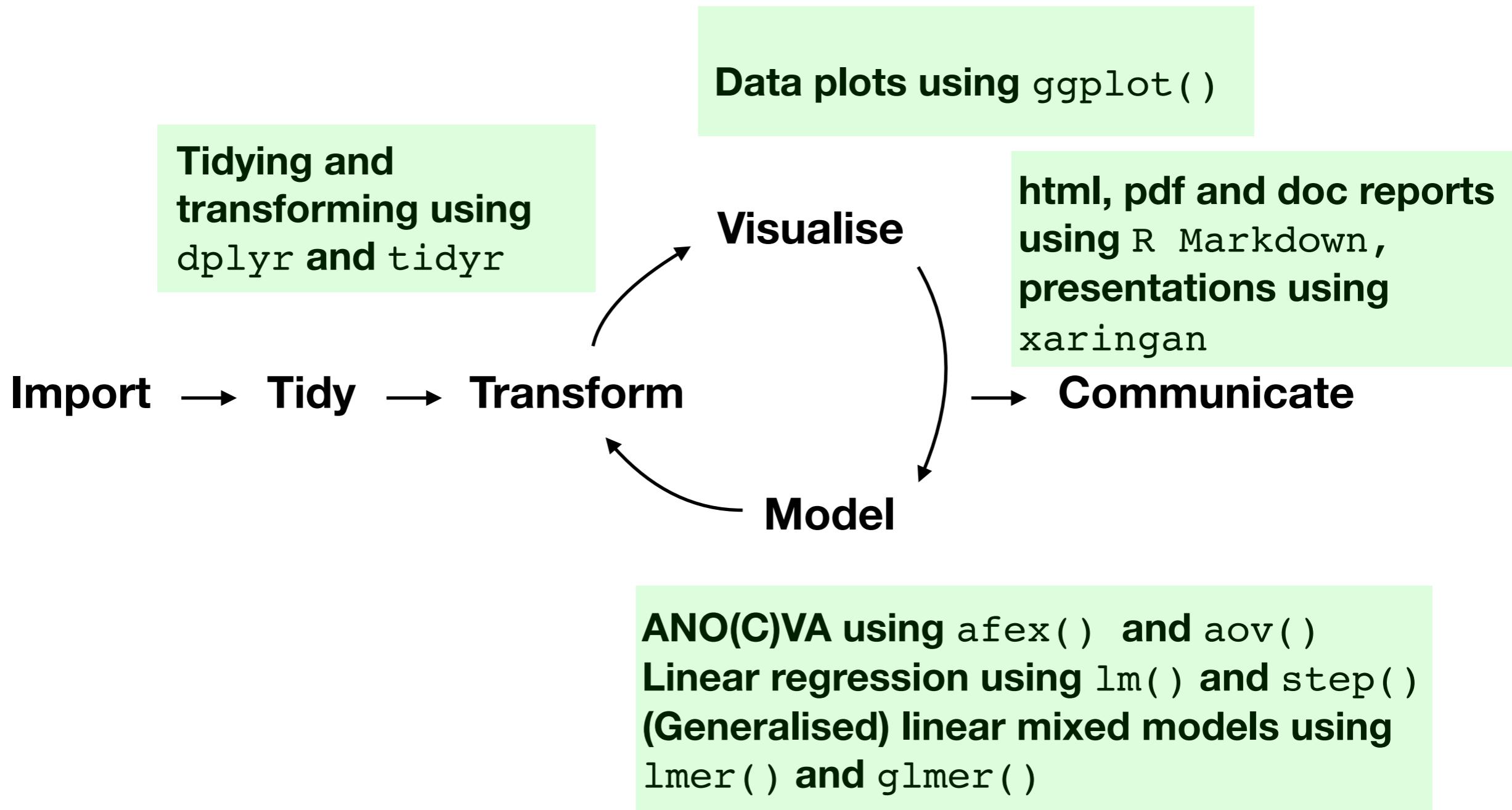
Install the complete tidyverse with:

```
install.packages("tidyverse")
```

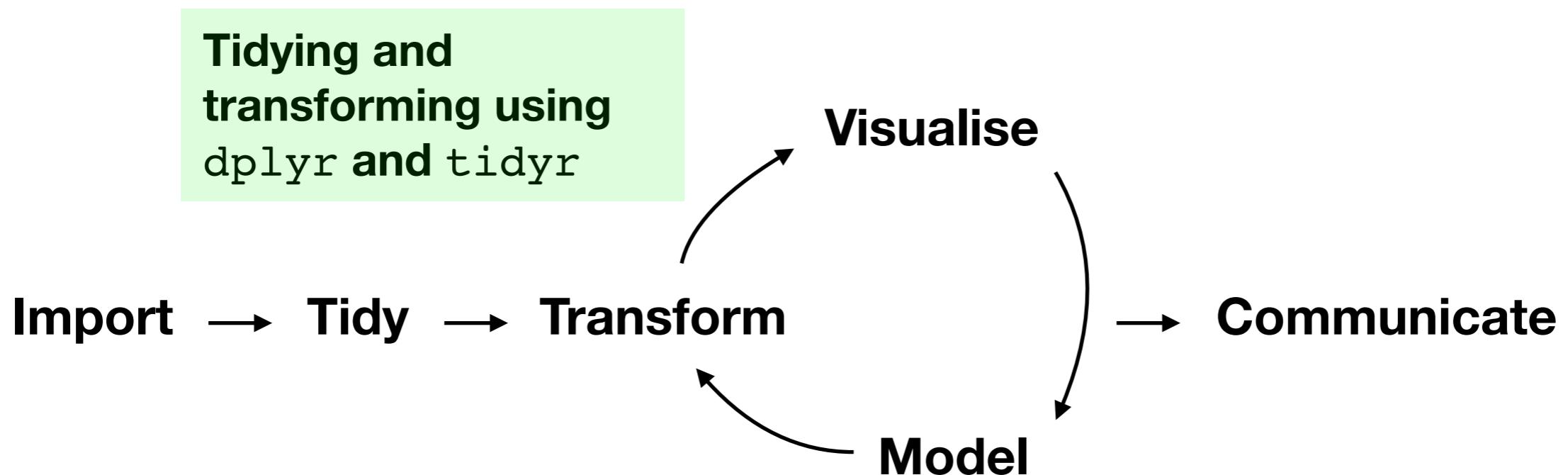
Tidyverse packages

- The Tidyverse contains a number of packages, all containing functions that are designed to ‘play well’ with each other. Packages include `ggplot2`, `dplyr`, and `tidyverse`.
- You can load each package separately with (e.g.)
 - > `library(ggplot2)`
- or load all tidyverse packages with
 - > `library(tidyverse)`

Workflow



Workflow



Let's get ready to code...

On your computer, fire up RStudio and install the tidyverse:

```
> install.packages("tidyverse")
```

Then create a new Project in a new folder, and fire up a new script...

On the first line of your script type:

```
library(tidyverse)
```

Then run the script...

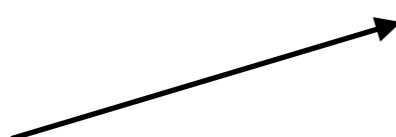
Let's now load two datasets that I've created:

```
data1 <- read_csv("https://bit.ly/31Te6HQ")  
dataRT <- read_csv("https://bit.ly/2Zf1WOr")
```

Tidying and Transforming Data

Imagine we have two rectangular datasets - one (called data1) contains a large number of records of individual participants with measures of Working Memory, IQ, and reading comprehension.

If we type into the Console:
> data1
the data frame is displayed like
this...



	id	wm	iq	comp
	<int>	<int>	<int>	<int>
1	1	43	72	16
2	2	51	109	18
3	3	55	107	18
4	4	38	102	20
5	5	52	121	17
6	6	52	92	16
7	7	47	68	21
8	8	47	97	23
9	9	47	93	22
10	10	45	101	17
# ... with 9,990 more rows				

To get more information about the structure of our data frame we can type:

```
> str(data1)
Classes 'tbl_df', 'tbl' and 'data.frame': 10000 obs. of 4 variables:
 $ id   : int  1 2 3 4 5 6 7 8 9 10 ...
 $ wm   : int  43 51 55 38 52 52 47 47 47 45 ...
 $ iq   : int  72 109 107 102 121 92 68 97 93 101 ...
 $ comp: int  16 18 18 20 17 16 21 23 22 17 ...
```

So we have 10,000 observations with 4 variables associated with each observation - all of them of type integer.

If you ever need help about a function (e.g. str), just type:

```
>?str
```

or

```
>help(str)
```

in the Console window.

Imagine that 48 of these 10,000 people also took part in a reading time experiment and we have their reading data (called `dataRT`) for Simple Sentence and Complex Sentence reading conditions:

```
> str(dataRT)
Classes 'tbl_df', 'tbl' and 'data.frame': 48 obs. of 3 variables:
 $ id           : int  6400 457 8291 4998 2579 9122 1138 5138 5244 3160 ...
 $ simple_sentence : int  1902 1797 2080 1856 1997 1868 2154 1933 1900 1929 ...
 $ complex_sentence: int  2341 2503 2731 2375 2177 2284 2441 2349 2371 2372 ...
```

We are maybe interested in analysing the data of these 48 people in the data frame called `dataRT` but covarying out the effect of IQ captured in our data frame called `data1`.

Problem - how can we combine these two data frames so that we end up with one data frame of 48 people, their reading times plus their individual difference measures?

Manually, in Excel we could open the two data frames as spreadsheets and cut and paste cases where the id number matches...

Probably ok for 48 participants, but what if you had 200 or 2,000?

In R, we can use the `inner_join` function from the `dplyr` package where we join the two data frames matched by ID.

```
> dataRT_all <- inner_join(data1, dataRT, by = (c("id")) )  
> dataRT_all  
    id   wm   iq  comp simple_sentence complex_sentence  
1   95   47   94    19                 2154                  2441  
2  400   45  118    18                 1824                  2456  
3  457   42  100    22                 1857                  2324  
4 1138   41   77    18                 1902                  2341  
5 1587   54   67    21                 1844                  2320  
6 1805   52  109    19                 2224                  2256  
7 1864   57  111    19                 1880                  2391  
8 2006   44  110    19                 2091                  2456  
9 2183   55  125    23                 1926                  2218  
10 2318  51   91    21                 1960                  2440
```

We can use the assignment symbol `<-` to assign the output of this `inner_join` function to a new variable I'm calling `dataRT_all`. We can ask for the structure of this new data frame using the `str()` function:

```
> dataRT_all <- inner_join(data1, dataRT, by = (c("id")))
> str(dataRT_all)
Classes 'tbl_df', 'tbl' and 'data.frame': 48 obs. of 6 variables:
 $ id           : int  95 400 457 1138 1587 1805 1864 2006 2183 2318 ...
 $ wm           : int  47 45 42 41 54 52 57 44 55 51 ...
 $ iq            : int  94 118 100 77 67 109 111 110 125 91 ...
 $ comp          : int  19 18 22 18 21 19 19 19 23 21 ...
 $ simple_sentence: int  2154 1824 1857 1902 1844 2224 1880 2091 1926 1960 ...
 $ complex_sentence: int  2441 2456 2324 2341 2320 2256 2391 2456 2218 2440 ...
```

So we have created a new data frame of 48 participants consisting of their reading times and their individual difference measures from two separate (and different sized) data frames...with one line of code...

Now imagine we find the distributions of reading times for our two conditions are positively skewed (and we discover the residuals are non-normal). We could log transform these two columns and have two new columns in our data frame - let's call them `log_simple` and `log_complex`. We can use the `mutate` function in the `dplyr` package to create two new columns.

```
> data_transformed <- mutate(dataRT_all, log_simple = log(simple_sentence) ,  
+ log_complex = log(complex_sentence))  
> data_transformed  
# A tibble: 48 x 8  
  id      wm      iq      comp simple_sentence complex_sentence log_simple log_complex  
  <int> <int> <int> <int>          <int>          <int>       <dbl>       <dbl>  
1 95     47     94     19          1960         2440    7.58     7.80  
2 400    45    118     18          2186         2200    7.69     7.70  
3 457    42    100     22          1797         2503    7.49     7.83  
4 1138   41     77     18          2154         2441    7.68     7.80  
5 1587   54     67     21          1936         2395    7.57     7.78  
6 1805   52    109     19          1864         2560    7.53     7.85  
7 1864   57    111     19          1930         2540    7.57     7.84  
8 2006   44    110     19          2230         2267    7.71     7.73  
9 2183   55    125     23          1857         2324    7.53     7.75  
10 2318   51     91     21          1918         2739    7.56     7.92  
# ... with 38 more rows
```

Perhaps we have a reason to exclude a particular participant - number 2006 for example. We can use the filter function in `dplyr` to keep those participants where the ID number does not equal 2006.

```
filtered_data <- filter(data_transformed, id != 2006)
```

`!=` stands for “not equal to”- here are other useful logical operators in R:

`<` less than

`<=` less than or equal to

`>` greater than

`>=` greater than or equal to

`==` exactly equal to

`!=` not equal to

We can now apply our logical vector to our dataRT_all data frame and create a new filtered data frame (which I am calling filtered_data):

```
> filtered_data <- filter(data_transformed, id != 2006)
> filtered_data
# A tibble: 47 x 8
  id    wm    iq   comp simple_sentence complex_sentence log_simple log_complex
  <int> <int> <int> <int>          <int>          <int>      <dbl>       <dbl>
1 95     47    94    19        1960        2440      7.58       7.80
2 400    45   118    18        2186        2200      7.69       7.70
3 457    42   100    22        1797        2503      7.49       7.83
4 1138   41    77    18        2154        2441      7.68       7.80
5 1587   54    67    21        1936        2395      7.57       7.78
6 1805   52   109    19        1864        2560      7.53       7.85
7 1864   57   111    19        1930        2540      7.57       7.84
8 2183   55   125    23        1857        2324      7.53       7.75
9 2318   51    91    21        1918        2739      7.56       7.92
10 2324  43   120    20        1891        2426      7.54       7.79
# ... with 37 more rows
```

We could then run an ANCOVA over the log transformed RTs while covarying out the individual participant effects...

Problem - imagine our data are in the wrong ‘shape’ - they are in **Wide format** (each row is one *participant*) but we need them in **Long** or **Tidy format** (each row is one *observation*).

In SPSS, most data will be in Wide format with each experimental condition its own column:

```
> dataRT
# A tibble: 48 x 3
      id simple_sentence complex_sentence
      <int>           <int>           <int>
1    6400            1902            2341
2     457            1797            2503
3   8291            2080            2731
4   4998            1856            2375
5   2579            1997            2177
6   9122            1868            2284
7   1138            2154            2441
8   5138            1933            2349
9   5244            1900            2371
10  3160            1929            2372
# ... with 38 more rows
```

For many analyses in R, data need to be in Long format with each row being one observation. So, we want to transform our dataRT data frame so it looks like this:

id	condition	rt
...

To do this we can use the `gather()` function in the `tidyverse` package.

```
> data_long <- gather(dataRT, "condition", "rt", c("simple_sentence",  
"complex_sentence"))
```

The first parameter is the name of the data frame we want to reshape, the second is the name of the new ‘Key’ column, the third is the name of the new value column and the fourth the names of the columns we want to collapse.

We can use this to create a new data frame called `data_long` which looks like this:

```
> data_long <- gather(dataRT, "condition", "rt", c("simple_sentence",
  "complex_sentence"))
> data_long
# A tibble: 96 x 3
  id condition       rt
  <int> <chr>     <int>
1 6400 simple_sentence 1902
2 457  simple_sentence 1797
3 8291 simple_sentence 2080
4 4998 simple_sentence 1856
5 2579 simple_sentence 1997
6 9122 simple_sentence 1868
7 1138 simple_sentence 2154
8 5138 simple_sentence 1933
9 5244 simple_sentence 1900
10 3160 simple_sentence 1929
# ... with 86 more rows
```

And in reverse we can use the `spread()` function to go from Long to Wide data format:

```
> data_wide <- spread(data_long, "condition", "rt")
> data_wide
# A tibble: 48 x 3
  id complex_sentence simple_sentence
  <dbl> <dbl> <dbl>
1 95     2441    2154
2 400    2456    1824
3 457    2324    1857
4 1138   2341    1902
5 1587   2320    1844
6 1805   2256    2224
7 1864   2391    1880
8 2006   2456    2091
9 2183   2218    1926
10 2318  2440    1960
# ... with 38 more rows
```

We're now back to where we started with data in Wide format:

This is just a small example of functions in the `dplyr` and `tidyverse` packages that allow you to tidy, transform, and reshape your data. All of your code for doing this should appear at the start of your analysis script so that others (and you in 5 years or 5 days time) can see exactly what you did.

This allows for fully reproducible data preparation in the first part of your analysis workflow (important for Open Science and reproducibility).

Generating Descriptives - using dplyr

- You can use the `group_by()` and `summarise()` functions in the `dplyr` package to generate descriptives.
- In the following example, we are also using the pipe operator `%>%` which passes a value into an expression or function call from left to right:

```
> data_long %>%
  group_by(condition) %>%
  summarise(Mean = mean(rt), Min = min(rt), Max = max(rt), SD = sd(rt))
# A tibble: 2 x 5
  condition      Mean   Min   Max     SD
  <chr>        <dbl> <int> <int>  <dbl>
1 complex_sentence 2405.  2177  2739  132.
2 simple_sentence 1957.  1694  2356  147.
```

Tidying Up Some Real World Messy Data

Let's load another dataset that I created:

```
my_data <- read_csv("https://bit.ly/2KPZEE9")
```

Tidying Up Some Real World Messy Data

- We ran a reaction time experiment with 24 participants and 4 conditions - they are numbered 1-4 in our datafile.

```
> my_data
# A tibble: 96 x 3
  participant condition     rt
        <int>      <int> <int>
  1             1          1    879
  2             1          2   1027
  3             1          3   1108
  4             1          4    765
  5             2          1   1042
  6             2          2   1050
  7             2          3    942
  8             2          4    945
  9             3          1    943
 10            3          2    910
# ... with 86 more rows
```

- But actually it was a repeated measures design where we had one factor (Prime Type) with two levels (A vs. B) and a second factor (Target Type) with two levels (A vs. B)
- We want to recode our data frame so it better matches our experimental design.
- First we need to recode our 4 conditions like this:

```
# Recode condition columns follows:  
# Condition 1 = prime A, target A  
# Condition 2 = prime A, target B  
# Condition 3 = prime B, target A  
# Condition 4 = prime B, target B  
  
my_data <- my_data %>%  
  mutate(condition = recode(condition,  
    "1" = "primeA_targetA",  
    "2" = "primeA_targetB",  
    "3" = "primeB_targetA",  
    "4" = "primeB_targetB"))
```

- Now our data frame looks like this:

```
> my_data
# A tibble: 96 x 3
  participant condition       rt
  <int> <chr>      <int>
1          1 primeA_targetA    879
2          1 primeA_targetB   1027
3          1 primeB_targetA   1108
4          1 primeB_targetB    765
5          2 primeA_targetA  1042
6          2 primeA_targetB  1050
7          2 primeB_targetA   942
8          2 primeB_targetB   945
9          3 primeA_targetA   943
10         3 primeA_targetB   910
# ... with 86 more rows
```

- We then need to separate out our Condition column into two - one for our first factor (Prime), and one for our second factor (Target).

```
> my_data <- separate(my_data, col = "condition", into = c("prime",
  "target"), sep = "_")
> my_data
# A tibble: 96 x 4
  participant prime  target       rt
  <int> <chr>  <chr>     <int>
1          1 primeA targetA    879
2          1 primeA targetB   1027
3          1 primeB targetA   1108
4          1 primeB targetB    765
5          2 primeA targetA  1042
6          2 primeA targetB  1050
7          2 primeB targetA   942
8          2 primeB targetB   945
9          3 primeA targetA   943
10         3 primeA targetB   910
# ... with 86 more rows
```

- This is looking good - we now have our two factors coded separately and our data are in tidy format (i.e., one observation per row).
- How long would this have taken you in Excel? Would it have been reproducible?

- Perhaps we want to go from the data in long format, to wide format.

```
> my_data <- unite(my_data, col = "condition", c("prime", "target"), sep = "_")
> wide_data <- spread(my_data, key = "condition", value = "rt")
> wide_data
# A tibble: 24 x 5
  participant primeA_targetA primeA_targetB primeB_targetA primeB_targetB
      <int>        <int>        <int>        <int>        <int>
1         1          879        1027        1108        765
2         2         1042        1050        942         945
3         3          943        910         952        900
4         4          922        1006        1095        988
5         5          948        908         916       1241
6         6         1013        950         955       1045
7         7          930        855        1057        897
8         8          998        906        1110        952
9         9          929        949         837        883
10        10          781        865         970        953
# ... with 14 more rows
```

- No matter what format your data are in originally, you can use functions from the `dplyr` and `tidyverse` packages to quickly get it into whatever format you need for analysis.

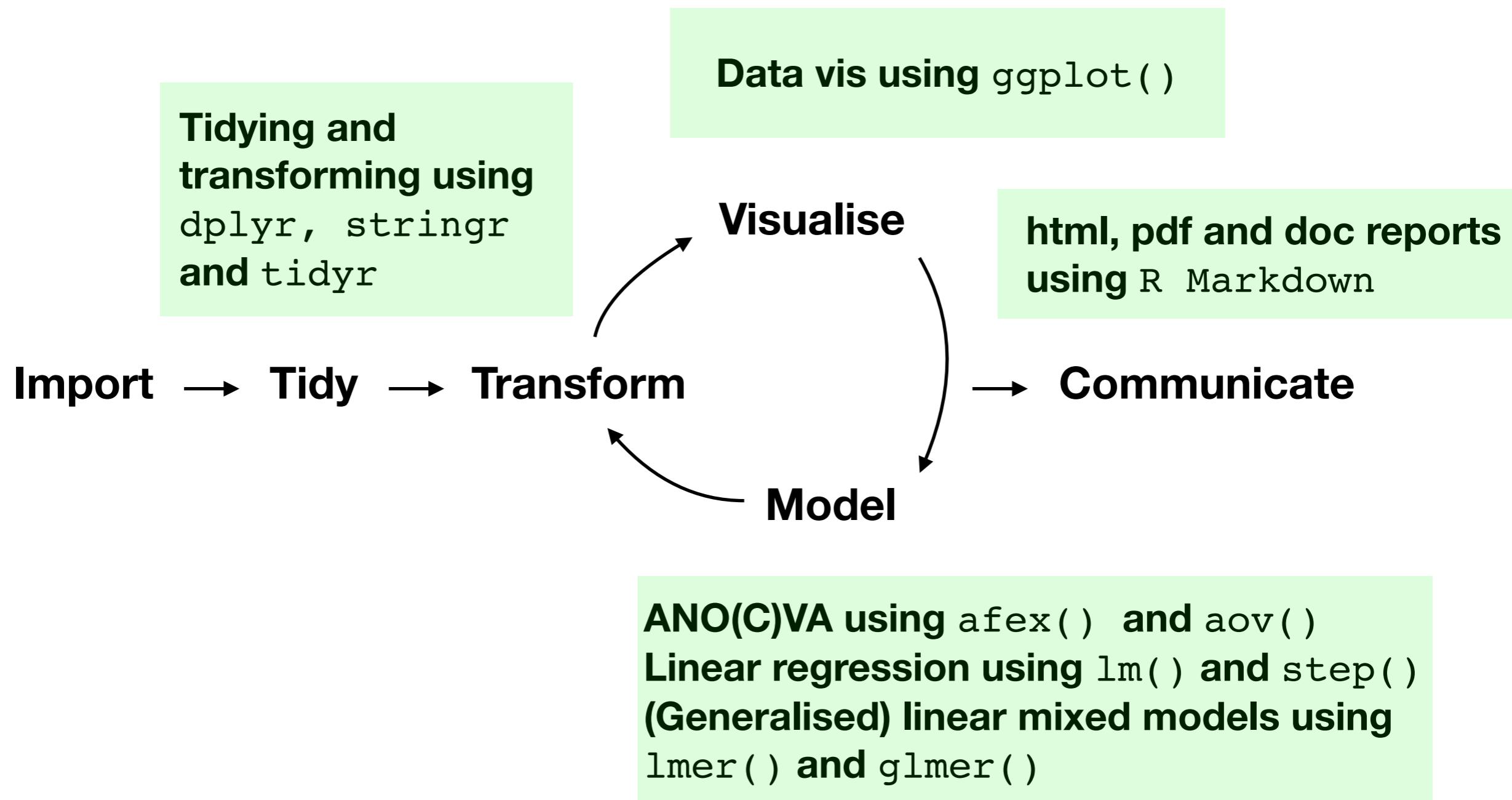
Or using the pipe...

- We could alternatively have used the `%>%` operator to combine all the last few operations which would have avoided the need to create temporary variables.

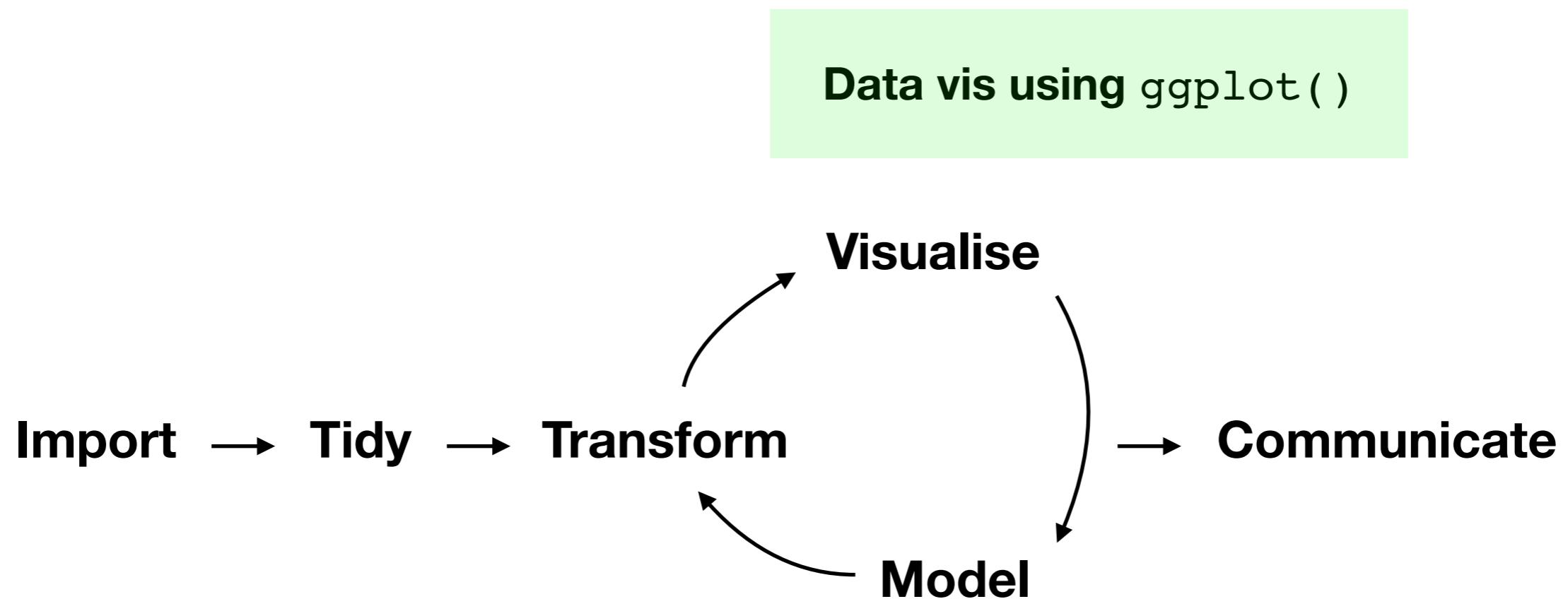
```
my_data %>%
  unite(col = "condition", c("prime", "target"), sep = "_") %>%
  spread(key = "condition", value = "rt")
```

- Take `my_data` and then `unite` and then `spread`...

A Reproducible Workflow



A Reproducible Workflow

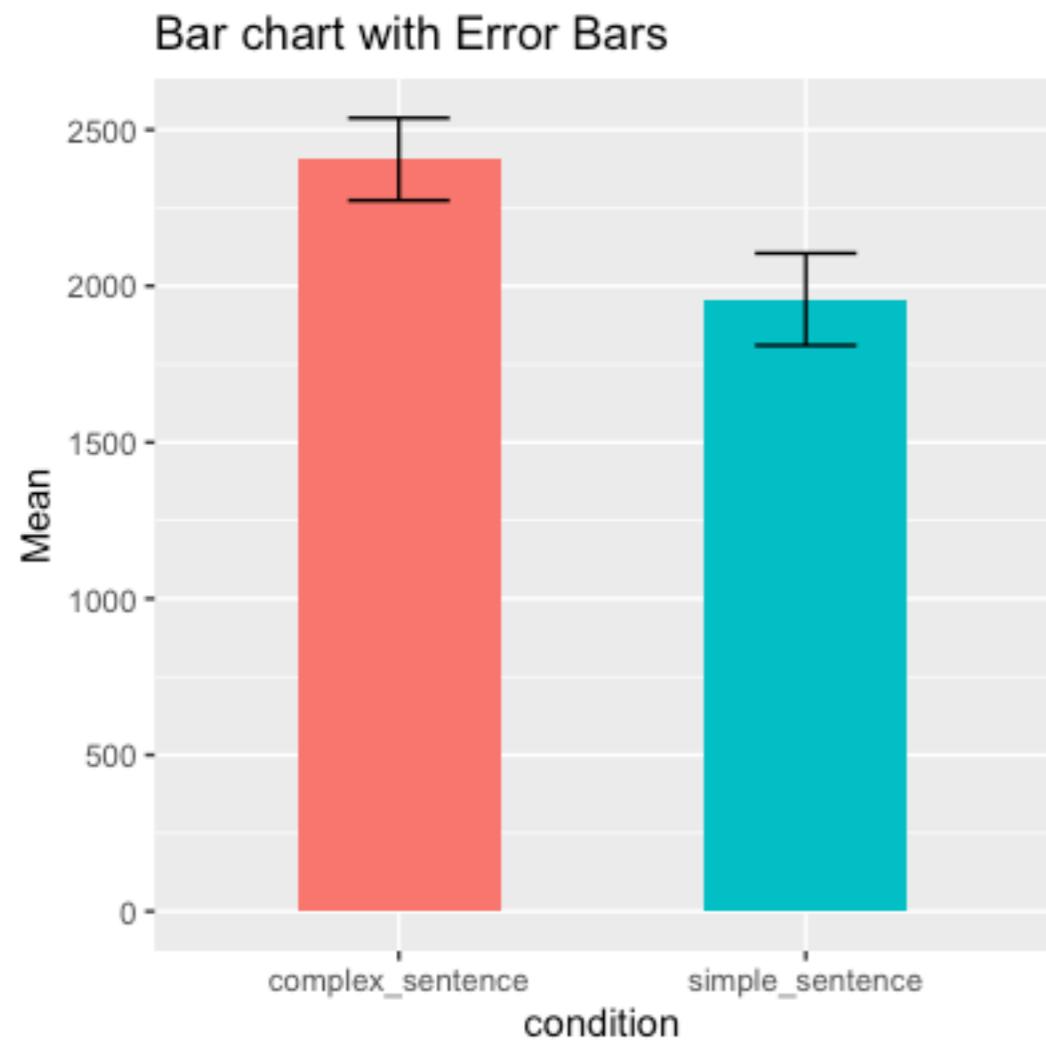


Visualising Your Data

- R has a number of in built (base) graphics functions, but you're more likely to use functions from within the `ggplot2` package. `ggplot2` is part of the `tidyverse` so if you have used `library(tidyverse)` then `ggplot2` will already be loaded.

```
> library(ggplot2)
```

Bar Graphs (bleurgh!)

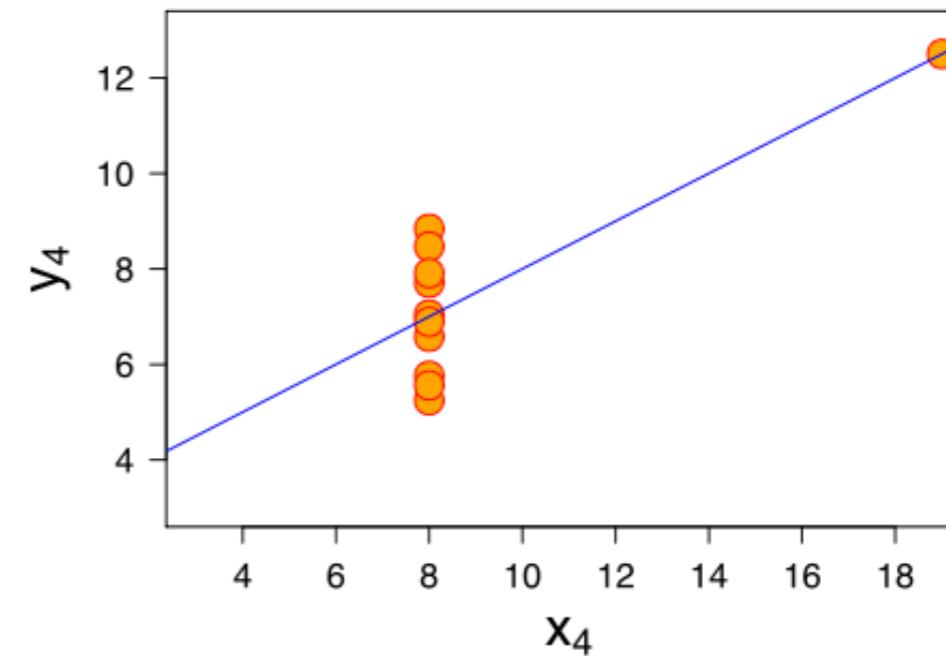
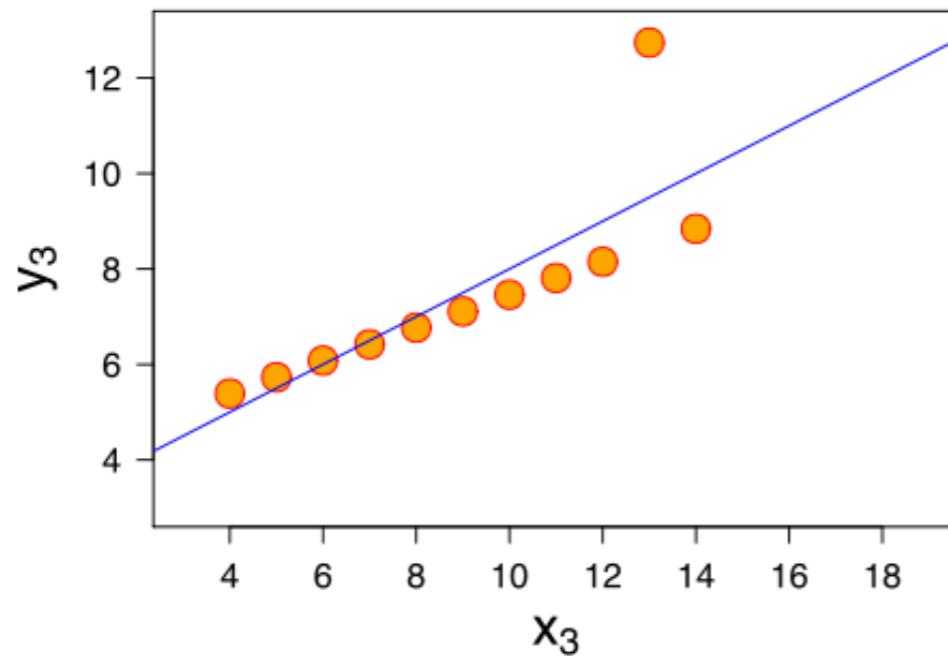
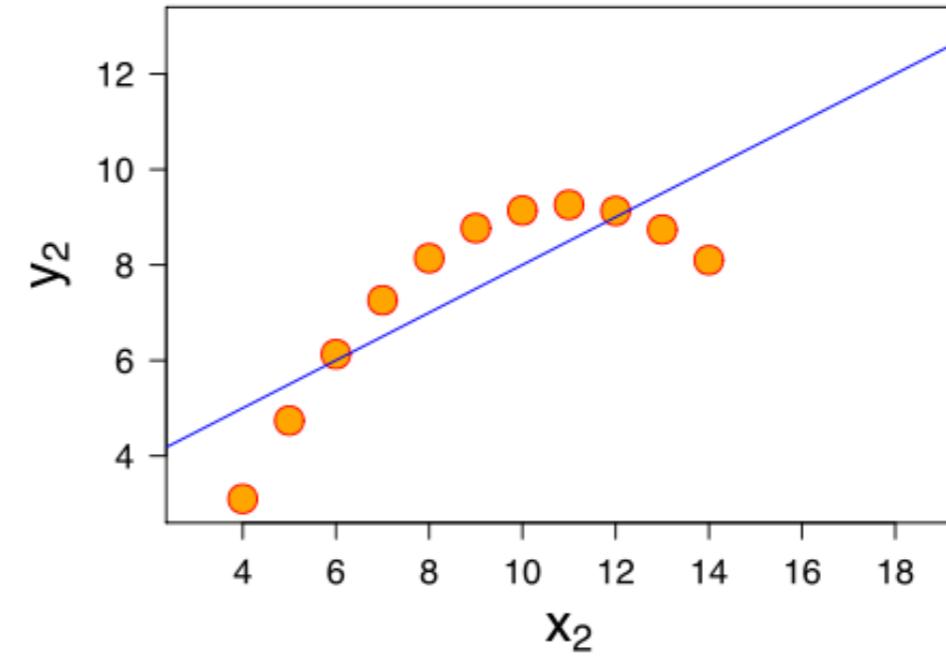
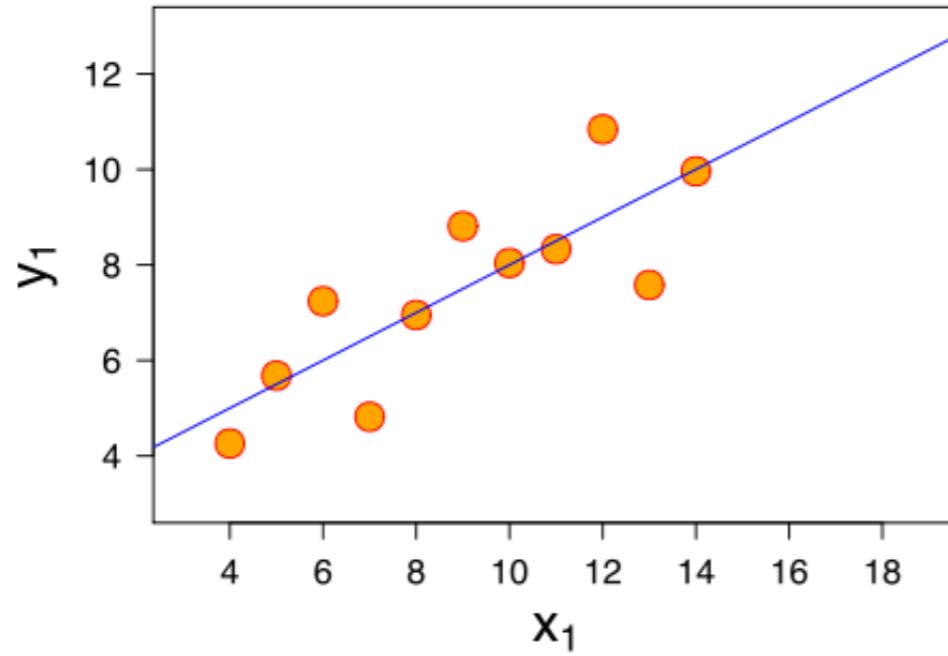


Bar graphs tend to be quite limited in terms of what they communicate. Here they communicate the means for levels of a factor and information about variance. But they don't tell us anything about the *distribution* of the data.

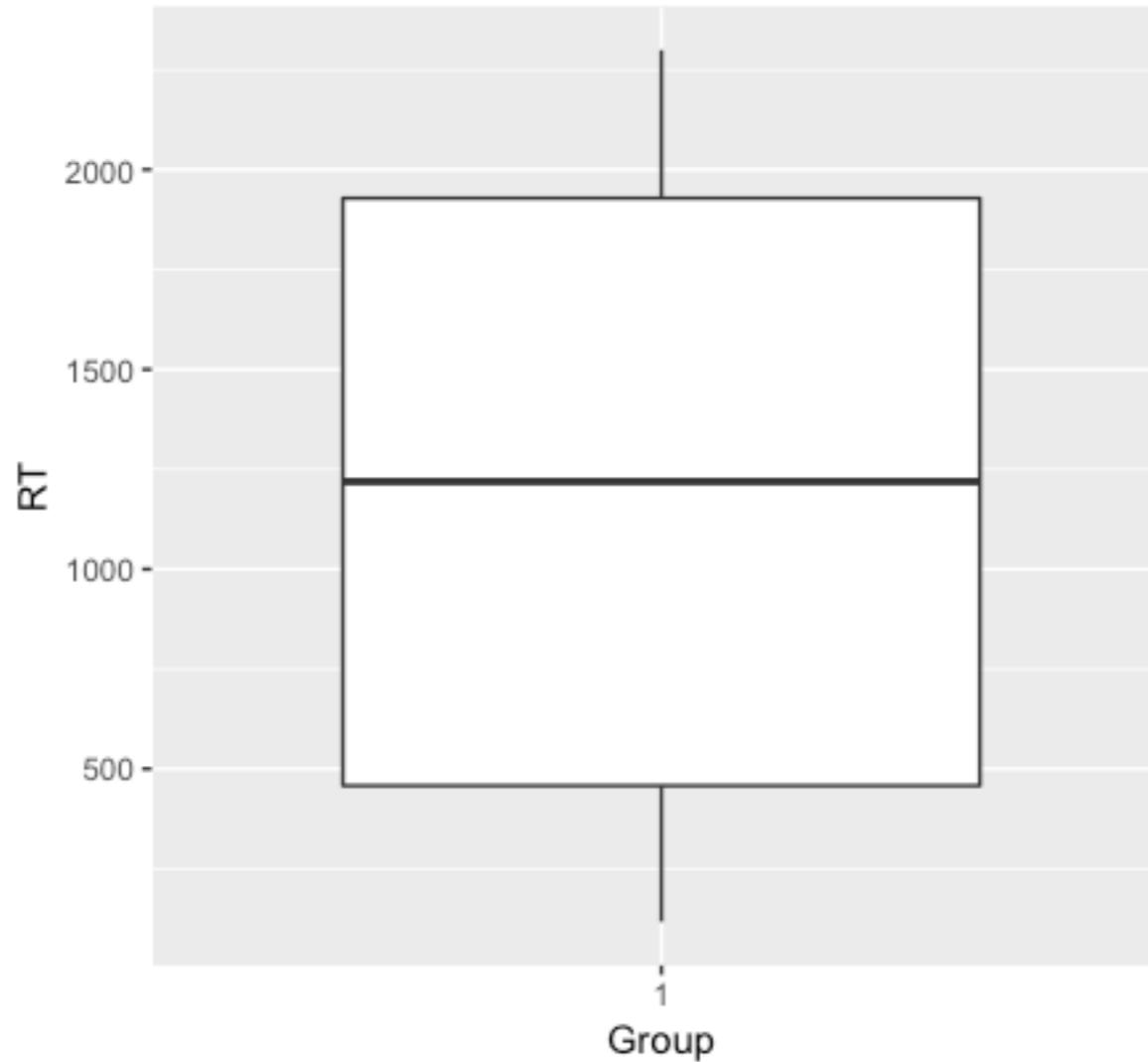
```
data_summ <- data_long %>% group_by(condition) %>% summarise(Mean = mean(rt), sd = sd(rt))

ggplot(data_summ, aes(x = condition, y = Mean, group = condition,
                      fill = condition, ymin = Mean - sd, ymax = Mean + sd)) +
  geom_bar(stat = "identity", width = .5) +
  geom_errorbar(width = .25) +
  ggtitle("Bar chart with Error Bars") +
  guides(fill = FALSE)
```

Anscombe's Quartet

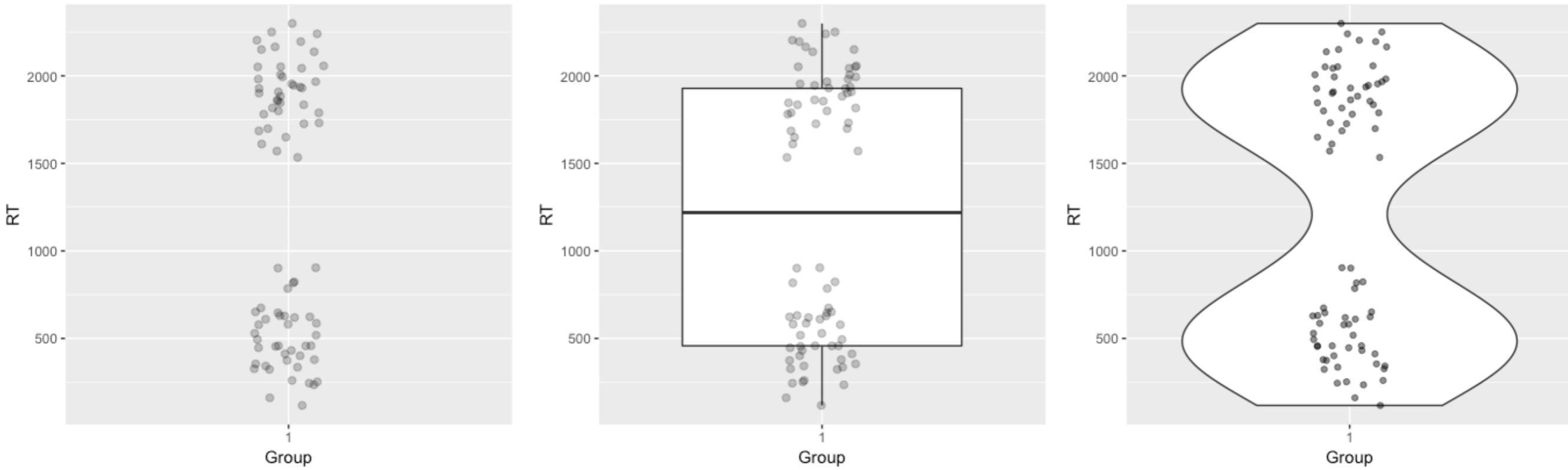


Plots Based on Aggregated Data Can Mislead...



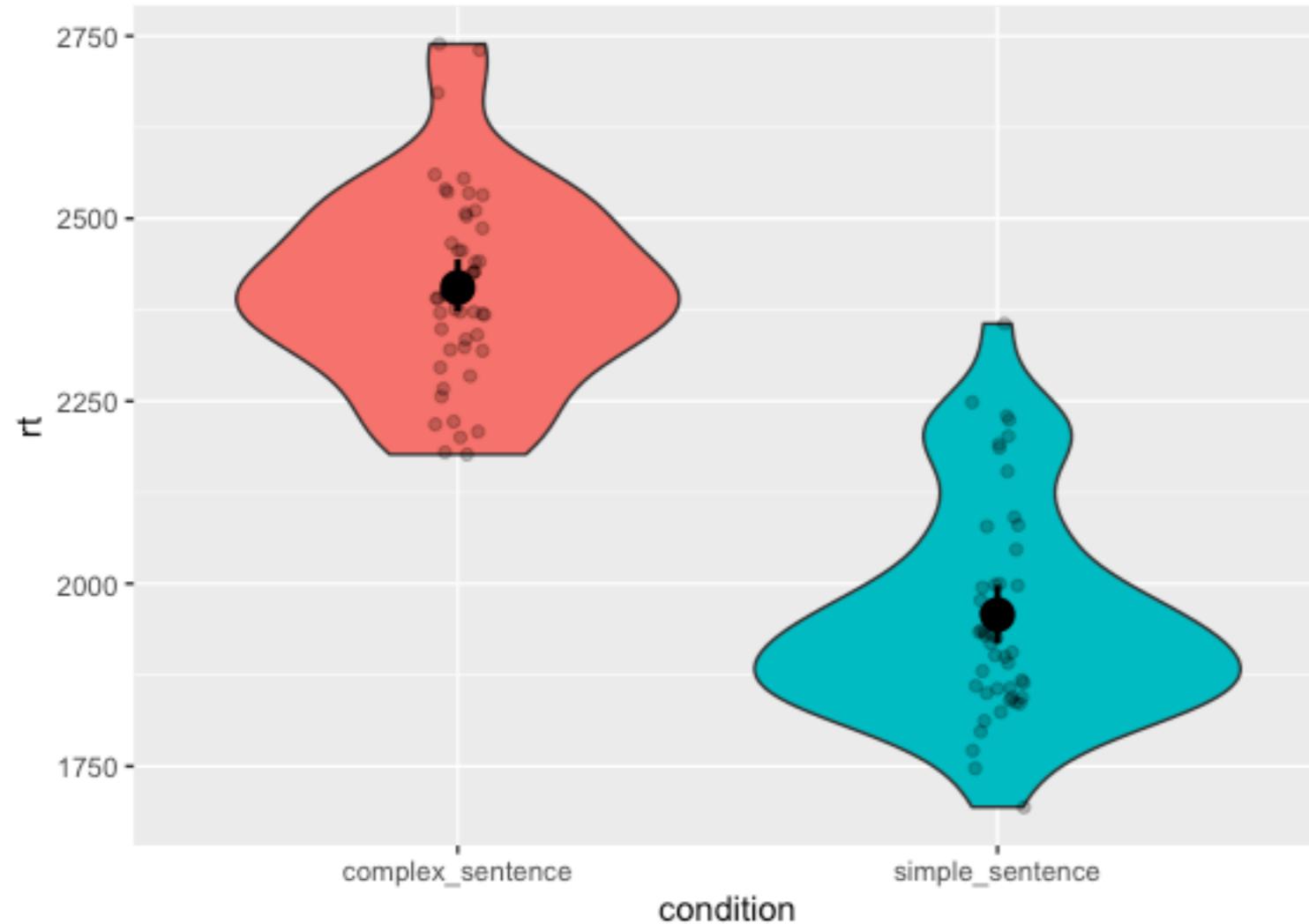
You might make one set of inferences based on this boxplot - maybe a median around 1,250 with the 25th and 75th percentiles being ~480 to ~1,980...

But look more closely at the actual data...



The data are clearly bimodal with no actual data point near the mean.
Distribution shape matters and we need to capture that in our data visualisations.

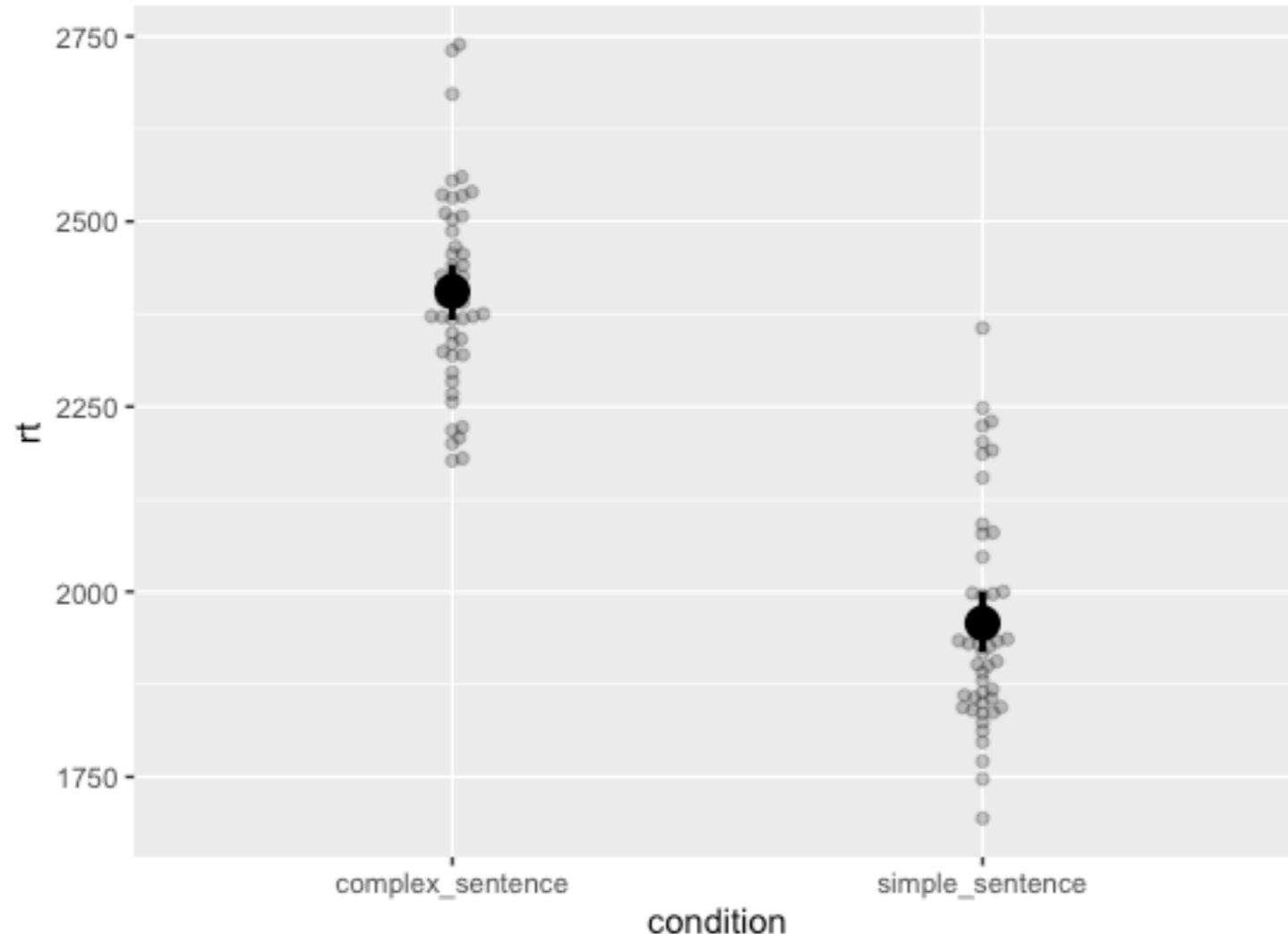
Violin Plots



Violin plots tell us about the distribution of the data. The width at any point corresponds to the *density* of the data at that value.

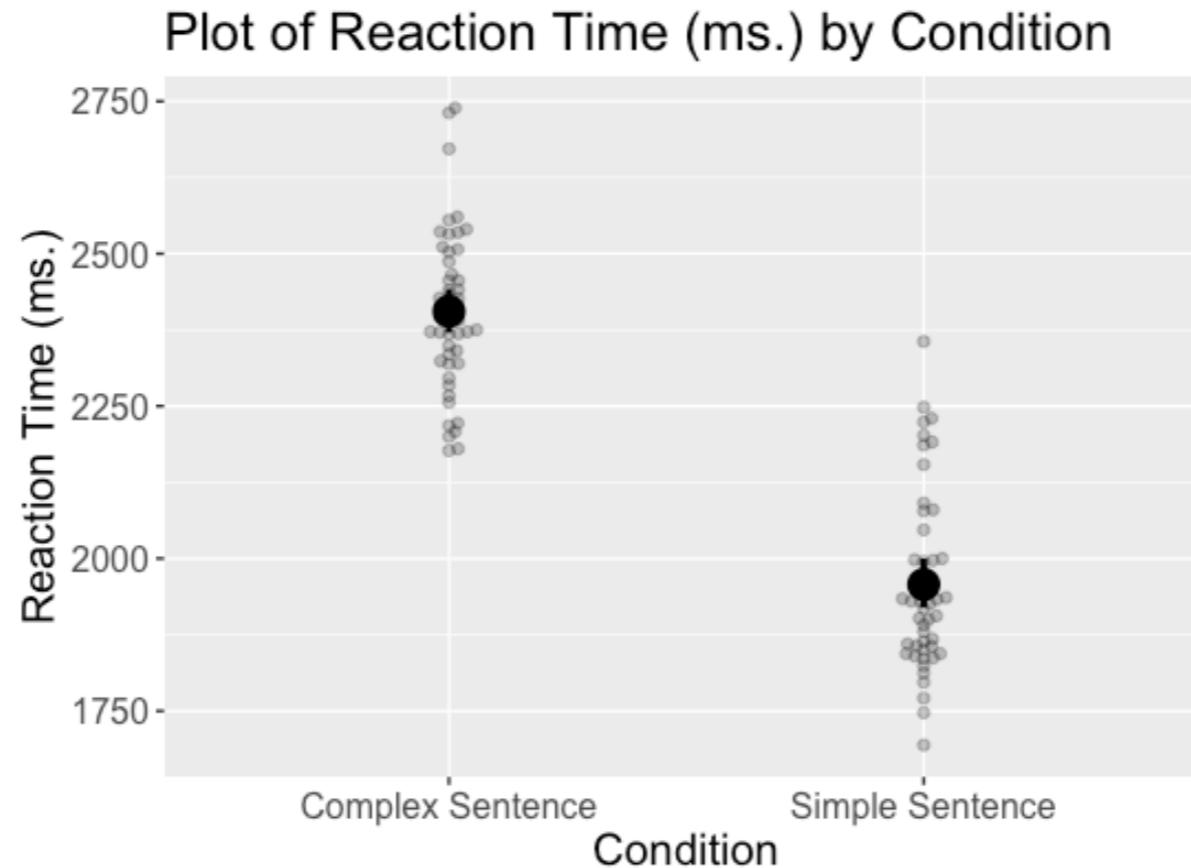
```
ggplot(data_long, aes(x = condition, y = rt,
  group = condition, fill = condition)) +
  geom_violin() +
  geom_jitter(alpha = .25, position = position_jitter(0.05)) +
  guides(colour = FALSE, fill = FALSE) +
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1)
```

Beeswarm Plots



```
ggplot(data_long, aes(x = condition, y = rt, group = condition, fill = condition)) +  
  geom_beeswarm(alpha = .25) +  
  guides(colour = FALSE, fill = FALSE) +  
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1)
```

Tidying up our labels



```
data_long %>%
  mutate(Condition = recode(condition,
    "complex_sentence" = "Complex Sentence",
    "simple_sentence" = "Simple Sentence")) %>%
  ggplot(aes(x = Condition, y = rt, group = Condition, fill = Condition)) +
  geom_beeswarm(alpha = .25) +
  guides(colour = FALSE, fill = FALSE) +
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1) +
  labs(title = "Plot of Reaction Time (ms.) by Condition",
       x = "Condition",
       y = "Reaction Time (ms.)") +
  theme(text = element_text(size = 15))
```

Themes

- The ggthemes package has lots of pre-built ggplot themes that we can apply to our ggplot visualisations.

```
>library(ggthemes)
```

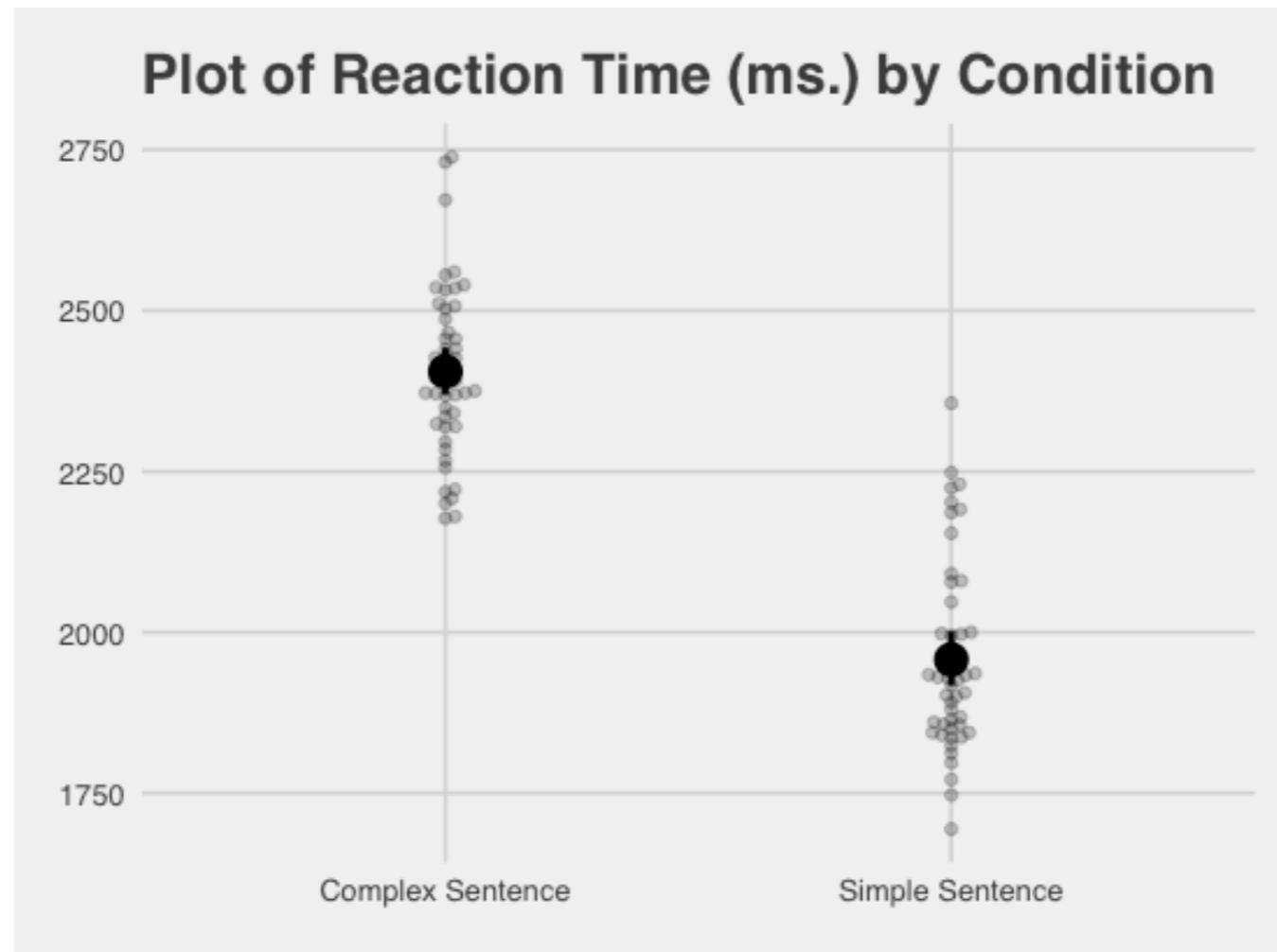
A screenshot of the RStudio interface showing the documentation for the ggthemes package. The code editor shows a snippet of R code:

```
+ "complex_sentence" = "Complex Sentence",
+ theme_excel          {ggthemes}
+ theme_excel_new      {ggthemes}
+ theme_few             {ggthemes}
+ theme_fivethirtyeight {ggthemes}  #<-- This is the currently selected theme
+ theme_foundation       {ggthemes}
+ theme_gdocs            {ggthemes}
+ theme_hc               {ggthemes}
```

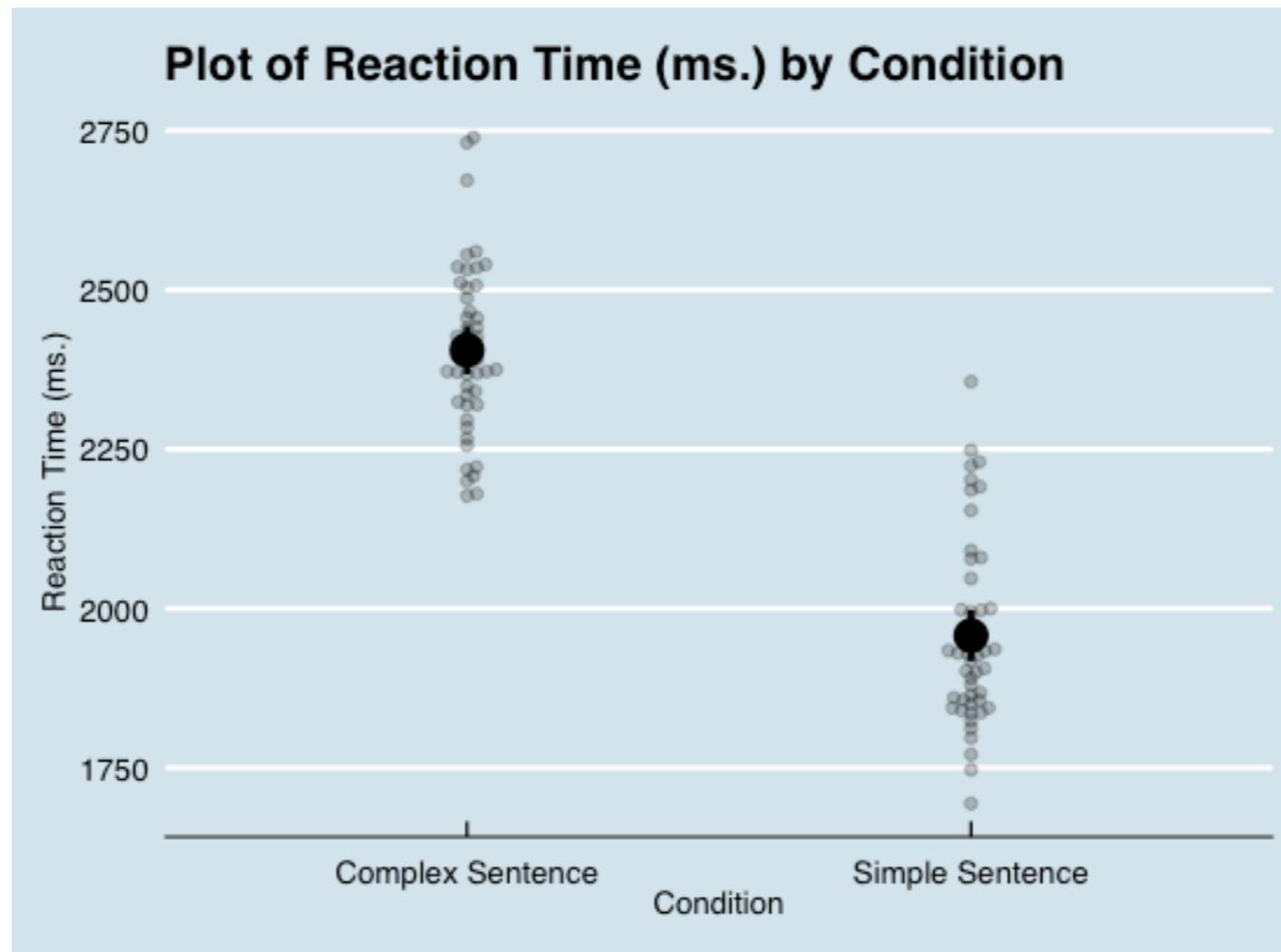
The `theme_fivethirtyeight` function is highlighted with a blue selection bar. A tooltip window is open over the function definition, containing the following text:

theme_fivethirtyeight(base_size = 12, base_family =
"sans")
Theme inspired by the plots on <http://fivethirtyeight.com>.
Press F1 for additional help

The RStudio interface includes a vertical scrollbar on the right side of the code editor.

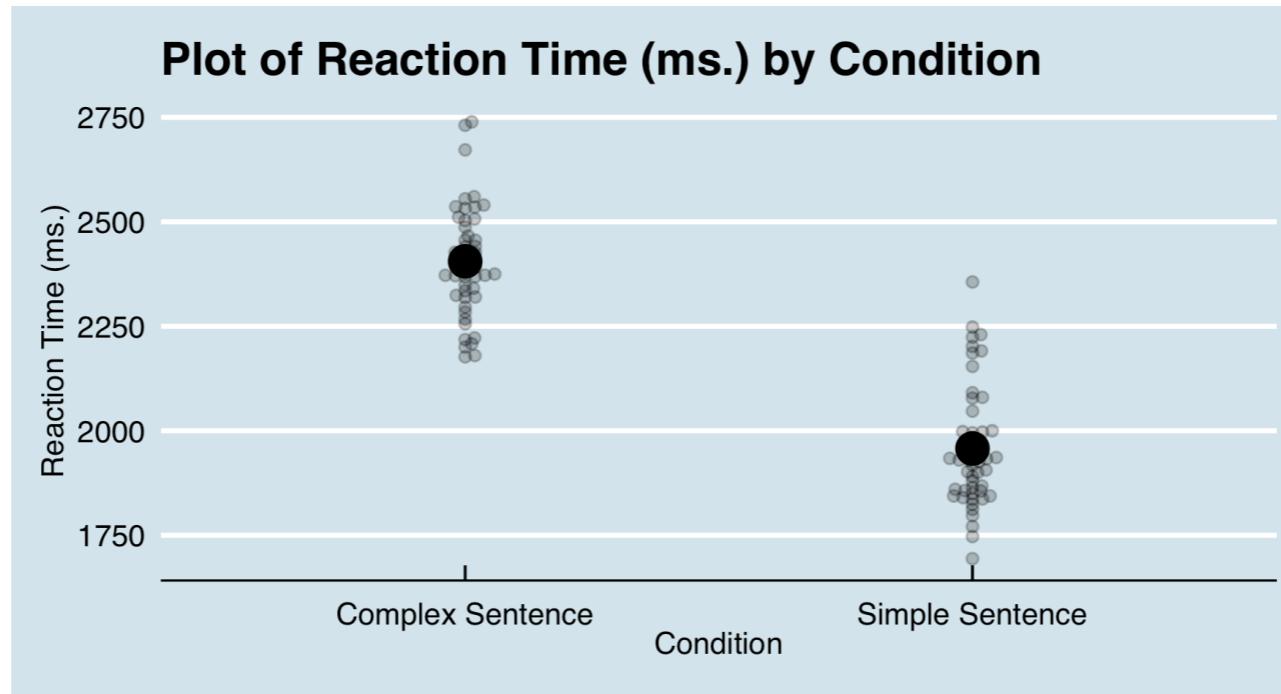


```
data_long %>%
  mutate(Condition = recode(condition,
    "complex_sentence" = "Complex Sentence",
    "simple_sentence" = "Simple Sentence")) %>%
  ggplot(aes(x = Condition, y = rt)) +
  geom_beeswarm(alpha = .25) +
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1) +
  labs(title = "Plot of Reaction Time (ms.) by Condition",
       x = "Condition",
       y = "Reaction Time (ms.)") +
  theme_fivethirtyeight()
```



```
data_long %>%
  mutate(Condition = recode(condition,
    "complex_sentence" = "Complex Sentence",
    "simple_sentence" = "Simple Sentence")) %>%
  ggplot(aes(x = Condition, y = rt)) +
  geom_beeswarm(alpha = .25) +
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1) +
  labs(title = "Plot of Reaction Time (ms.) by Condition",
       x = "Condition",
       y = "Reaction Time (ms.)") +
  theme_economist()
```

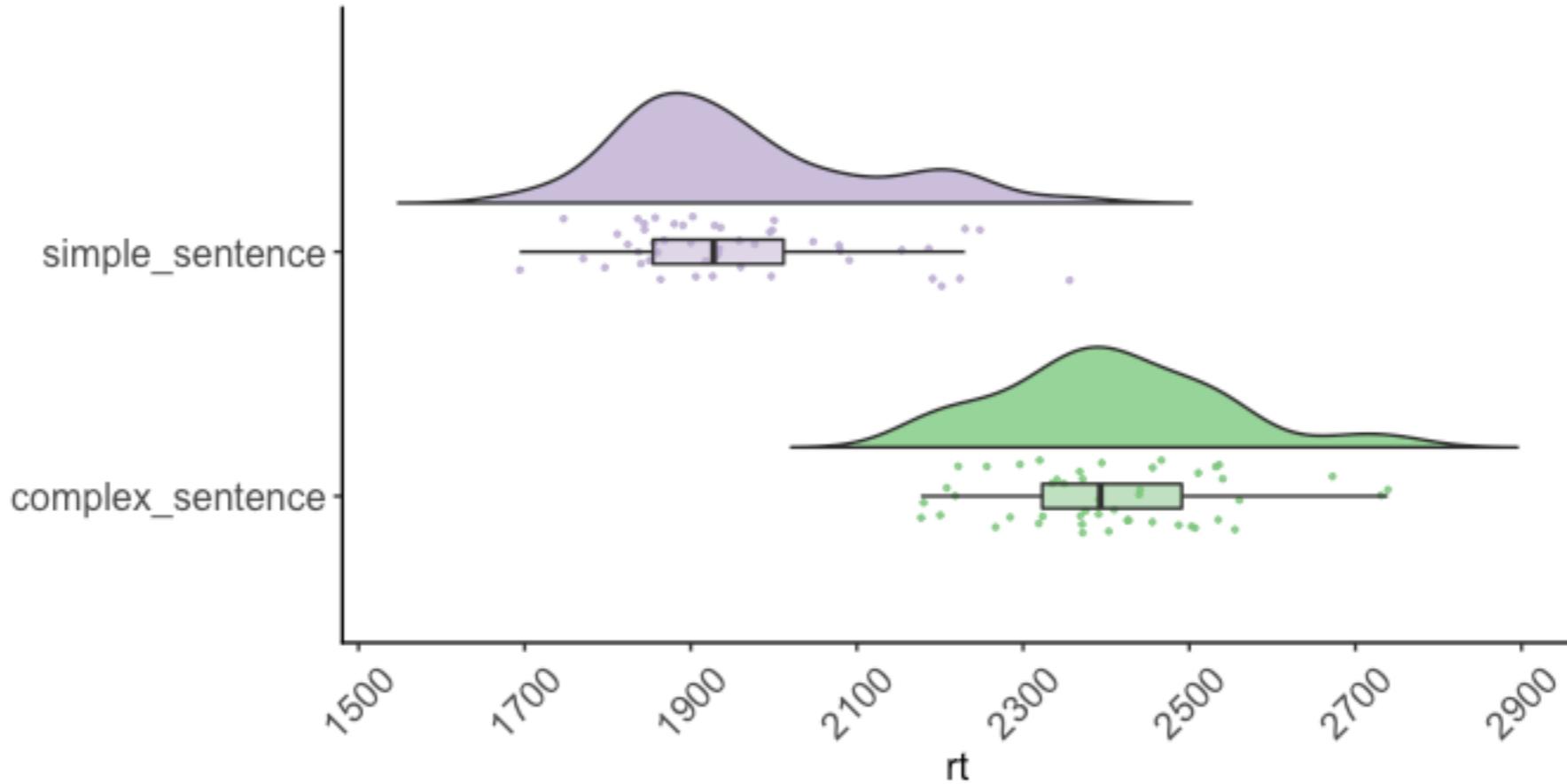
Changing Plot Dimensions



```
my_plot <- data_long %>%
  mutate(Condition = recode(condition,
    "complex_sentence" = "Complex Sentence",
    "simple_sentence" = "Simple Sentence")) %>%
  ggplot(aes(x = Condition, y = rt)) +
  geom_beeswarm(alpha = .25) +
  stat_summary(fun.data = "mean_cl_boot", colour = "black", size = 1) +
  labs(title = "Plot of Reaction Time (ms.) by Condition",
       x = "Condition",
       y = "Reaction Time (ms.)") +
  theme_economist()

ggsave("my_plot.png", my_plot, height = 8, width = 15, units = "cm")
```

Raincloud Plots



Developed by Micah Allen (UCL), raincloud plots allow you to see the raw data, and the shape of the distribution alongside a box plot (capturing the median, 25th and 75th percentiles as hinges, and $1.5 * \text{IQR}$ from the hinges as the whisker length.)

A Variety of Plots Using the Same Dataset

We're going to use the built-in dataset 'mpg' to build a variety of plots. First, let's find out about the data by using the head function to view the first part of the data.

```
> head(mpg)
# A tibble: 6 x 11
  manufacturer model displ year cyl trans   drv   cty   hwy fl class
  <chr>        <chr> <dbl> <int> <int> <chr> <chr> <int> <int> <chr> <chr>
1 audi         a4     1.8  1999     4 auto(15) f      18    29  p   compact
2 audi         a4     1.8  1999     4 manual(m5) f      21    29  p   compact
3 audi         a4     2.0  2008     4 manual(m6) f      20    31  p   compact
4 audi         a4     2.0  2008     4 auto(av)   f      21    30  p   compact
5 audi         a4     2.8  1999     6 auto(15) f      16    26  p   compact
6 audi         a4     2.8  1999     6 manual(m5) f      18    26  p   compact
```

We can explore the data further by asking for all the possibilities in each column using the `unique` function. For example, we can check to see how many different types of cars there are. Note that the \$ after the dataset name allows us to refer to a column in the mpg dataset.

```
> unique(mpg$manufacturer)
[1] "audi"        "chevrolet"    "dodge"       "ford"        "honda"       "hyundai"     "jeep"
[8] "land rover" "lincoln"      "mercury"     "nissan"     "pontiac"     "subaru"     "toyota"
[15] "volkswagen"
```

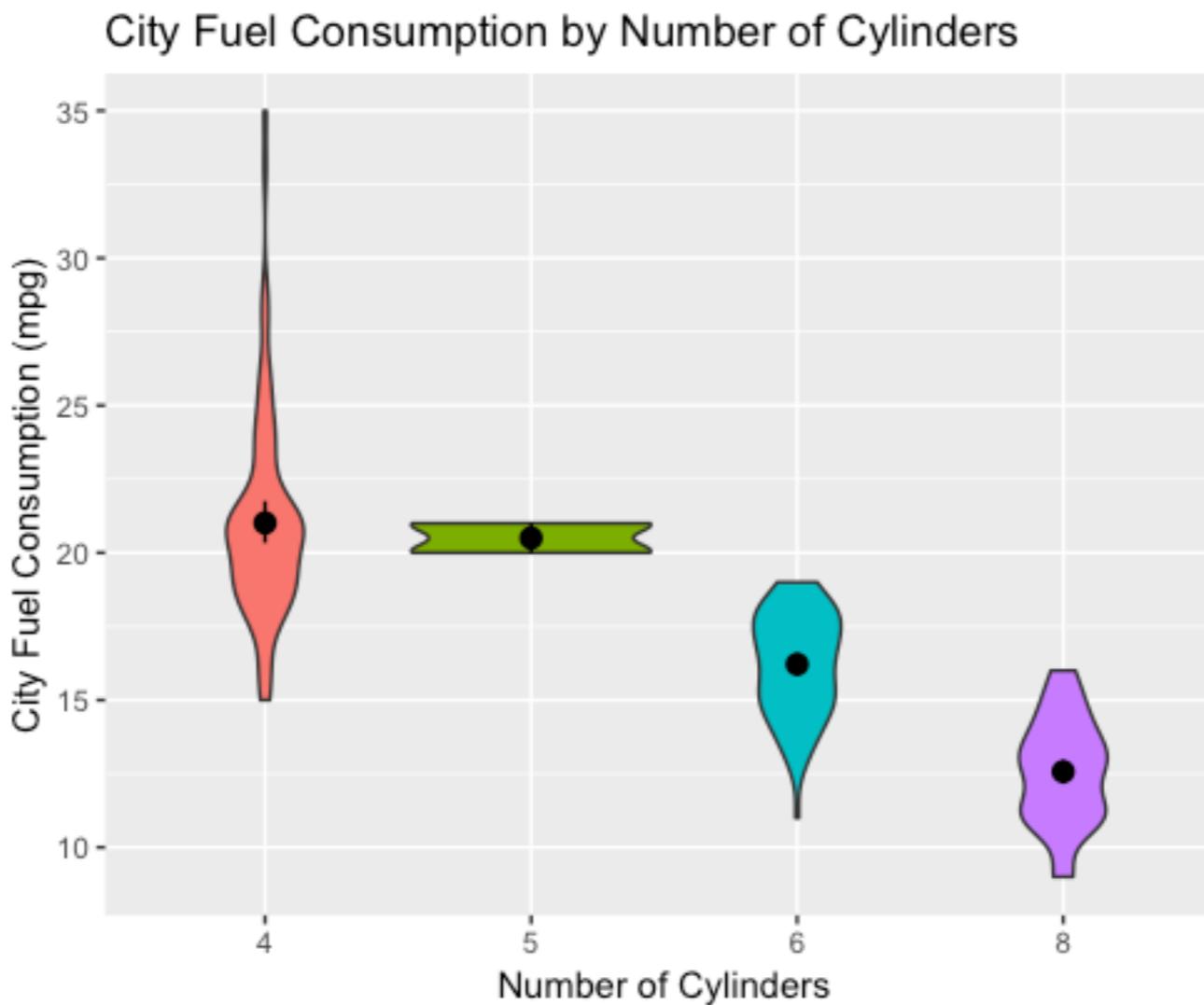
We can use the `length` function to give us the total number of unique possibilities:

```
> length(unique(mpg$manufacturer))
[1] 15
```

Let's look at a whole bunch of different visualisations using the mpg data set...

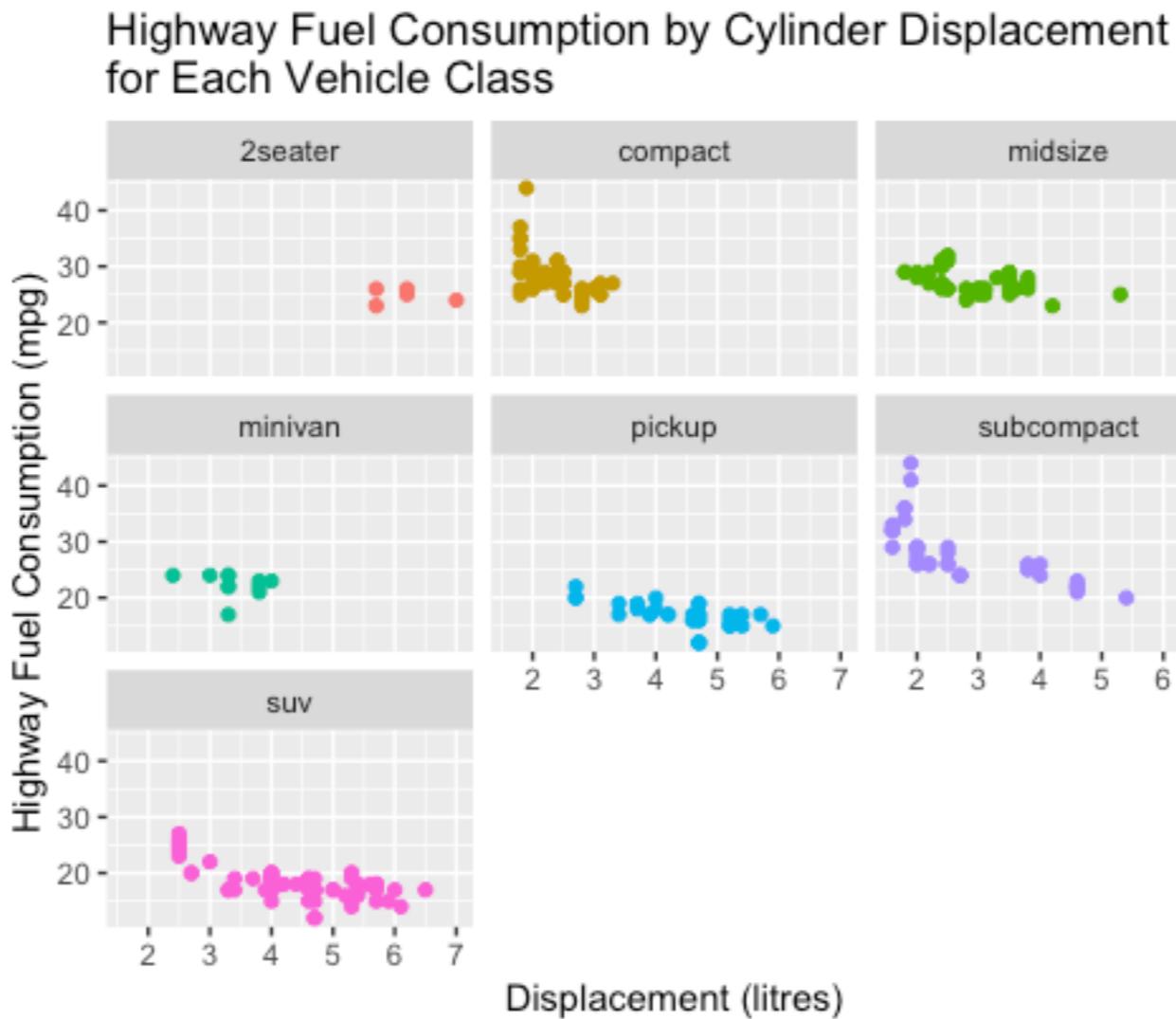
This illustrates the idea that there is not one 'correct' way to visualise the data, but rather that your choice of visualisation will be influenced by the question you're investigating, or the story you're wanting to tell...

Violin Plots



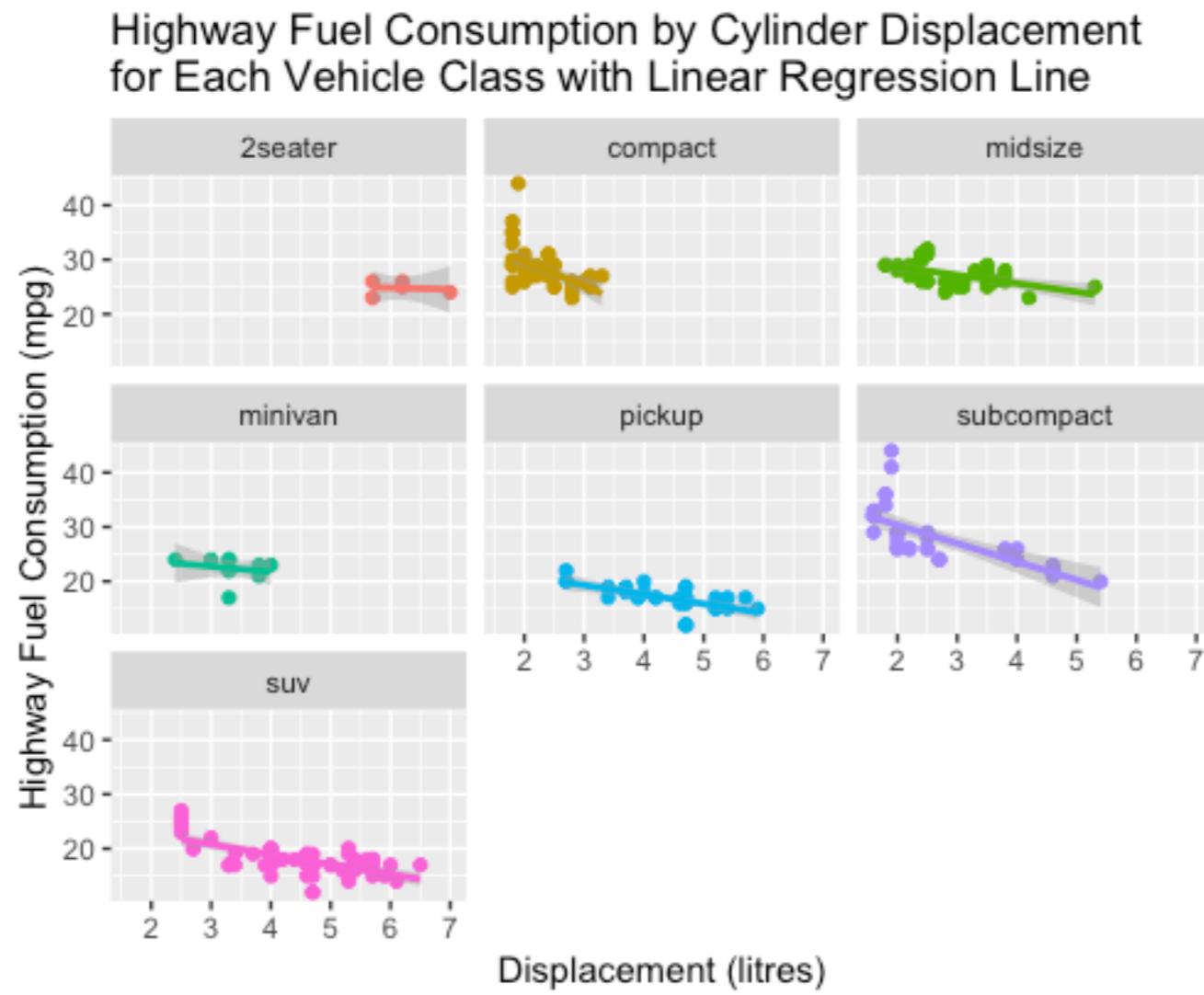
```
ggplot(mpg, aes(x = factor(cyl), y = cty, fill = factor(cyl))) +  
  geom_violin() +  
  guides(colour = FALSE, fill = FALSE) +  
  stat_summary(fun.data = mean_cl_boot, colour = "black", size = .5) +  
  labs(title = "City Fuel Consumption by Number of Cylinders",  
       x = "Number of Cylinders",  
       y = "City Fuel Consumption (mpg)")
```

Faceting



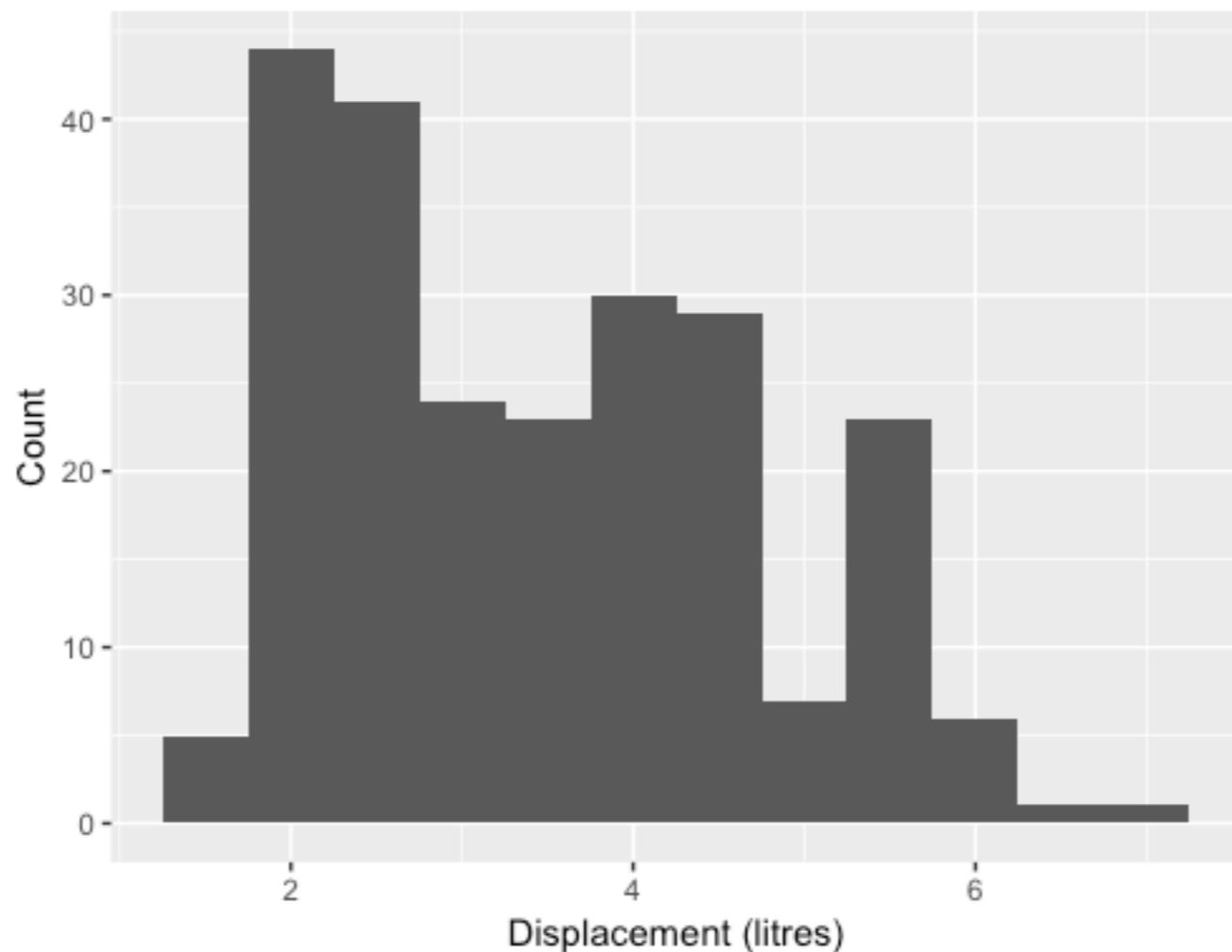
```
ggplot(mpg, aes(x = displ, y = hwy, colour = class)) +  
  geom_point() +  
  facet_wrap(~ class) +  
  guides(colour = FALSE) +  
  labs(title = "Highway Fuel Consumption by Cylinder Displacement \nfor Each Vehicle Class",  
       x = "Displacement (litres)",  
       y = "Highway Fuel Consumption (mpg)")
```

Adding a regression line



```
ggplot(mpg, aes(x = displ, y = hwy, colour = class)) +  
  geom_point() +  
  facet_wrap(~ class) +  
  guides(colour = FALSE) +  
  geom_smooth(method = "lm") +  
  labs(title = "Highway Fuel Consumption by Cylinder Displacement \nfor Each Vehicle Class",  
       x = "Displacement (litres)",  
       y = "Highway Fuel Consumption (mpg)")
```

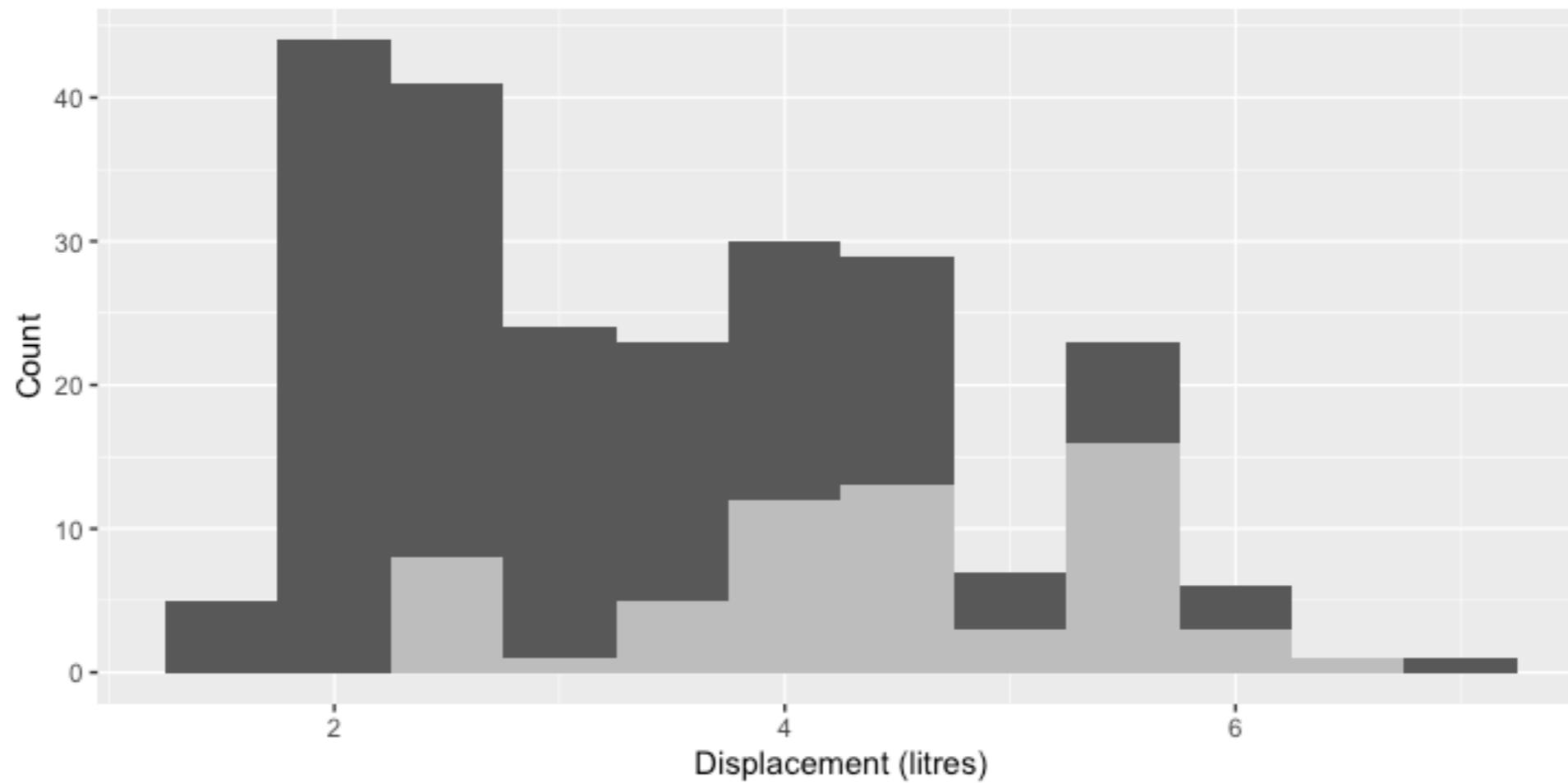
Histogram of Cylinder Displacement



```
ggplot(mpg, aes(x = displ)) +  
  geom_histogram(binwidth = .5) +  
  guides(fill = FALSE) +  
  labs(title = "Histogram of Cylinder Displacement",  
       x = "Displacement (litres)",  
       y = "Count")
```

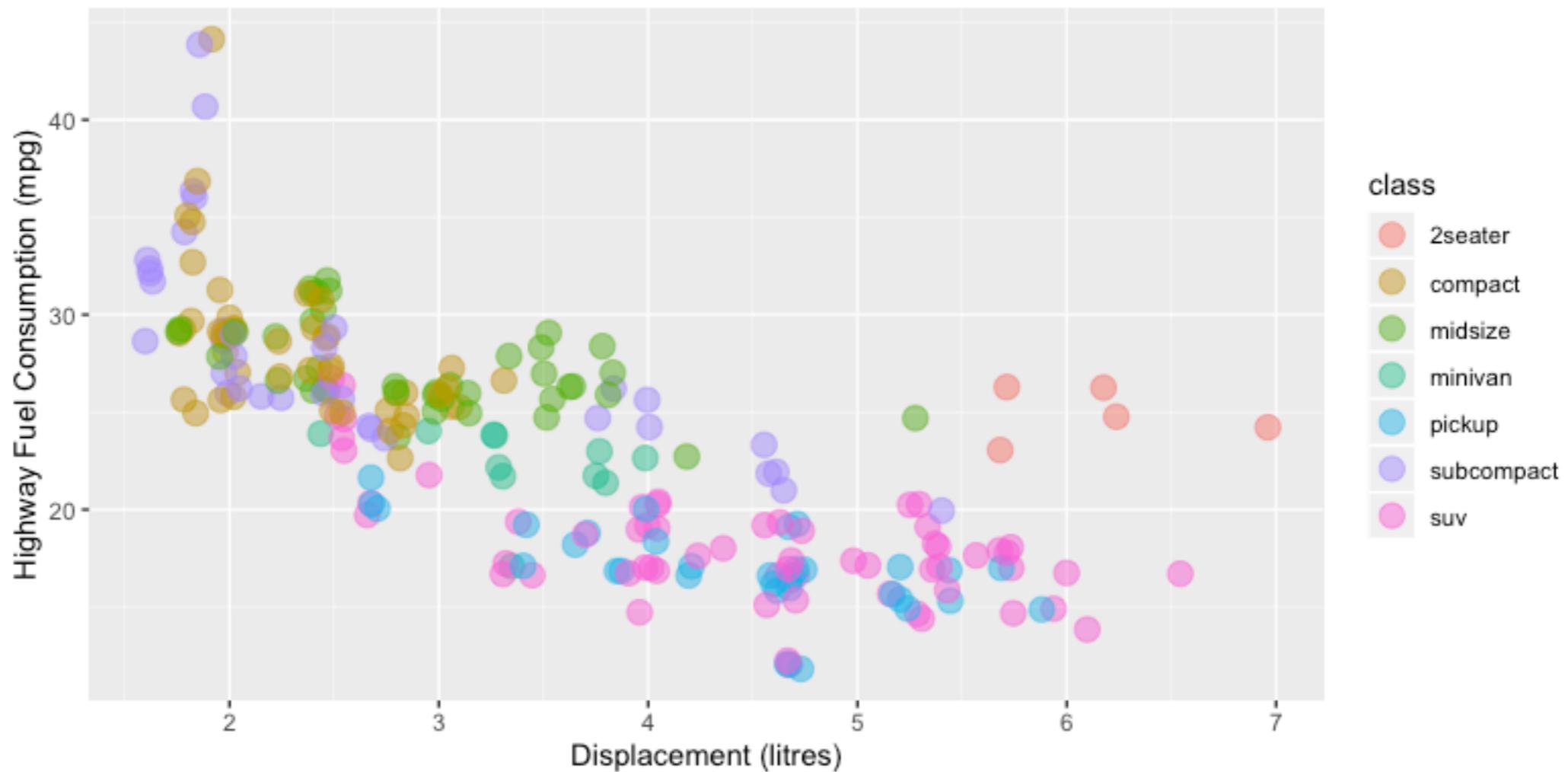
Histogram of Cylinder Displacement

SUVs highlighted



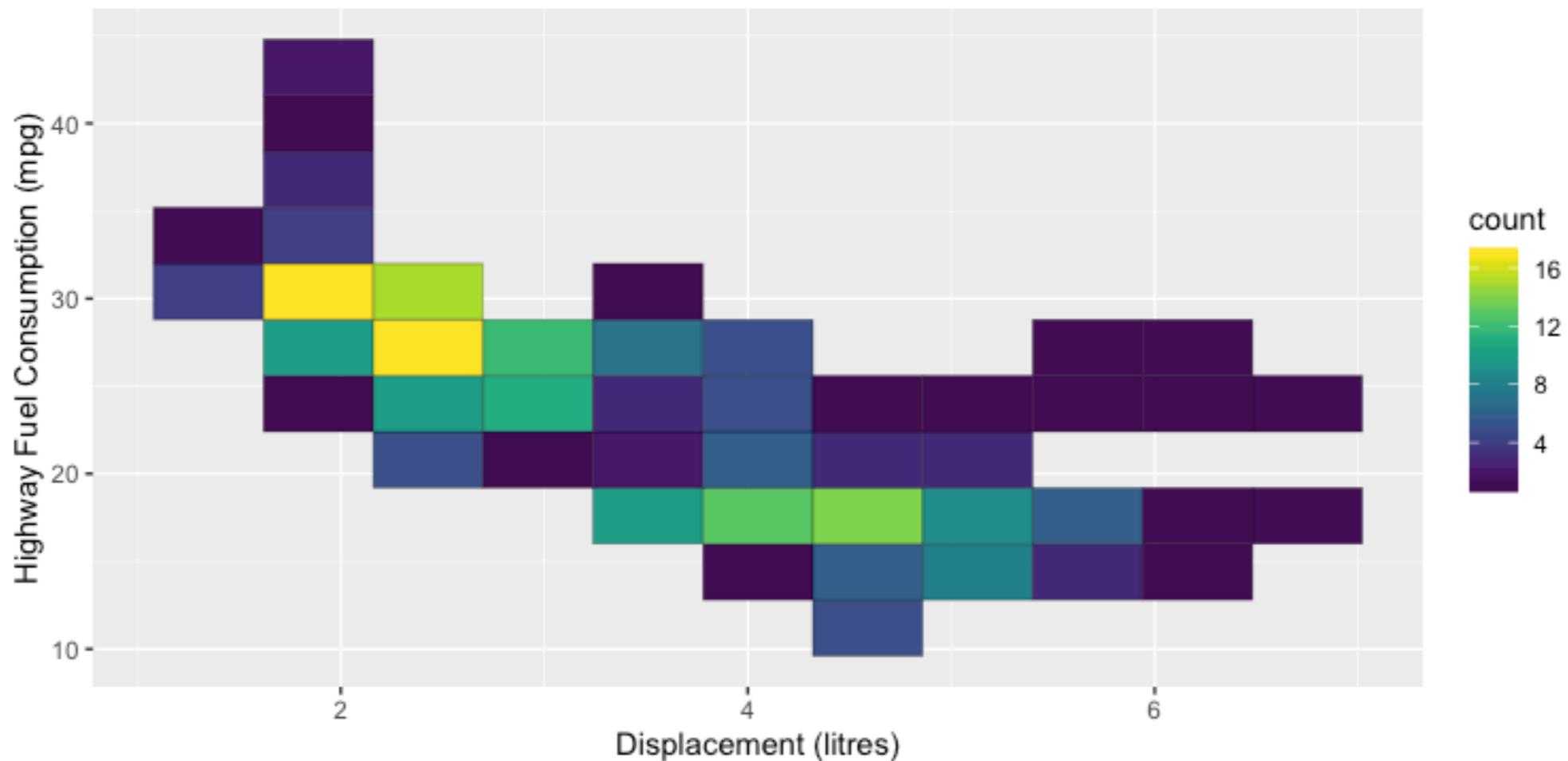
```
ggplot(mpg, aes(x = displ)) +  
  geom_histogram(binwidth = .5) +  
  geom_histogram(data = filter(mpg, class == "suv"), fill = "grey", binwidth = .5) +  
  guides(fill = FALSE) +  
  labs(title = "Histogram of Cylinder Displacement",  
       subtitle = "SUVs highlighted",  
       x = "Displacement (litres)",  
       y = "Count")
```

Scatterplot of Highway Fuel Consumption against
Engine Displacement Grouped by Class



```
ggplot(mpg, aes(x = displ, y = hwy, colour = class)) +  
  geom_jitter(width = 0.05, alpha = .5, size = 4) +  
  labs(title = "Scatterplot of Highway Fuel Consumption against\nEngine  
Displacement Grouped by Class",  
    x = "Displacement (litres)",  
    y = "Highway Fuel Consumption (mpg)")
```

Density heatmap of Highway Fuel Consumption against Engine Displacement



```
ggplot(mpg, aes(x = displ, y = hwy)) +  
  stat_bin2d(bins = 10, colour = "black") +  
  scale_fill_viridis() +  
  labs(title = "Density heatmap of Highway Fuel Consumption against Engine Displacement",  
       x = "Displacement (litres)",  
       y = "Highway Fuel Consumption (mpg)")
```

Plotting Time Series Data

We can install the `gapminder` package which contains lots of interesting data about life expectancy, population size, GDP for lots of countries collected over lots of years.

```
> gapminder
# A tibble: 1,704 x 6
  country      continent    year  lifeExp      pop  gdpPercap
  <fct>        <fct>     <int>   <dbl>     <int>      <dbl>
1 Afghanistan Asia       1952    28.8     8425333    779.
2 Afghanistan Asia       1957    30.3     9240934    821.
3 Afghanistan Asia       1962    32.0    10267083    853.
4 Afghanistan Asia       1967    34.0    11537966    836.
5 Afghanistan Asia       1972    36.1    13079460    740.
6 Afghanistan Asia       1977    38.4    14880372    786.
7 Afghanistan Asia       1982    39.9    12881816    978.
8 Afghanistan Asia       1987    40.8    13867957    852.
9 Afghanistan Asia       1992    41.7    16317921    649.
10 Afghanistan Asia      1997    41.8    22227415    635.
# ... with 1,694 more rows
```

Animated visualisations

For datasets with time series information, we might think it could be easier to tell our data story if we animate by time.

The `ggridge` package allows us to create animated visualisations from within R which we can import or embed in R Markdown documents.

We need to install it via the usual route:

```
> install.packages("ggridge")
```

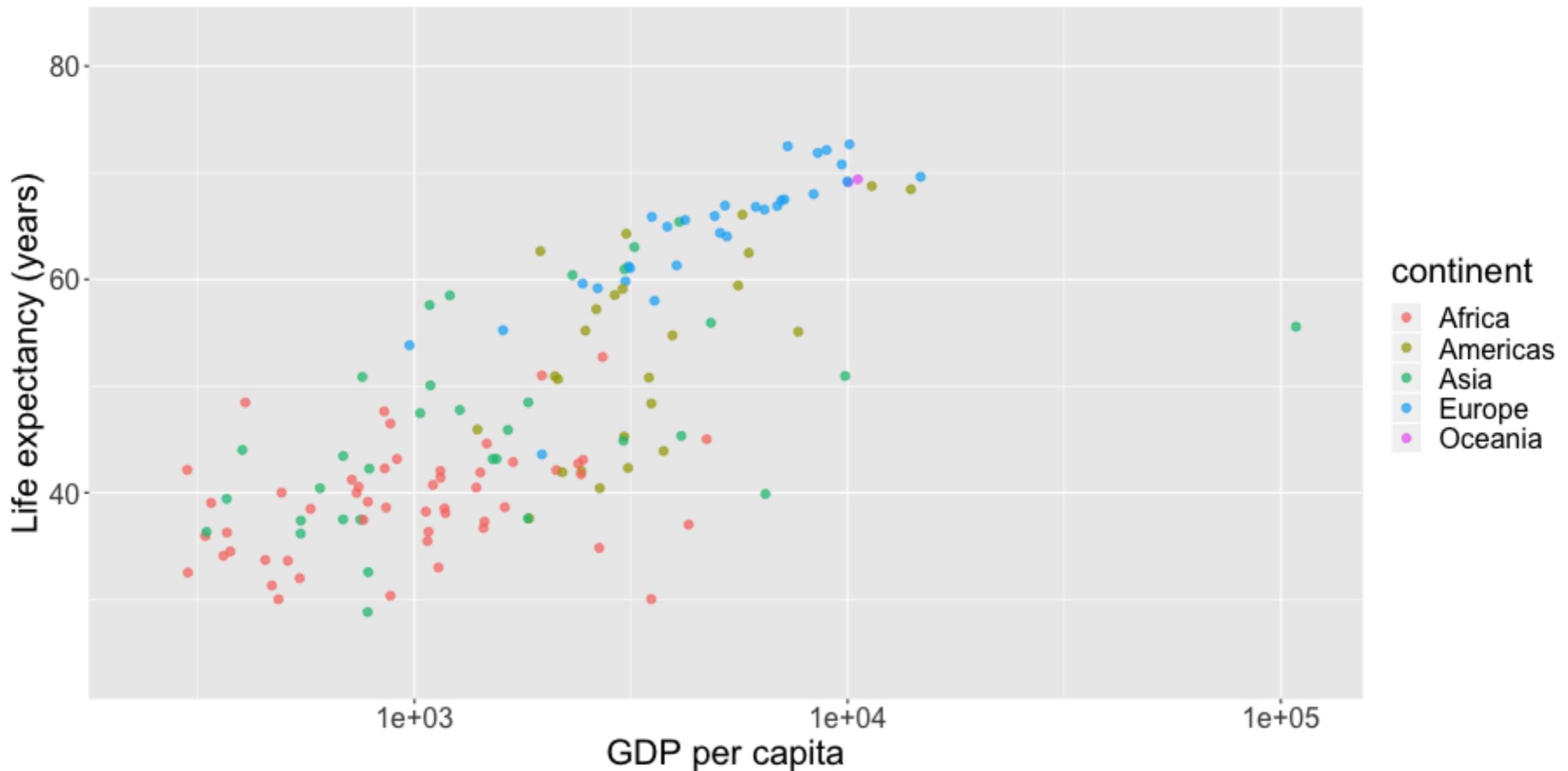
and then load it when we want to use it:

```
> library(ggridge)
```

Animated Time Series Data

Gapminder dataset

Year: 1952



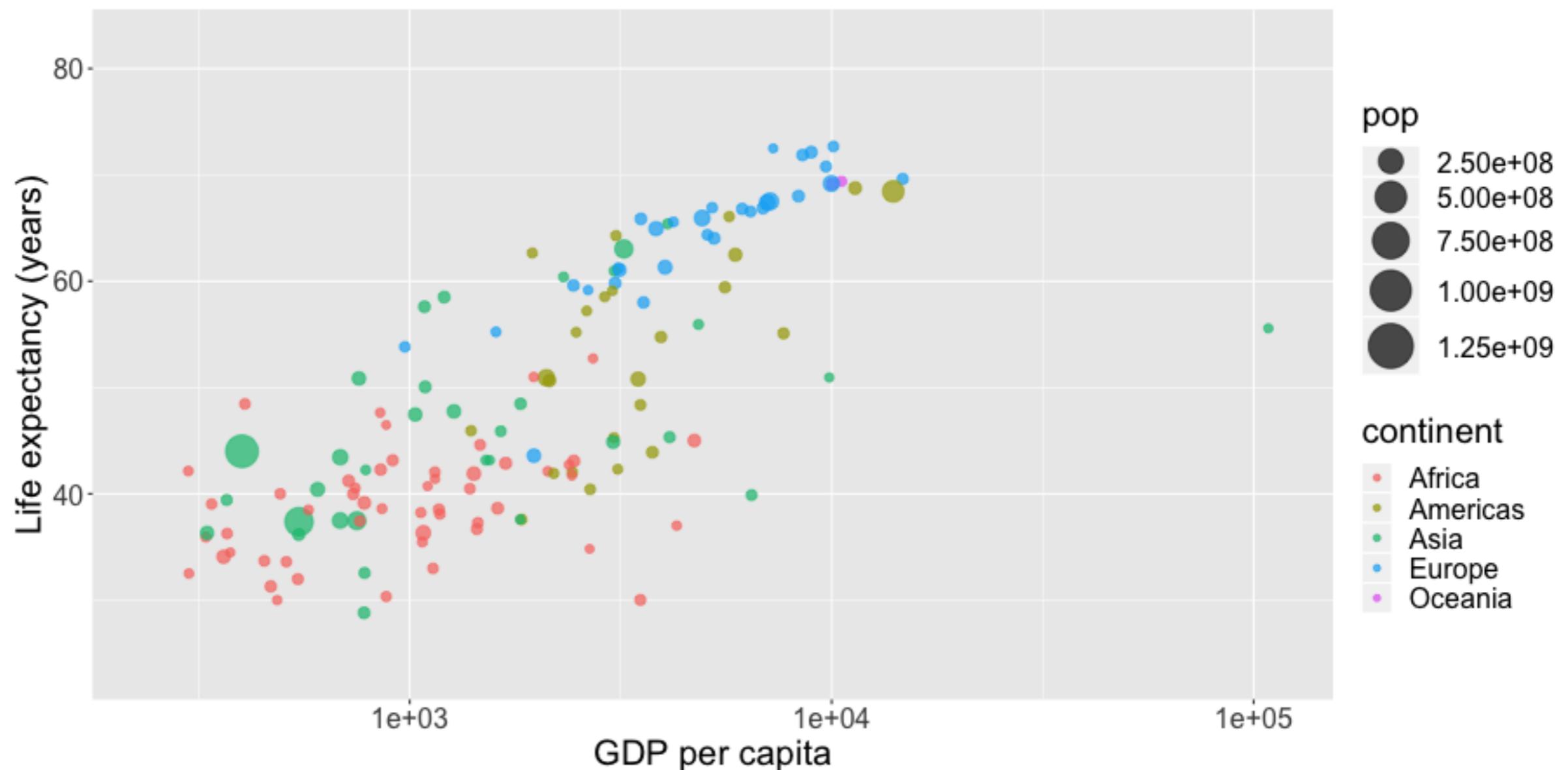
Visualising Data with 5 Variables Simultaneously

Animated Time Series Data

Now with a representation of population size.

Gapminder dataset

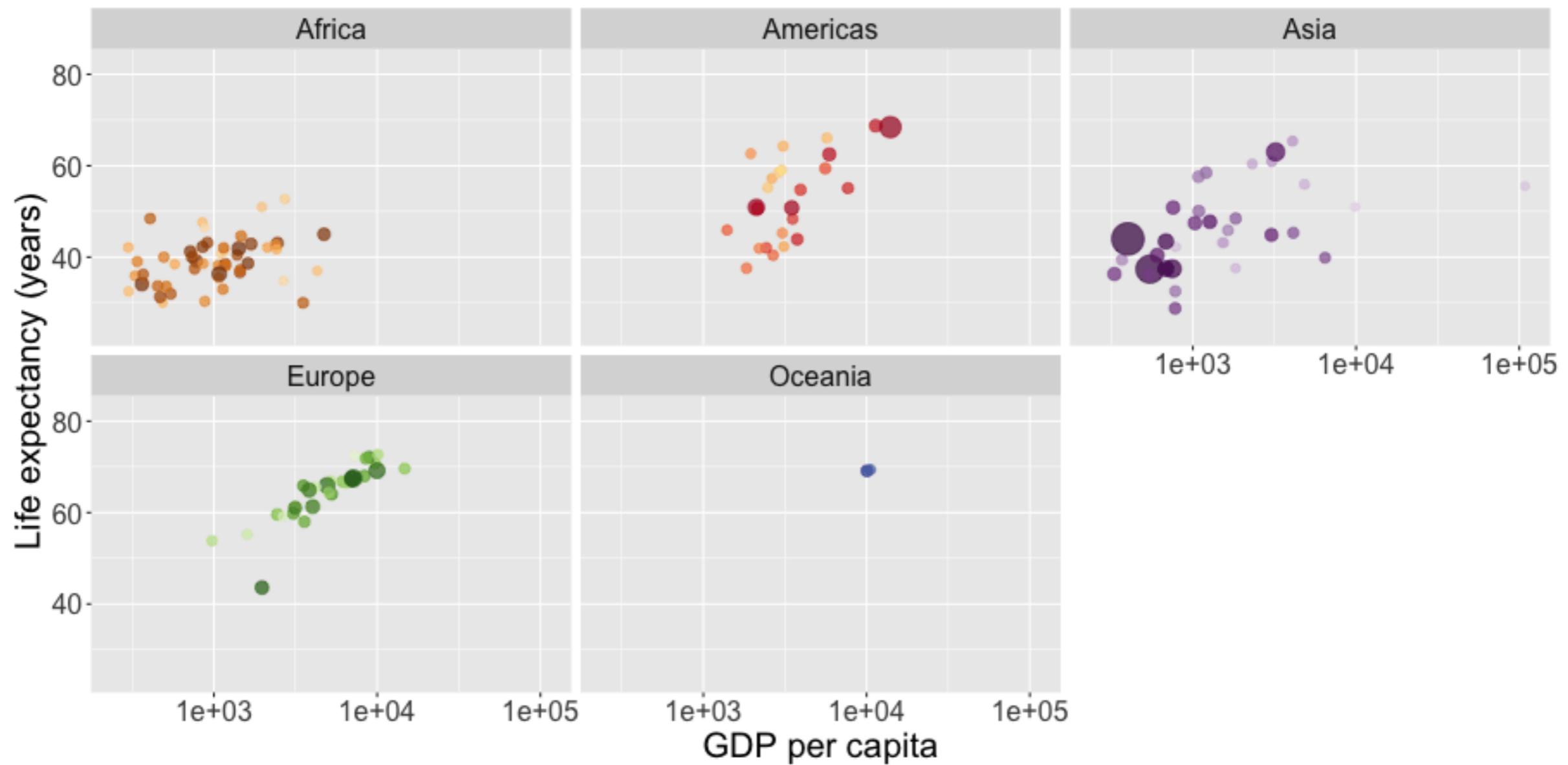
Year: 1952



Separately by Continent

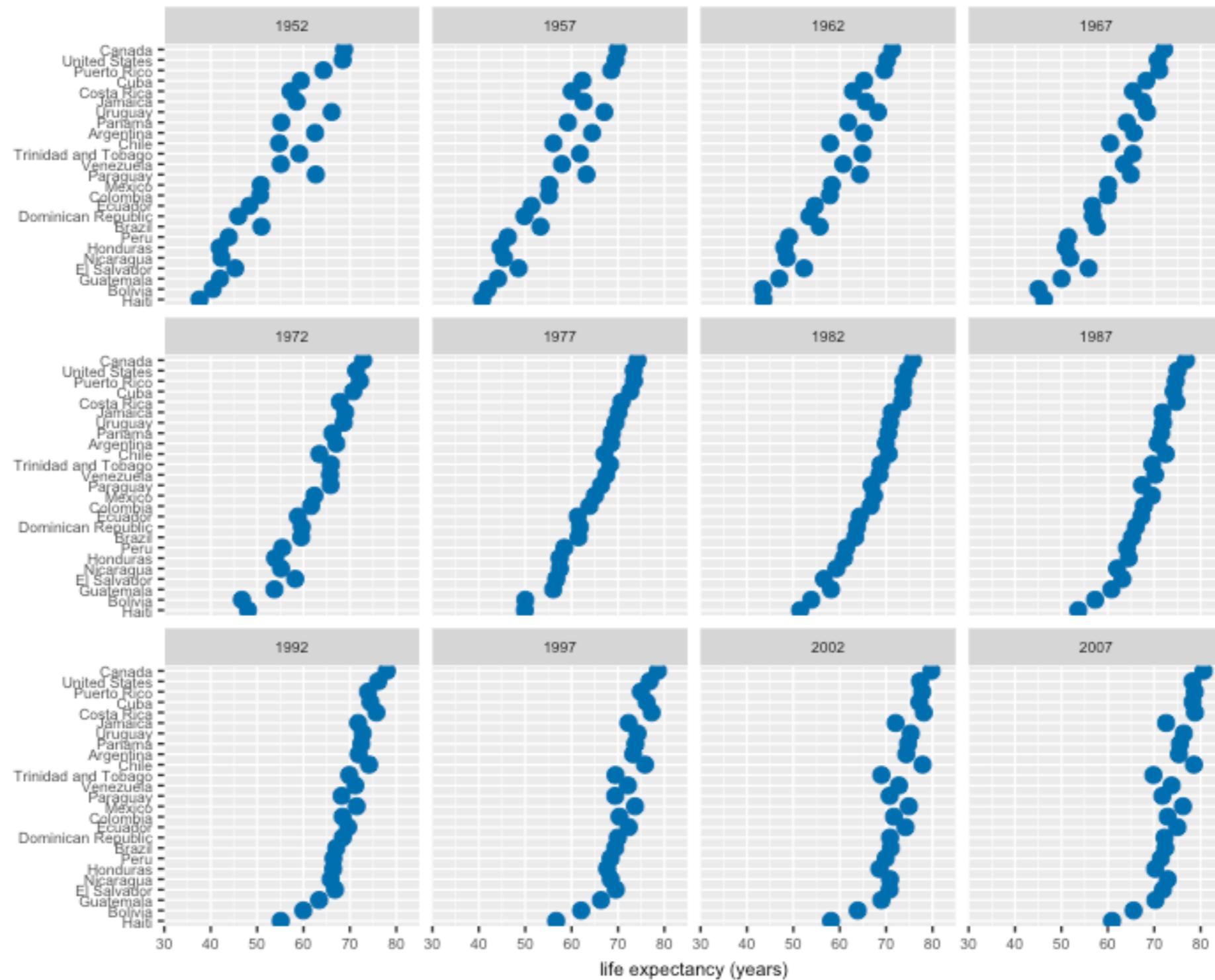
Gapminder dataset

Year: 1952

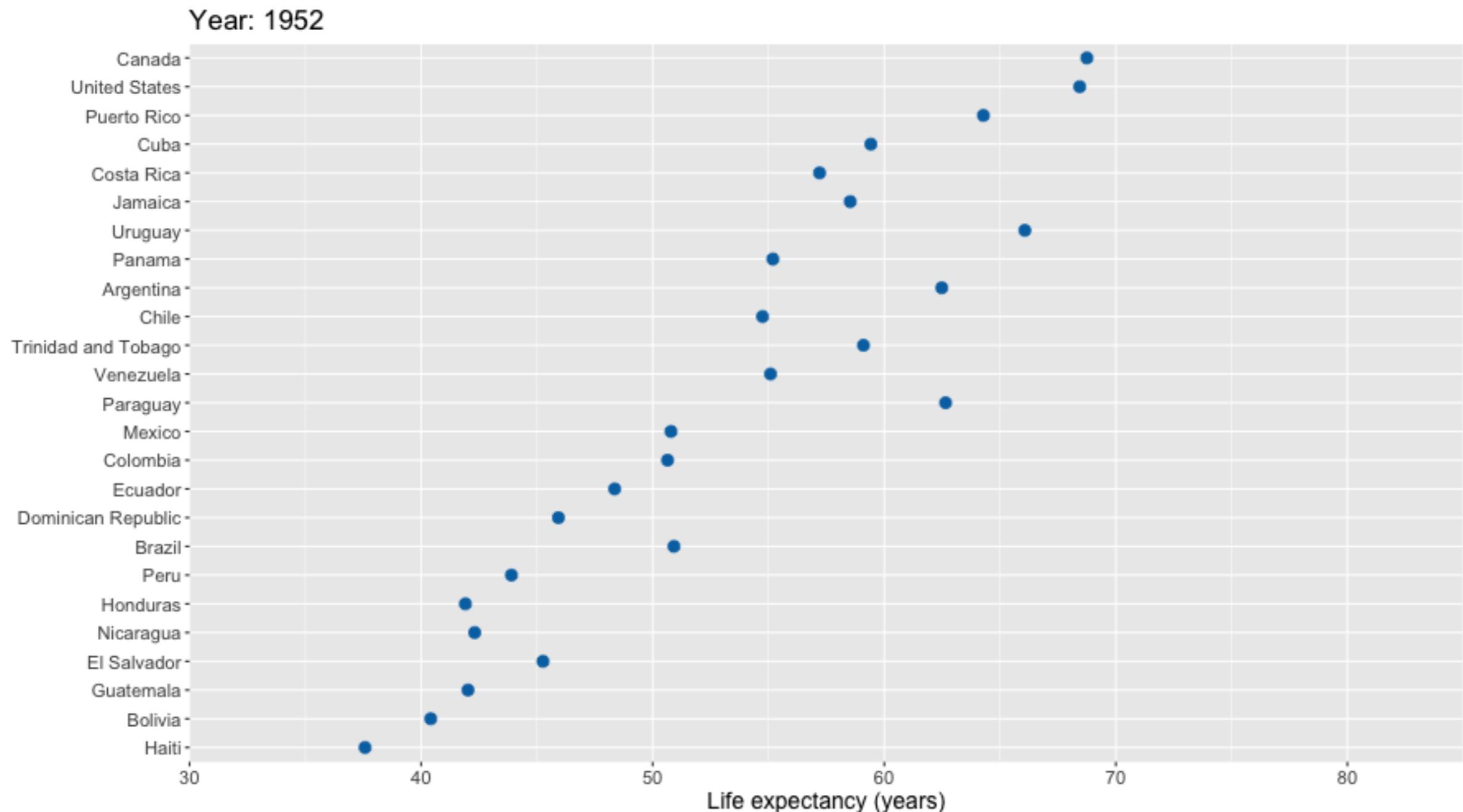


<https://github.com/thomasp85/gganimate>

Life Expectancy - Americas - Static



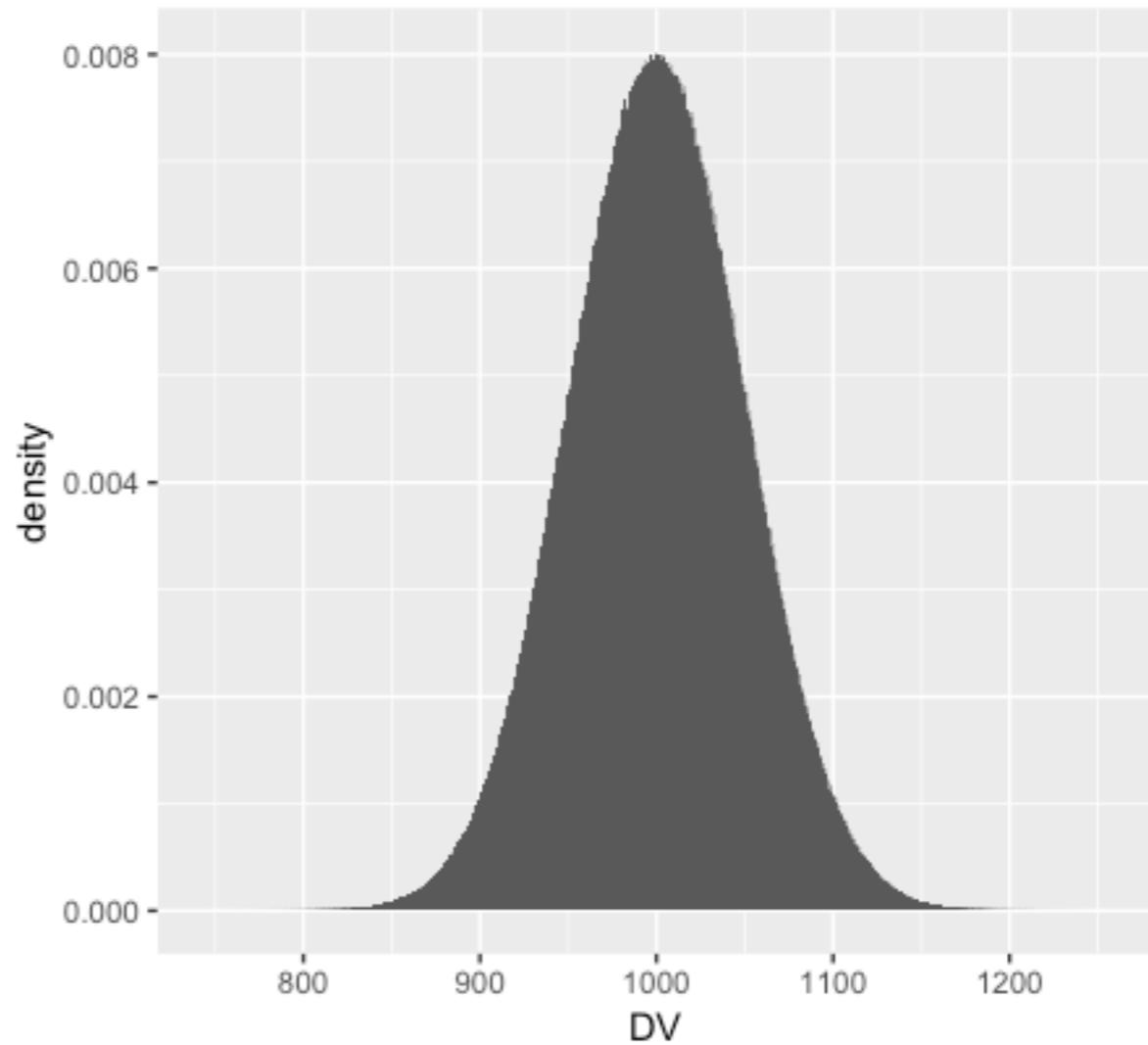
Life Expectancy - Americas - Animated



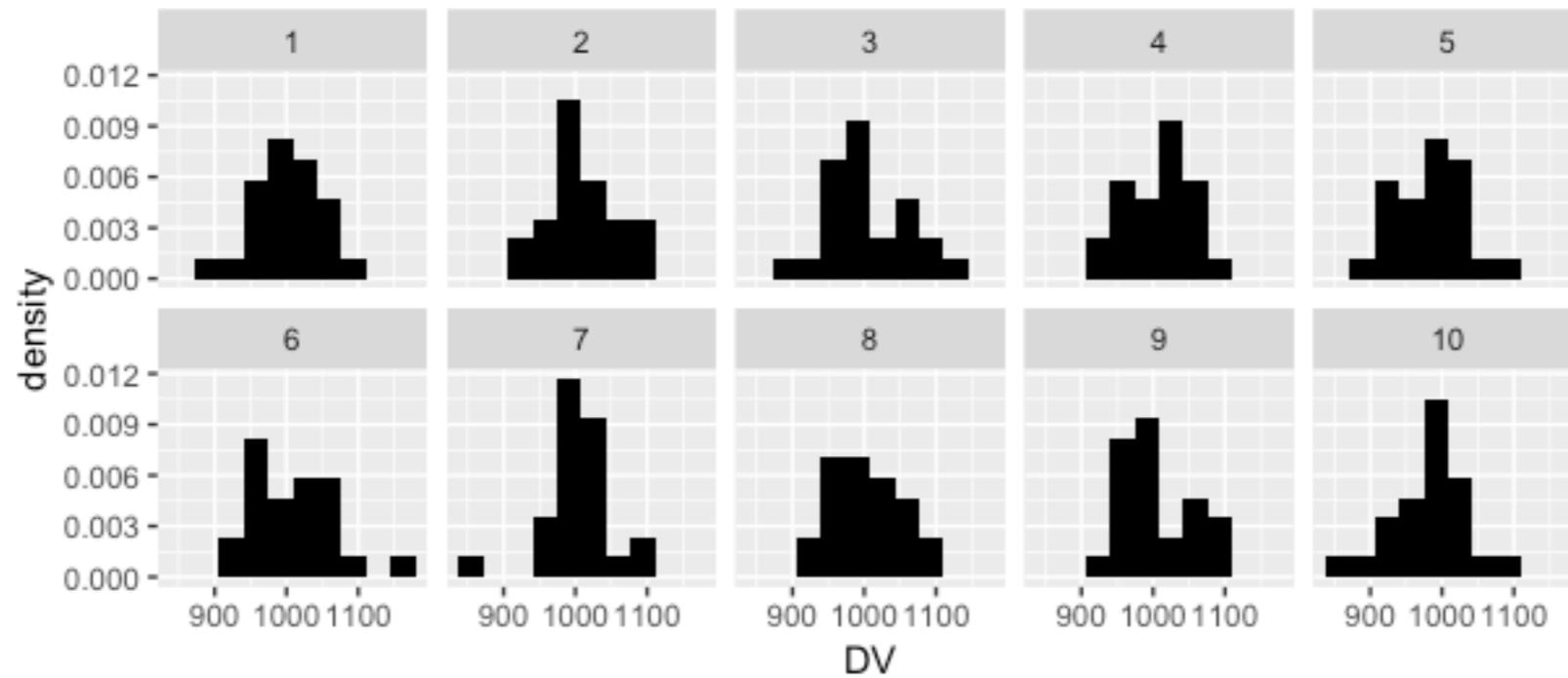
Although we have data only once every 5 years, the `gganimate` package interpolates between each census date to provide a smoother animation.

Using animation and data vis. to understand statistical concepts

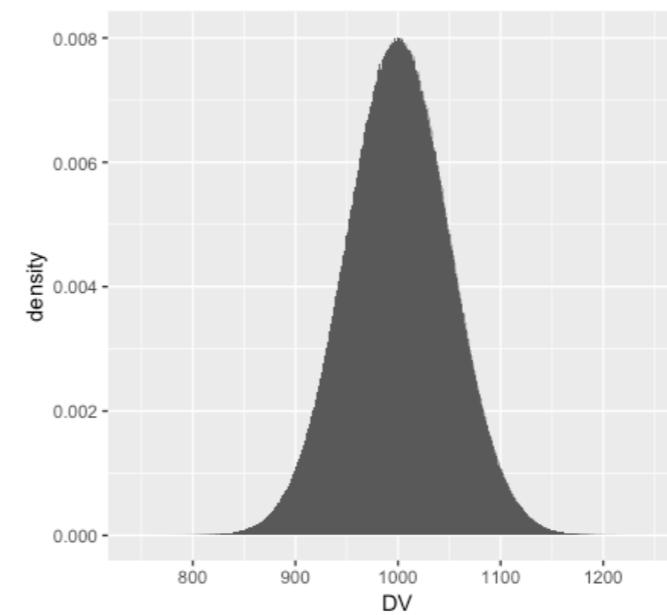
When we sample from a population, we are taking a sample of data points from the population distribution - the population could look something like this:



If our sample sizes are small, few sample distributions actually look like the population from which they're drawn and most sample means are a little different from the population mean:



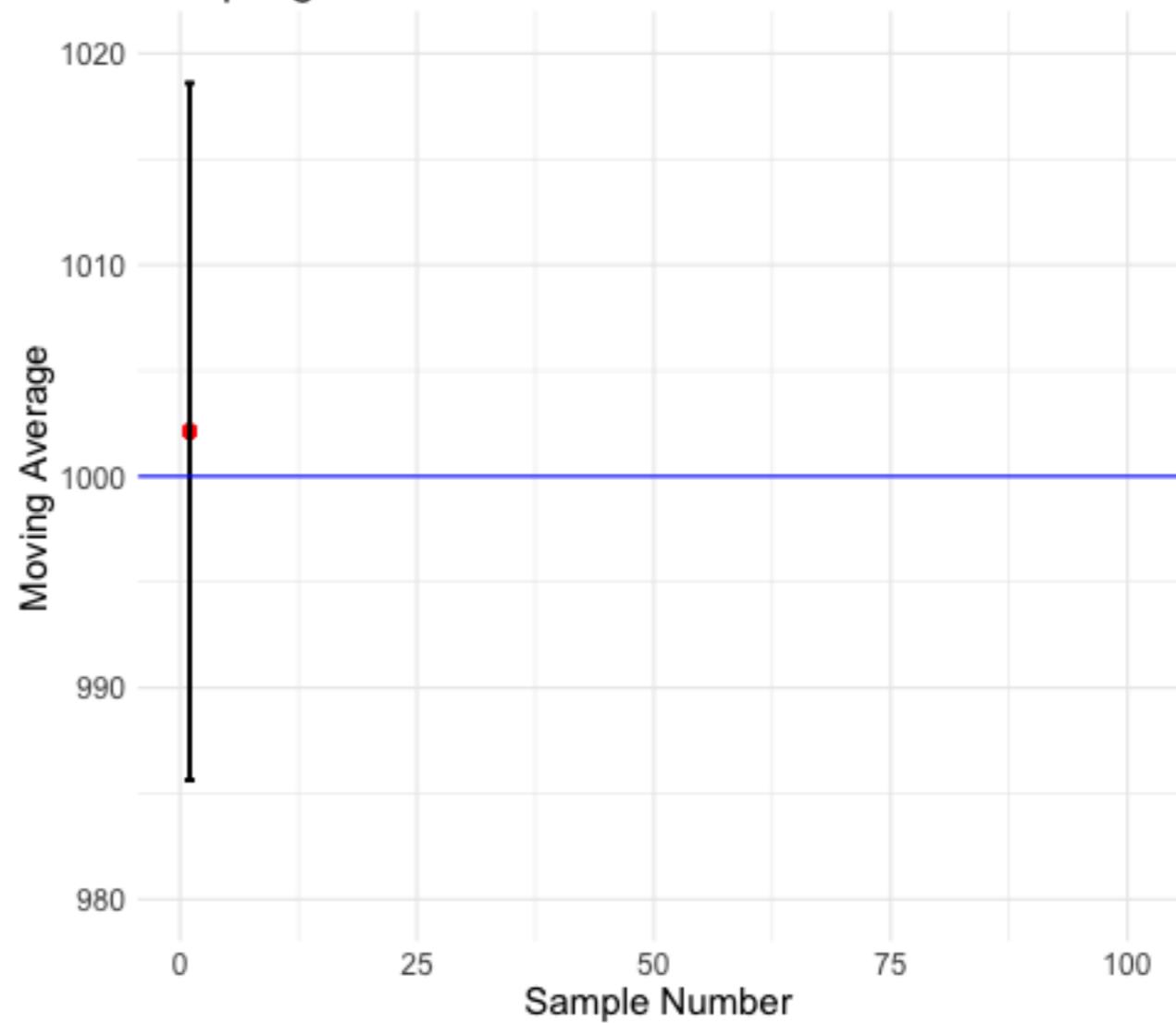
None of these look much like:



- When we take a sample from a population the mean of the sample may be quite different from the mean of the population (aka sampling error).
- If I take two samples, and work out the mean of these two samples, I should have a better estimate of the population mean than if I just looked at the mean of one of the samples.
- If I take three samples, work out the mean of these three samples etc. etc.

- The more samples (each with their own mean) we draw from the population, the closer we get to the true mean of the population. As sampling increases, the 95% CI bands around the mean narrow.
- So, animation can be used not just to communicate information, but also principles...

Moving average gets closer to the population mean (blue line) and CI bands narrow as sampling increases.



The Key Question

There is no such thing as the *best* way of visualising data - the method you choose will be determined by (e.g.) the type of data you have, the message you want to communicate, and the type of audience you will be communicating with.

Animations can be helpful, but they involve data being presented at a pace that might not suit the viewer - probably best suited for communicating time series data.

Podcast

The screenshot shows the homepage of the 'NOT SO STANDARD DEVIATIONS' website. The header includes links to Google Scholar, Scopus, jobs.ac.uk, BBC News, Chester Weather, The Grauniad, Chester Weather Station, GitHub, and My Website. A search bar and navigation links for About and Episodes are also present. The main banner features a stylized illustration of two hosts, Roger Peng and Hilary Parker, smiling. The title 'NOT SO STANDARD DEVIATIONS' is prominently displayed on the left, and the subtitle 'with Roger Peng and Hilary Parker' is on the right. Below the banner is a dark rectangular area containing social media sharing icons for Twitter, Email, and RSS. The episode page for '87 - Leadership Skills for the Apocalypse' is shown below, featuring a player interface with playback controls, episode details, and show notes.

To worksheet_1...