



Making your Workflow work for you

Chloe Yap

MD-PhD candidate

Cognitive Health Genomics (Gratten)

MRI-UQ, PCTG-IMB

I begin today by acknowledging the Traditional Custodians of the land on which we meet, and pay respects to their elders past, present and emerging.

About me



About me

Background

Undergraduate: BSc (BiomedSc), spent more time in lab than in lectures ...

Postgraduate: MD-PhD, PhD in autism omics (boils down to data analysis)

R experience

Been “in the game” since end-2014

Organic, on-the-job learning

Eternally indebted to my “teachers” ...

Priorities born of necessity

Multiple projects, distractors, on/off work → readability, organised directories, “centralised” results, working with manuscript output in mind

Running many iterations of similar analyses → automate as much as possible



Aims of this talk

- I *think* I'm relatively efficient(?) → some insights into how I operate
 - Disclaimer: many things I do are pretty hacky and not best practice! But they're “good enough”...

The screenshot shows a RStudio interface with several tabs open, illustrating a complex bioinformatics pipeline. The tabs include:

- methylation_qc.sh
- find_dmr_chisq_Mike.R
- cg20533952
- methylation_metaanalysis_metal_hsu_minerva
- controltraits
- .Rmd x osca
- combp_dmr_pipeline.sh
- individualID methylationArray snpArray wholeGeno
- metagenomics_qc.sh
- gzcat_pairedend.sh
- pgs_sbayersr_sbayersreigen_azd.R
- metabolomics_osca_analysis.R
- merge_aab_qtab_grm_pcs.sh
- metabolomics_osca_analysis.sh
- objID ProbeChr GeneProbe_bp topSNP topSNP_chr t
- missing_877_vs_878_vs_885_etc.txt
- untitled
- gwas_downstream.sh
- pgs_sbayersr_mid.msh
- metabolomics_smr_driver.sh
- untitled
- untitled
- assoc_sleep_adjclin <- fread(paste(oreml_dir, "/lipids", transform_choice, "_sleep_covdemo.linear", sep = " ", na.strings = NA), fileEncoding = "UTF-8")
- colnames(assoc_sleep_adjclin)[which(colnames(assoc_sleep_adjclin) == "p")] <- "p_adjclinical"
- assoc_sleep.tmp <- data.frame(inner_join(assoc_sleep, assoc_sleep_adjclin, by = "Probe"), Pheno = "sleep problems")
- assoc_age <- fread(paste(oreml_dir, "/lipids", transform_choice, ".age_covdemo.linear", sep = "")) %>% dplyr::select(-c("Probe"))
- colnames(assoc_age_adjclin) <- fread(paste(oreml_dir, "/lipids", transform_choice, ".age_covdemoclinical.linear", sep = ""))
- colnames(assoc_age_adjclin)[which(colnames(assoc_age_adjclin) == "p")] <- "p_adjclinical"
- assoc_age.tmp <- data.frame(inner_join(assoc_age, assoc_age_adjclin, by = "Probe"), Pheno = "age")
- assoc_tanner <- fread(paste(oreml_dir, "/lipids", transform_choice, ".tannergenital_covdemo.linear", sep = "")) %>% dplyr::select(-c("Probe"))
- assoc_tanner_adjclin <- fread(paste(oreml_dir, "/lipids", transform_choice, ".tannergenital_covdemoclinical.linear", sep = ""))
- colnames(assoc_tanner_adjclin)[which(colnames(assoc_tanner_adjclin) == "p")] <- "p_adjclinical"
- assoc_tanner.tmp <- data.frame(inner_join(assoc_tanner, assoc_tanner_adjclin, by = "Probe"), Pheno = "Tanner (genital)")
- assoc_bmi <- fread(paste(oreml_dir, "/lipids", transform_choice, ".bmi_covdemo.linear", sep = "")) %>% dplyr::select(-c("Probe"))
- assoc_bmi_adjclin <- fread(paste(oreml_dir, "/lipids", transform_choice, ".bmi_covdemoclinical.linear", sep = ""))
- colnames(assoc_bmi_adjclin)[which(colnames(assoc_bmi_adjclin) == "p")] <- "p_adjclinical"
- assoc_bmi.tmp <- data.frame(inner_join(assoc_bmi, assoc_bmi_adjclin, by = "Probe"), Pheno = "BMI")
- assoc.tmp <- rbind(assoc_asd.tmp, assoc_iqdq.tmp, assoc_sleep.tmp, assoc_age.tmp, assoc_tanner.tmp, assoc_bmi.tmp)
- colnames(assoc.tmp) <- c("Probe", "p_covdemo", "p_covdemoclinical", "Pheno")
- assoc.tmp\$Lipid.species_header <- assoc.tmp\$Probe
- assoc.tmp <- inner_join(assoc.tmp, anno, by = "Lipid.species_header")
- assoc.tmp\$label <- ifelse(assoc.tmp\$p_covdemo <= (0.05/302) & assoc.tmp\$p_covdemoclinical <= (0.05/302), assoc.tmp\$Lipid.species_header, "Other")
- ggplot(assoc.tmp, aes(x = -log10(p_covdemo), y = -log10(p_covdemoclinical), label = label)) +
- geom_point()
- geom_smooth(method = "lm") +
- geom_vline(xintercept = -log10(0.05/302), colour = "red") +
- geom_hline(yintercept = -log10(0.05/302), colour = "gold") +
- geom_text_repel(size = 2, force = 5, colour = "black", min.segment.length = 0) +
- facet_wrap(~ Pheno, scales = "free", ncol = 1)
- ggsave(paste(oreml_dir, "/figures/lwas_p_covdemo_v_covdemoclinical.png", sep = ""), height = 12, width = 8)

Bottom status bar: Line 630, Column 4 Tab Size: 4 Plain Text

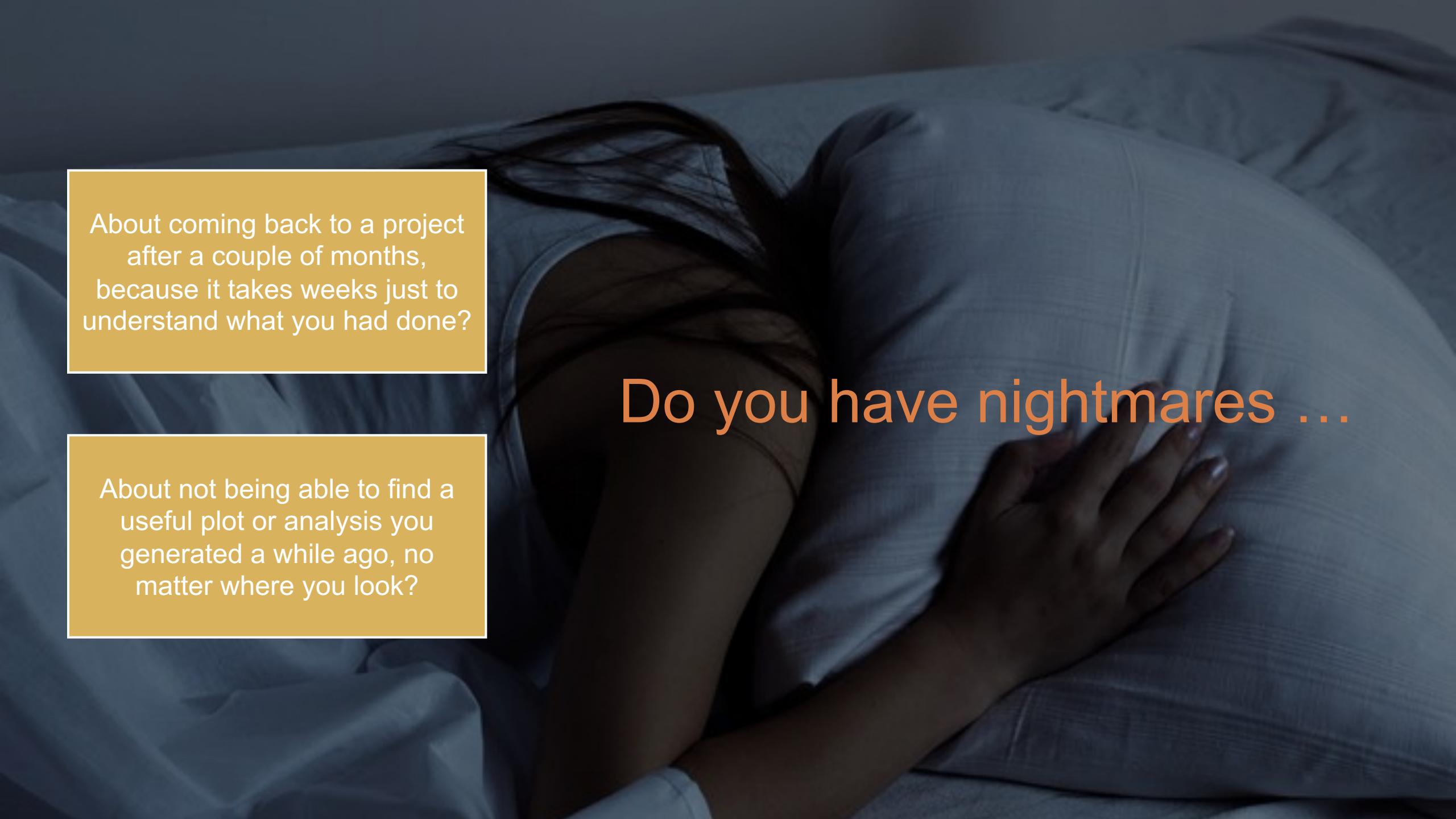
Bottom right corner: Tab Size: 4 R

Bottom right corner: -11.35, $\hat{\rho}_{\text{Pearson}} = -0.74$, $C_{95\%}^{HDI} [-0.88, -0.54]$, $r_{\text{beta}}^{JZS} = 1.41$

Outline

- Directory set-up, readability
- ~Automating repetitive analyses
- Results output: RMarkdown
- Keeping manuscript in mind
 - Plots



A dark, moody photograph of a person sleeping in bed. The person's head is turned to the side, and their hand is resting near their face. The lighting is low, creating a somber and contemplative atmosphere.

About coming back to a project
after a couple of months,
because it takes weeks just to
understand what you had done?

About not being able to find a
useful plot or analysis you
generated a while ago, no
matter where you look?

Do you have nightmares ...

Directory set-up: Principles

- Why?
 - Readability, much easier to start back where you left off
 - Helps to trace outputs to the data and scripts used
- Principles
 - Modular organisation, 1 folder per analysis
 - *If it will be a new figure or results section → new folder!*
 - *Try to keep a driver script for an analysis (.sh script with qsub or all .Rmd)*
 - *If you've put enough effort in typing code, it's probably worth filing away*
 - Consistency
 - Readability (?hackable?) → make and update README files
- See: Noble 2009 **PLOS Computational Biology A Quick Guide to Organizing Computational Biology Projects**
<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1000424>

OPEN ACCESS

EDUCATION

A Quick Guide to Organizing Computational Biology Projects

William Stafford Noble 

Published: July 31, 2009 • <https://doi.org/10.1371/journal.pcbi.1000424>

Article	Authors	Metrics	Comments	Media Coverage
				

Introduction

Principles

File and Directory Organization

The Lab Notebook

Carrying Out a Single Experiment

Handling and Preventing Errors

Command Lines versus Scripts versus Programs

The Value of Version Control

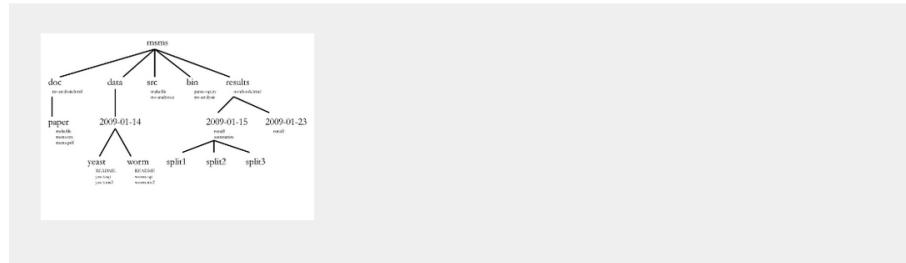
Conclusion

Acknowledgments

References

Reader Comments

Figures



Citation: Noble WS (2009) A Quick Guide to Organizing Computational Biology Projects. PLoS Comput Biol 5(7): e1000424. <https://doi.org/10.1371/journal.pcbi.1000424>

Editor: Fran Lewitter, Whitehead Institute, United States of America

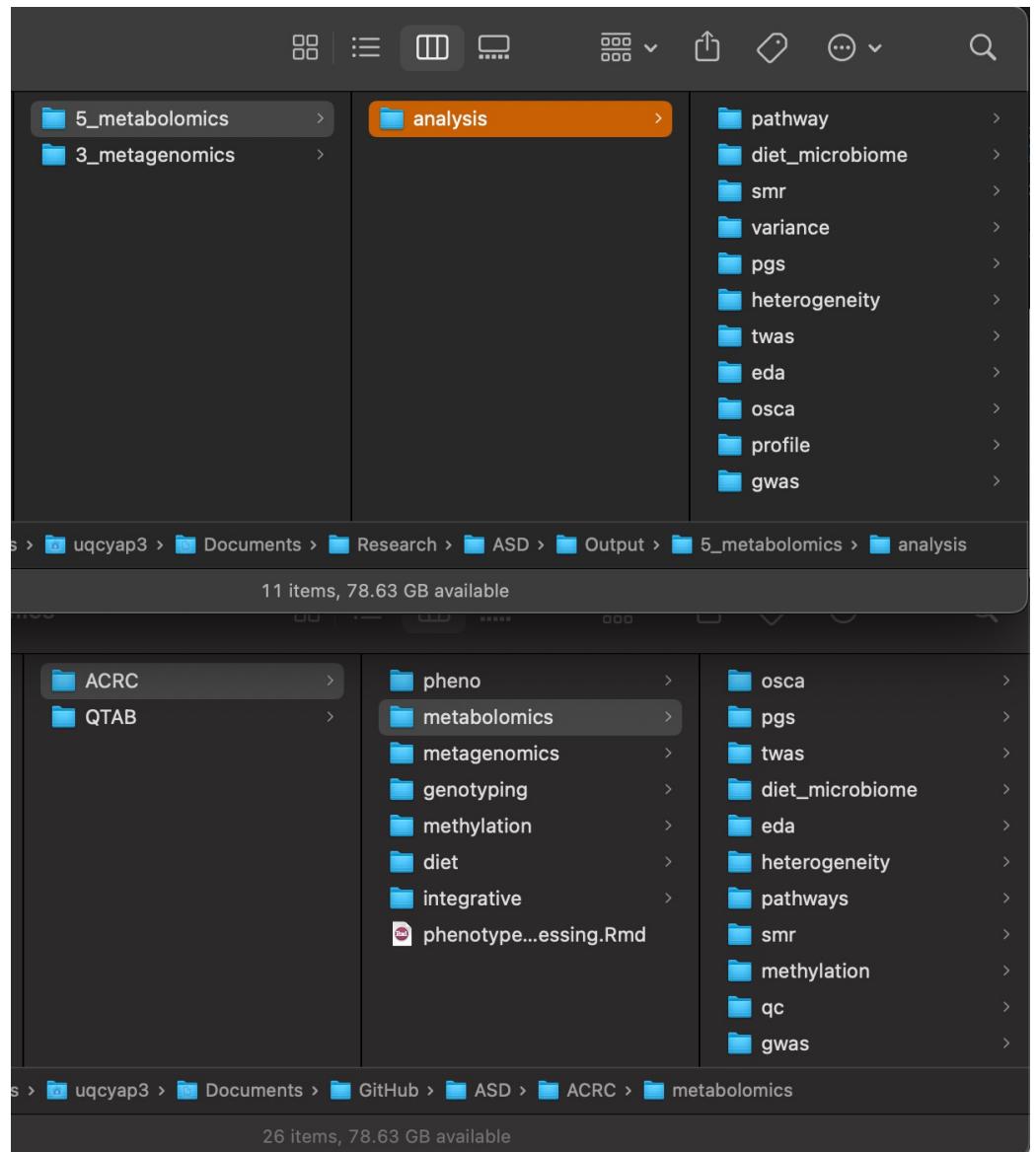
Published: July 31, 2009

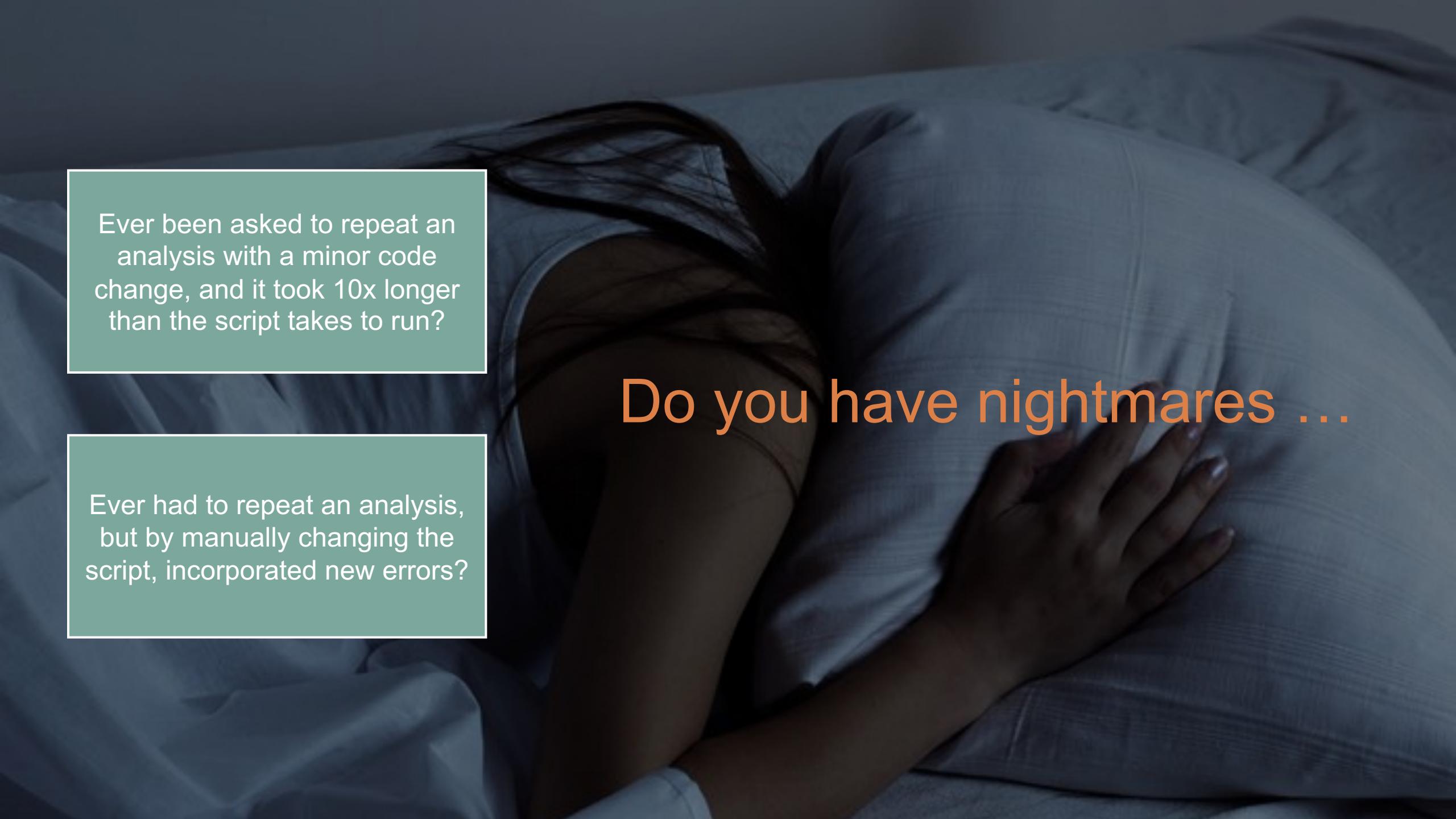
Copyright: © 2009 William Stafford Noble. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits

An example:

For each project ...

- Data/ Output/ Scripts/
 - Within each, aim for ~identical folder structure eg.
 - Data/gwas
 - Output/gwas
 - Scripts/gwas
- Write-up/
 - Manuscripts/
 - Presentations/
 - Applications/, etc.
- File naming
 - Logical
 - Date: YYMMDD so auto-sorts!



A dark, moody photograph of a person sleeping in bed. The person's head is turned away from the camera, and their hand is resting near their face. The lighting is low, creating a somber and intimate atmosphere.

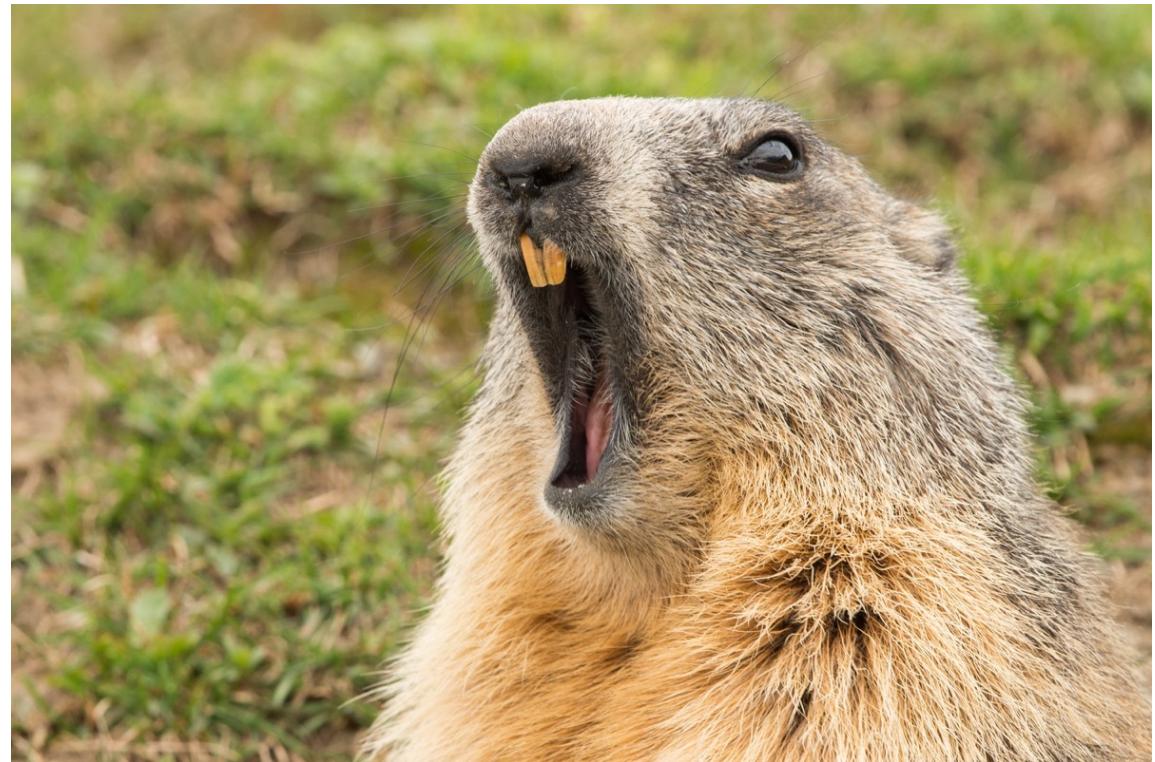
Ever been asked to repeat an analysis with a minor code change, and it took 10x longer than the script takes to run?

Ever had to repeat an analysis, but by manually changing the script, incorporated new errors?

Do you have nightmares ...

Automating repetitive analyses

- Why?
 - No matter how well you do it the first time, there's a 99% chance* you'll have to rerun the analysis ...
 - Massive time-saver (and hair-saver!)
 - Reduces errors
- Principles
 - Generalisability
 - Flexibility



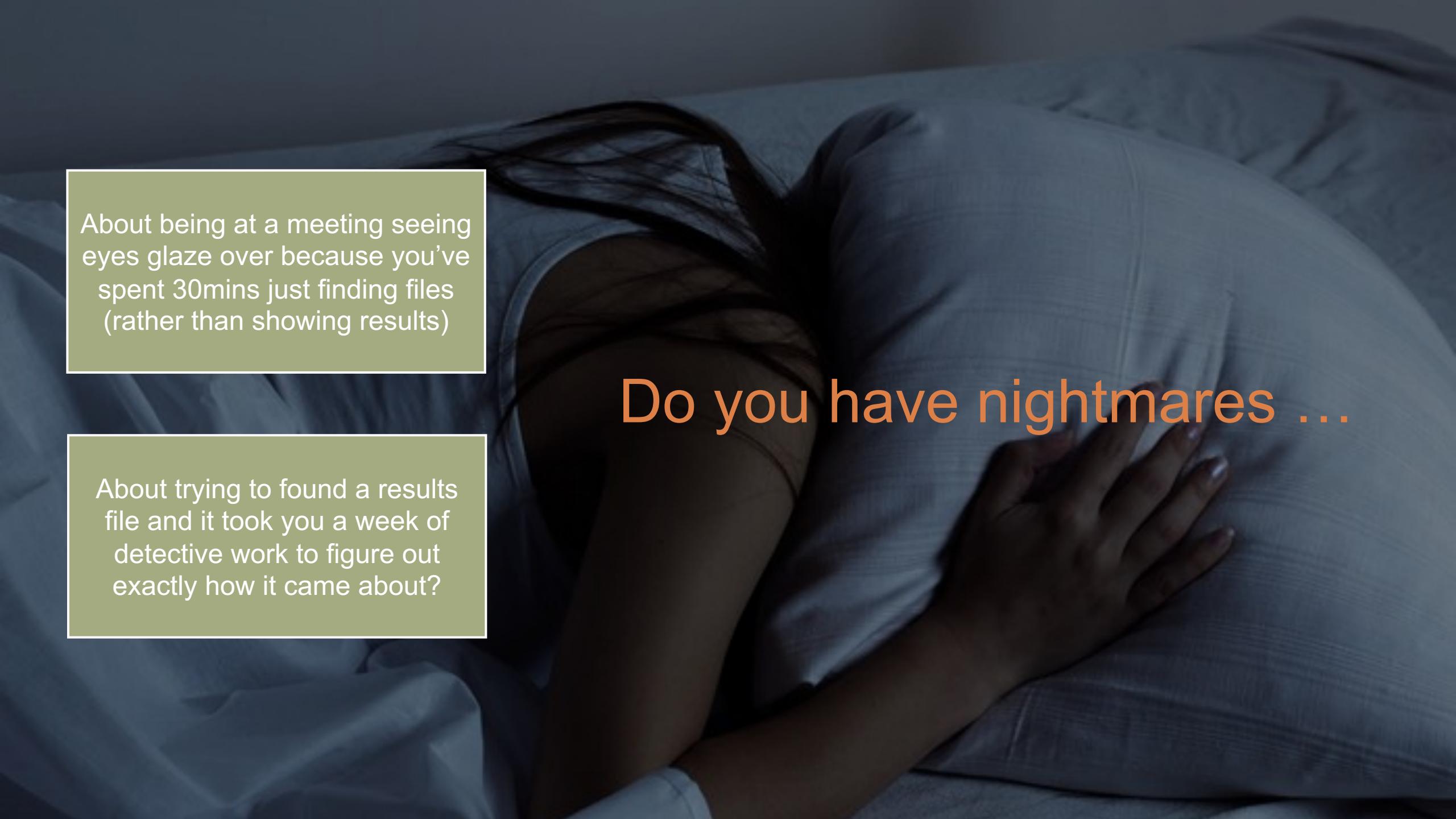
*not backed by empirical data

Automating repetitive analyses

- Create “settings”
 - Generate separate .Rmd with different settings
 - Settings reflected in output file name
 - Eg. “filen” variable
- `source()` scripts at the start of analysis
 - If there is common code at the start of multiple scripts
 - If input file is being updated
- On HPC → “pipeline” script in bash
 - .sh wrappers of .R scripts

Miscellaneous fun things in analysis

- `group_by()`
 - `%>% summarise()`
 - `%>% pivot_wider()`
 - Performing operations in parallel eg. `lm()`

A dark, moody photograph of a person sleeping peacefully in bed. The person is lying on their side, facing away from the camera, with their head resting on a pillow. A hand is visible, resting on the person's back. The lighting is low, creating a calm and restful atmosphere.

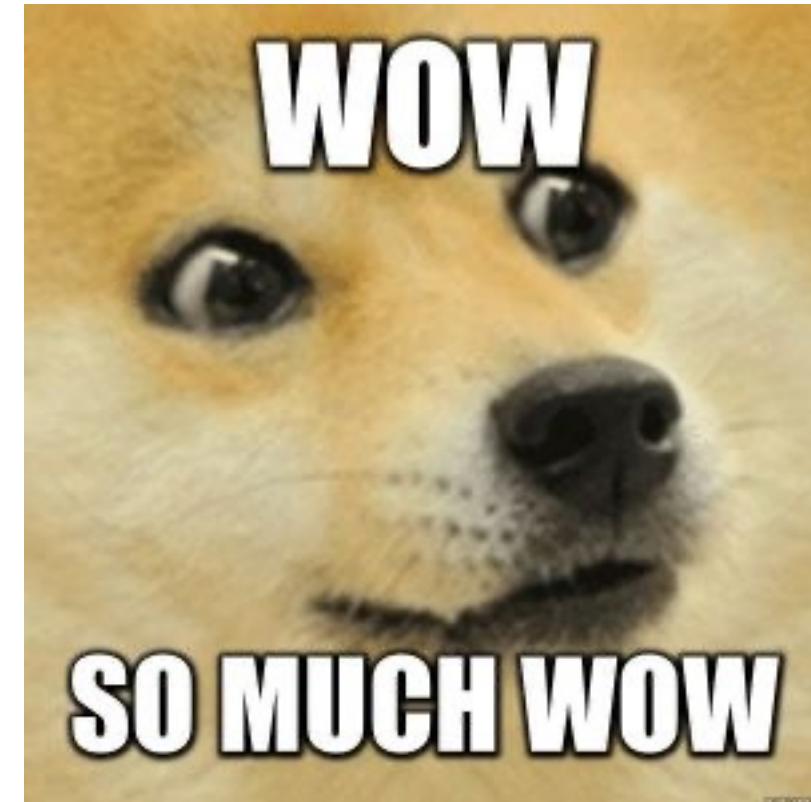
About being at a meeting seeing
eyes glaze over because you've
spent 30mins just finding files
(rather than showing results)

About trying to find a results
file and it took you a week of
detective work to figure out
exactly how it came about?

Do you have nightmares ...

Results output: RMarkdown

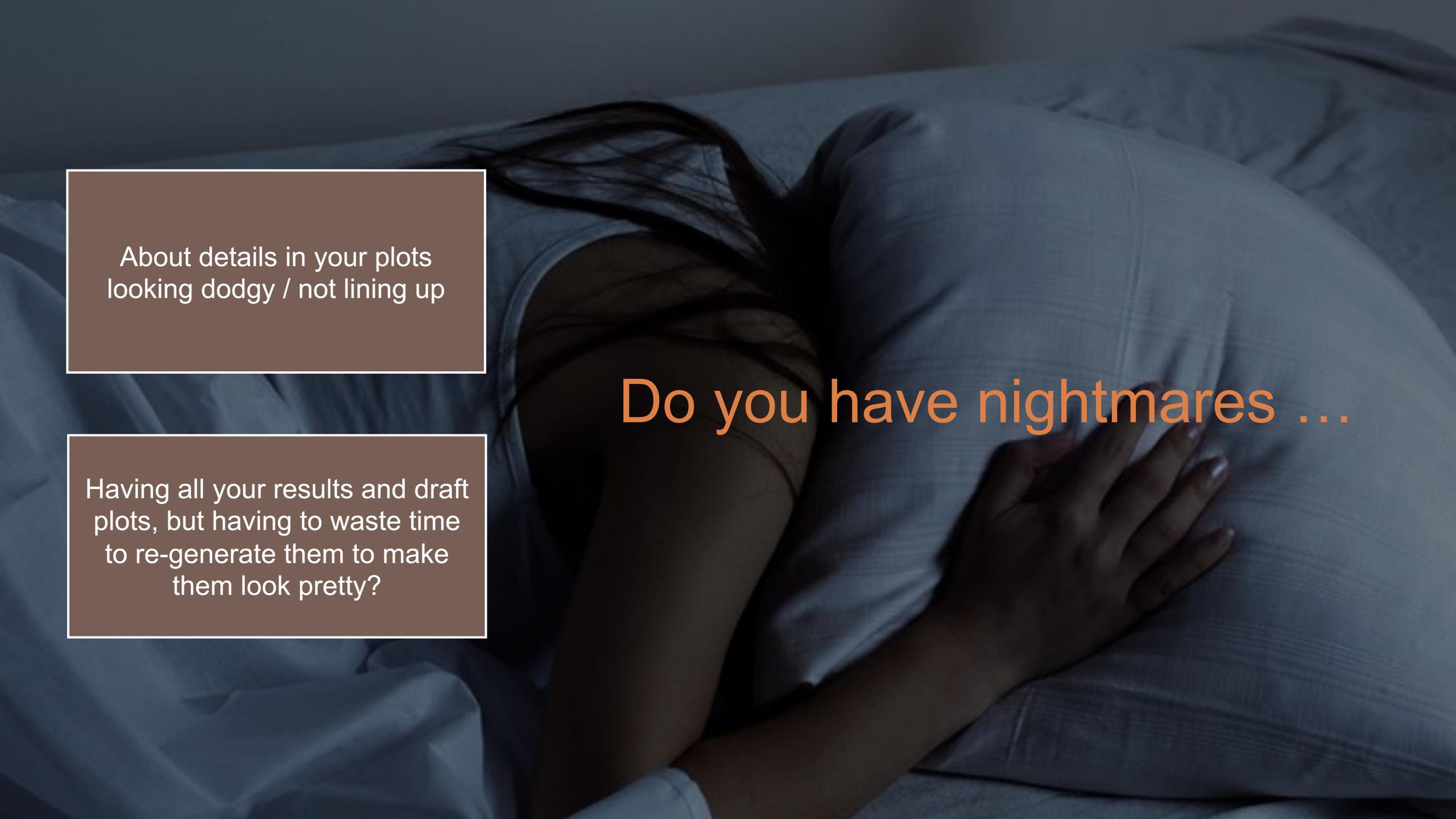
- Why?
 - Reproducible research (runnable)
 - Re-run analysis with a single click
 - Easier to present results
 - “Centralised” results → easier to pull into manuscript
 - Much wow, little effort
 - ... or amazing things with more effort
 - Supplementary Materials
 - “Website” of results
 - Interactivity, etc.



Results output: RMarkdown



- “Keeping manuscript in mind”
 - Set up
 - QC .Rmd → exploratory data analysis, separate so doesn’t clutter results
 - analysis .Rmd → results that are more geared towards the manuscript
- HPC users
 - I can never get RStudio to reliably work on HPC ...
 - Solutions:
 - GitHub on HPC + local → update HPC version
 - Write code locally as an .Rmd and test interactively on HPC → generate .Rmd html by copying across .Rmd from local → HPC

A dark, moody photograph of a person sleeping in bed. The person's head is turned away from the camera, and their hand is covering their face, conveying a sense of distress or噩梦 (nightmares).

About details in your plots
looking dodgy / not lining up

Having all your results and draft
plots, but having to waste time
to re-generate them to make
them look pretty?

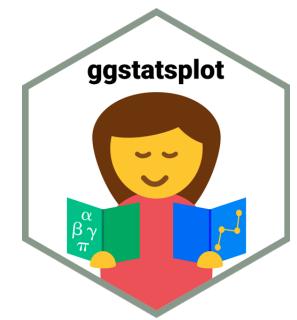
Do you have nightmares ...

Plots

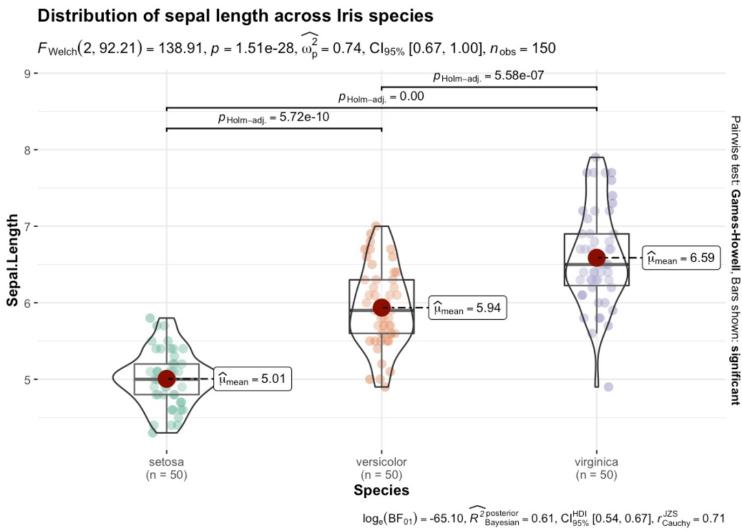
- Preferred packages
 - ggplot
 - ggstatsplot
 - ggrepel for labels
- Design preferences
 - Colour palette: Dark2, various wesanderson palettes
 - Choose early (so you don't have to go back and spend time re-making plots)
 - Keep it simple
 - Keep colour scheme consistent across figures
 - Remove default grey plot background
 - Export as .svg → fiddly editing on Inkscape (or Illustrator)

ggstatsplot

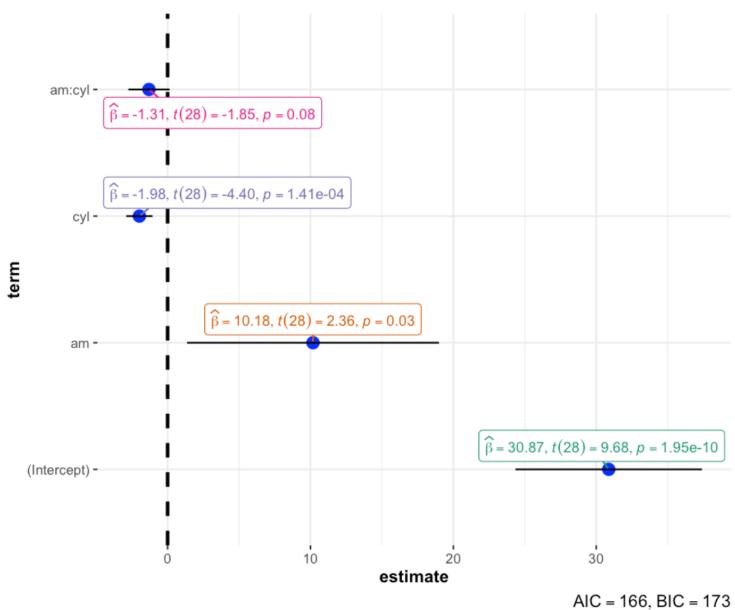
<https://indrajeetpatil.github.io/ggstatsplot/>



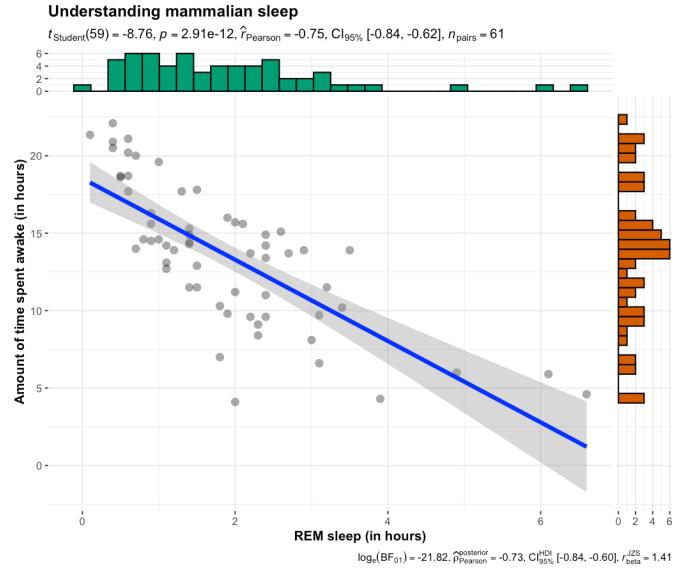
```
ggbetweenstats(  
  data = iris,  
  x = Species,  
  y = Sepal.Length,  
  title = "Distribution of sepal length across Iris species"  
)
```



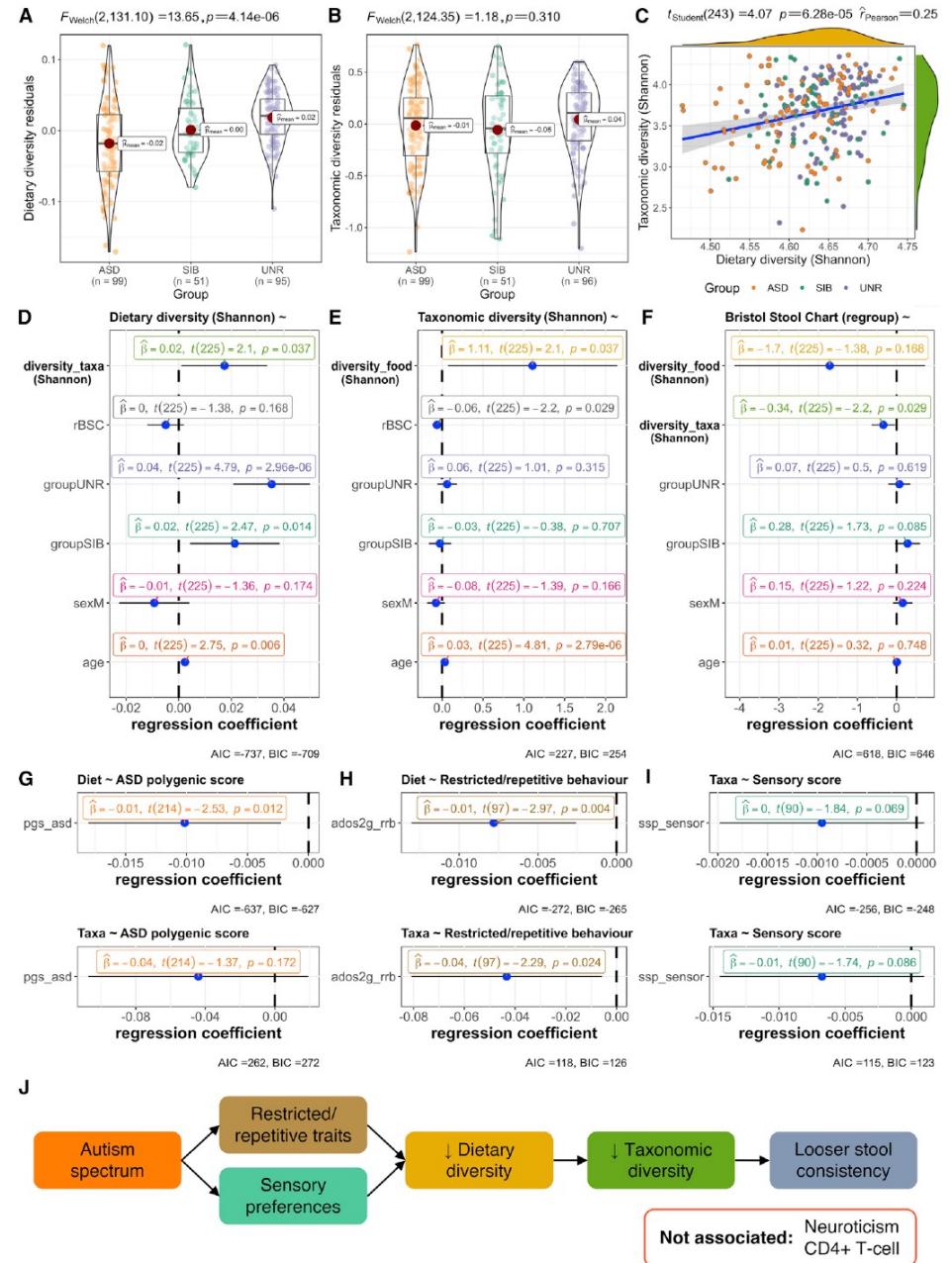
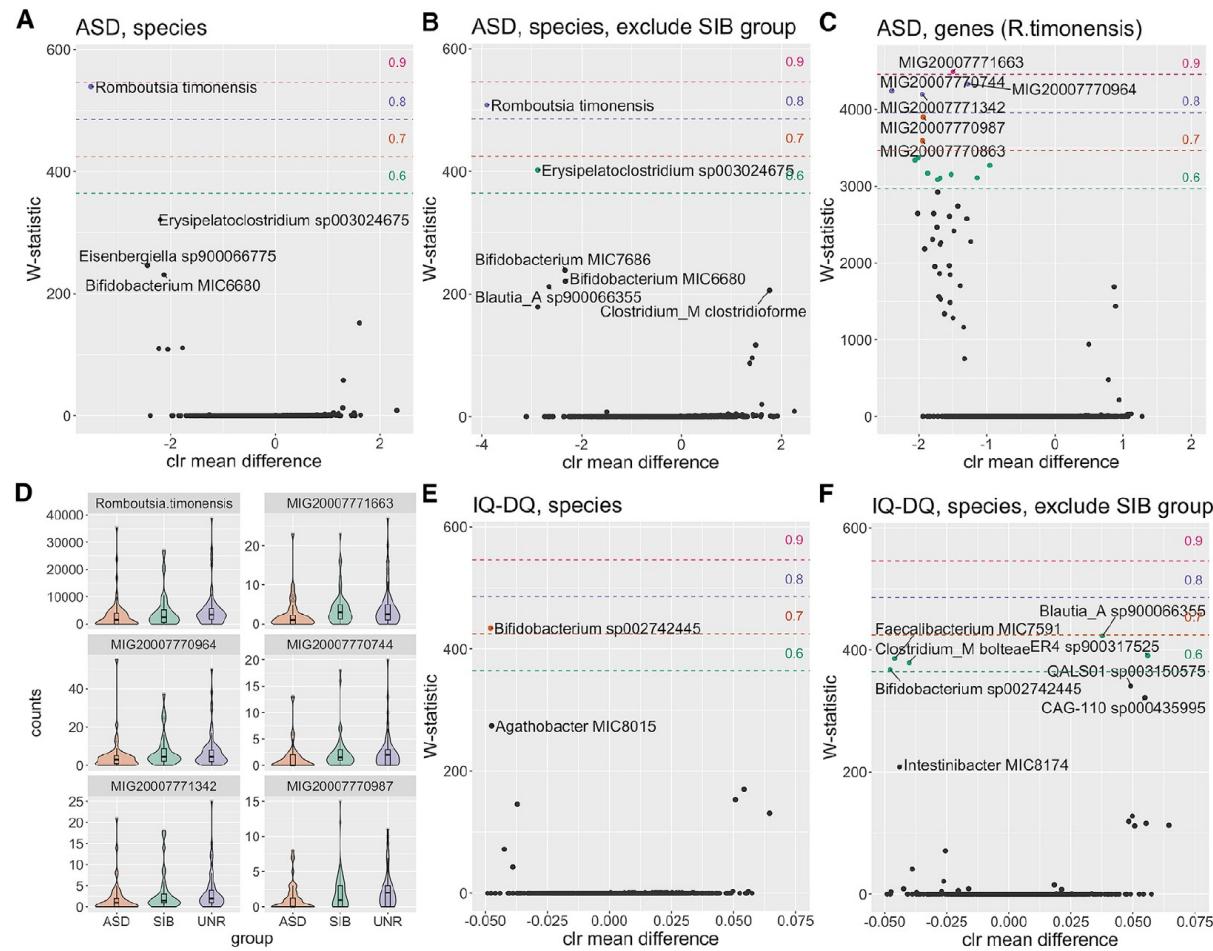
```
## model  
mod <- stats::lm(formula = mpg ~ am * cyl, data = mtcars)  
  
ggcoefstats(mod)
```



```
ggscatterstats(  
  data = ggplot2::msleep,  
  x = sleep_rem,  
  y = awake,  
  xlab = "REM sleep (in hours)",  
  ylab = "Amount of time spent awake (in hours)",  
  title = "Understanding mammalian sleep"  
)
```



Colour coding



.svg → Inkscape

What is Inkscape?

- Vector graphics software, like Illustrator
 - “infinite zoom”
- Free and open-source! (unlike Illustrator)

Use

- Complicated figures
 - Eg. R figures + drawings/illustrations
- Figure editing that is fiddly/annoying in R
 - Eg. overlapping labels, annotations



Handy resources

- StackOverflow
- R Bloggers
- STHDA
- R Graph Gallery
- R Memes for Statistical Fiends (light content)
- swirl, etc.
- Noble 2009 *PLOS Computational Biology A Quick Guide to Organizing Computational Biology Projects*

<https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1000424>

