Reproducible science + 10 recommendations

Brice Ozenne^{1,2} - brice.mh.ozenne@gmail.com

 1 Section of Biostatistics, Department of Public Health, University of Copenhagen

² Neurobiology Research Unit, University Hospital of Copenhagen, Rigshospitalet.

23 March 2023

1 / 20



Good practices



1 / 20



```
1 / 20
```



1 / 20



1 / 20

Why making my analysis reproducible?

Why making my analysis reproducible?

For your future you:

- re-generate results/tables/graphs for review (months later)
- apply the same/similar approach on new projects (years later)

Why making my analysis reproducible?

For your future you:

- re-generate results/tables/graphs for review (months later)
- apply the same/similar approach on new projects (years later)

For collaborators:

- facilitate review by statistician/programmer
- explaining and sharing is facilitated

Why making my analysis reproducible?

For your future you:

- re-generate results/tables/graphs for review (months later)
- apply the same/similar approach on new projects (years later)

For collaborators:

- facilitate review by statistician/programmer
- explaining and sharing is facilitated

For science:

- by 'cleaning' your code you may spot mistakes
- you provide easier access to your methodology

How to make my analysis reproducible The are many ways ... I'll illustrate mine on an example.

	code									
	report									
	results									
۵	Readme.org									
Rea	idme.org									
Contain the R code (folder) code used to generate the results for the cortisol AUC article.										
analysis-AUC-cortisol.R : simulation study										
• data-management.R : data management for the real studies (export processed data in the folder data)										
• re-analysis.R : re-analysis of the real studies using 3 vs. 5 measurements to compute the CAR-AUCi										
R and package versions:										
	sessionInfo()									
R version 4.1.1 (2021-08-10) Platform: x86_64-v6d-ming02/x64 (64-bit) Running under: Windows 10 x64 (bulld 19044)										
	other attached packages:									
	[1] officer_0.4.1 xlsx_	0.6.5	lubridate_1.8.0							
	<pre>[4] plyr_1.8.6 kmlSh</pre>	ape_0.9.5	lattice_0.20-45							
	[7] kml 2.4.1 longi	tudinalData 2.4.1	misc3d 0.9-1							

Good practices

Objective

Create a new folder containing:

- the data
- the programming instructions (R script, SAS script, ...)
- specifics about the program used (version)
- to reproduce results, tables, and figures of an article
 - only those but all of those
 - I usually do that at the end of a project when the structure of the article is 'stable'

▲ avoid unnecessary programming instruction (old analyses)!
 ▲ if you share the folder consider carefully if you share the data or not

Programming instructions

I like to have:

- a data management file
- data analysis files: one per research question or simulation/real data
- one file per table
- one file per figure



Programming instructions

I like to have:

- a data management file
- data analysis files: one per research question or simulation/real data
- one file per table
- one file per figure



Saving intermediate results:

save time when long analyses (avoid recomputation)

Relation to the method section (1/2)

Reading code is often harder than reading english

- people use different programming languages
- require familiarity with the tools used

The method section should (still) fully describe:

- data manipulation during data management (exclusion, transformation, imputation, ...)
- statistical models and tests

Check the correspondence between your (clean) code and what you describe!

Relation to the method section (2/2)

Saying that your analysis has been done in R 2.4.0 does not help much for reproducing the results.

Sharing your code (and if possible data), with a description of your software version (e.g. sessionInfo) is useful.

Relation to the method section (2/2)

Saying that your analysis has been done in R 2.4.0 does not help much for reproducing the results.

Sharing your code (and if possible data), with a description of your software version (e.g. sessionInfo) is useful.

Side note: "A Cox model was used to relate the survival time of patients to their vaccination and age group." miss one key information

• what is the parameter of interest: hazard ratio, difference in 1 year risk, ...

Challenges

Making the analysis reproducible is not really time consuming

- you spend some time now but save time later. It is mostly copy pasting.
- it is a good sanity check

Making the code readable can be time consuming

- commenting the code, having meaningful names, appropriate spacing, avoid long lines ¹
- bugs can be created by renaming, or modifying lines of code. So better to start clean than to fix afterwards.

Good practexample of coding style https://style.tidyverse.org/

Some recommandations

Reproducible science

1. Start with a plan

- First define the research questions(s)
- Then look at the data

Why?

- save time: avoid to get lost in the possibilites
- validity: keep track of how many "tests" have been performed
- avoid temptation: hard to think "independently" after having looked at the data

2. Respect time-varying exposures



Immortal time bias when comparing never-switchers to switchers.

3. Use matching with care

Matching can be used to:

- (try to) adjust for unobserved confounders in an observational study.
- efficiently sample based on important factors

Data analysis must take account of the matching:

- more complex to understand and carry out (e.g. individual matching: conditional logistic regression)
- some coefficients cannot be interpreted (e.g. frequence matching on age: age coefficient)
 - not always clear what one really adjust on
 - ⚠ risk of overmatching
 - e.g. twin study relating diet to obesity

4. Inspect your data before fitting models

Table:

set								exposure							
	1	3	4	5	7	8	9	10	11	12	14	16	17	0	1
control	2	2	1	3	3	2	3	2	3	3	3	3	3	46	0
case	1	1	1	1	1	1	1	0	1	1	1	1	1	15	0

Graph: Histogram, Lexis diagram, Kaplan Meier curves, ...



5. Be concerned of what you don't observe

Drop-out / censoring:

(I could have observed it)

- completely at random?
- related to the outcome? to the exposure?
- \rightarrow investigate possible cause of drop-out/missingness?

Competing risks: (it cannot happen anymore)

- Specialized tools (e.g. cause specific methods)
- the treatment should not reduce the risk of a disease by killing people

Reproducible science

10 recommandations

6. Model specification matters (1/2)

The name of a model often matters less than how you use it:

```
## 1
glm(event \sim exposure + age,
    family = poisson, offset = log(time), data = dSplit)
## 2
glm(event \sim period + exposure + age,
    family = poisson, offset = log(time), data = dSplit)
## 2.bis
coxph(Surv(event, time) \sim exposure + age,
      data = d
## 3
glm(event \sim period * exposure + age,
    family = poisson, offset = log(time), data = dSplit)
```

• "We used a Poisson model to ..." is a rather vague statement

6. Model specification matters (2/2)

When fitting a Poisson regression

... remember the offset

- and the log-transform
- ... remember to specify the family argument
 - otherwise glm will perform a "standard linear regression"
- ... remember to properly model the time effect (Lexis macro)
 - otherwise too unflexible model

7 Vizualize the fitted model

In each strata or for specific individuals:

- Logisitic model: fitted prevalence/risk
- Poisson model: fitted rate/survival
- Cox model: fitted survival

Is that what you "wanted"? Does that match the observed data?



Good practices

8. Don't get trapped by the software

When possible, do examples "by hand"

softwares do calculation for you (+,-,*,/), not magic!

Labels of regression coefficients may be misleading

• in presence of interactions

	coef	<pre>exp(coef)</pre>	<pre>se(coef)</pre>	Z	Pr(z)	
bcg	-0.56991	0.56557	0.20445	-2.788	0.00531	**
dtp	0.17430	1.19041	0.72170	0.242	0.80916	
bcg:dtp	0.21183	1.23594	0.74364	0.285	0.77576	

Regression coefficients are tools:

• not (necessarily) parameters of interest

	Estimate	Std. Error	Z	Pr(z)	
exp(bcg)	0.565582	0.115605	-2.788	0.0157	*
exp(dtp)	1.190413	0.859121	0.242	0.9926	
<pre>exp(bcg + dtp + bcg:dtp)</pre>	0.832102	0.145368	-1.052	0.6362	
Good Adjusted p values reporte	18 / 20				

9. Absence of evidence is not evidence of absence

To "accept" the null hypothesis:

- look at the confidence intervals
- not at p-values!



9. Absence of evidence is not evidence of absence

To "accept" the null hypothesis:

- look at the confidence intervals
- not at p-values!



10. Know when to ask for help

Not too early:

• you know what you want

Not too late:

• you are keen on implementing suggestions!

Try to write down as precisely as possible:

your question
 what you have tried to address it

Statistical advisory service at KU:

- over the phone consultancy
- 20 minute meeting

statistical café