Reproducible Research with R

D G Rossiter

Section of Soil & Crop Sciences, Cornell University ISRIC-World Soil Information 南京师范大学地理科学学院 中国科学院南京土壤研究所

November 20, 2019

Outline

- 1 Reproducible Research
- 2 The R Project
- **3** Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem
- 6 Final words

3

Outline

1 Reproducible Research

- 2 The R Project
- 3 Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem
- 6 Final words

3

・ロン ・四 と ・ ヨ と ・ ヨ と …

Reproducible research 可再现的研究

Definition: "Research papers with accompanying software tools that allow the reader to directly reproduce the results and employ the computational methods that are presented in the research paper." 1

¹Gentleman, R., & Lang, D. T. (2007). Statistical analyses and reproducible research. Journal of Computational and Graphical Statistics, 16(1), 1–23. https://doi.org/10.1198/106186007X178663 < > < > < > <

イロト 不得 とくほと くほとう

Why reproducible research?

- **Readers** can both **verify** and **adapt** the claims in the paper
 - no need to trust that the author has performed computations correctly
 - no hidden assumptions
- Authors can reproduce the results in the future e.g., when a study is extended to a new area
- Authors can adapt the document to a different medium journal paper, web page, tutorial . . .
- In-text computations, figures, tables, can be recalculated if the data changes
 - e.g., reading in a different climate dataset each day and generating a report

D G Rossiter

Elements of reproducible research

- **1** Datasets (tabular, geographic, time-series ...)
- 2 Computer code
- **3 Text** explaining processing
- 4 Text showing results generated by the computer code
- **5** Graphics showing results generated by the computer code
- 6 Text discussing results

Literate Data /

Example reproducible research document



daily mean PM10 concentrations

Figure 7: Interpolated maps for daily PM10 concentration from May 1 to 9, 2005.

https://www.researchgate.net/publication/236011349_Spatio-temporal_ analysis_and_interpolation_of_PM_10_measurements_in_Europe < I > < I > <

Tools for reproducible research

1 A programmable **computing language**

1 e.g., the **R** project for statistical computing; **Python**

2 A literate programming language to mix code, output and text

e.g., Markdown², knitr³

3 A good user interface to use these

e.g., **R Studio**, Jupyter Notebook⁴ for Python

```
<sup>2</sup>https://www.markdownguide.org
<sup>3</sup>https://yihui.name/knitr/
<sup>4</sup>https://jupyter.org
```

D G Rossiter

Outline

- 1 Reproducible Research
- 2 The R Project
- 3 Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem
- 6 Final words

3

イロン イロン イヨン イヨン

- R is an open-source environment for statistical computing, data manipulation and visualisation;
- Statisticians have contributed over 15 000 specialised statistical procedures as packages;
- R and its packages are freely-available over the internet;
- R runs on Microsoft Windows, Unix[©] and derivatives Mac OS X and Linux;
- **R** is **fully programmable**, with modern computer language, **S**;
- Automation by user-written scripts, functions or packages;

...

Why R? (2)

...

- R has comprehensive technical documentation, user-contributed tutorials and textbooks showing R code;
- R syntax for model formulas are a standard, found in many documents;
- R can import and export in MS-Excel, text, fixed and delineated formats (e.g., CSV), databases (e.g., SQLite), vector and raster geographic coverages ...;
- R is fully supported in the R Markdown literate programming environment.

D G Rossiter

3

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 うの()

- sp, sf structures for spatial objects, spatial overlay, subsetting, sampling ...
- rgdal, proj4 coordinate reference systems, projections rgeos vector GIS operations
- raster manipulate spatial data in **raster format** gstat, fields,RandomFields, geoR **model-based geostatistics** (also Bayesian)
- spcosa, spsann spatial sampling, simulated annealing (Brus) spdep spatial dependence on lattice data (like Geoda) spatstat Spatial Point Pattern Analysis
- randomForest, ranger, Cubist, caret, nnet ... machine learning

Image: A math a math

R packages for pedology (soil science)

aqp Algorithms for Quantitative Pedology soiIDB access to USA soils databases

FreeFe 0 • 5 • 10 • 15 • 20 • 25 • 30



D G Rossiter

Reproducible Research with R

Outline

- 1 Reproducible Research
- 2 The R Project

3 Some simple R

- 4 Literate Data Analysis
- 5 The R ecosystem

6 Final words

・ロン ・四 と ・ ヨ と ・ ヨ と

R can be used as an **interactive calculator** at the > command prompt, with all the ususal operators.

> 2*pi/360 # degrees to radians
[1] 0.01745329

Any expression can be **saved** as an **object** in the **workspace** with the **assignment operator** <- and then re-used:

```
> deg2rad <- 2*pi/360
> 30*deg2rad
[1] 0.5235988
```

D G Rossiter

3

・ロト ・回ト ・ヨト ・ヨト

Functions

- \blacksquare Most computation in R depends on functions: arguments \rightarrow function \rightarrow results
- Packages all define functions, each function has some help (information on how to use) and examples
- Some simple functions: mean ,var, range
- Some complex functions: lm, gls, krige,
- Some graphics functions: plot, ggplot, hist
- Some functions change their behaviour depending on the class of their argument, e.g. summary, predict

Examples of functions



128 random numbers

Dataframes

A matrix with named columns (fields, attributes) and (optionally) named **rows** (cases, tuples).

```
> data(trees); class(trees) # a built-in example dataset
[1] "data.frame"
> dim(trees); colnames(trees)
[1] 31 3
[1] "Girth" "Height" "Volume"
> trees[1:3,] # access rows by row number
 Girth Height Volume
1 8.3
          70 10.3
2 8.6 65 10.3
3 8.8 63 10.2
> trees[. 1] # access columns by column number
 [1] 8.3 8.6 8.8 10.5 10.7 10.8 11.0 11.0 11.1 11.2 11.3 11.4 11.4 11.7 12.
[17] 12.9 13.3 13.7 13.8 14.0 14.2 14.5 16.0 16.3 17.3 17.5 17.9 18.0 18.0 20.
> trees[which.max(trees$Volume), ] # find largest tree, note $ operator
  Girth Height Volume
31 20.6
            87
                   77
                                              ◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 うの()
```

Vectorized computations

Almost never need to write **for** loops; most operations work over vectors **element-wise**.

```
> length(trees$Girth); length(trees$Height)
[1] 31
[1] 31
> summaryt(trees$Girth) # circumfrence inches, see metadata
    Min. 1st Qu. Median Mean 3rd Qu. Max.
    8.30 11.05 12.90 13.25 15.25 20.60
> summary((trees$Girth*2.54)/pi) # convert to diameter cm
    Min. 1st Qu. Median Mean 3rd Qu. Max.
    6.711 8.934 10.430 10.711 12.330 16.655
> summary(trees$Height/trees$Girth) # thinness/thickness index
    Min. 1st Qu. Median Mean 3rd Qu. Max.
    4.223 4.705 6.000 5.986 6.838 8.434
```

D G Rossiter

I naa

Model formulas

A compact way to describe the components of a model (linear, random forest, time series \dots)

- "depends on"
- + additive
- interaction
- / nested
- remove a term

Examples:

log(Zn) ~ flood.frequency + distance*elevation
weight ~ grade/height
(these names must be defined in the data used for the model)

3

> data(trees)

> summary(trees)

Girth	Hei	ght V	Volume		
Min. : 8.3	30 Min.	:63 Min.	:10.20		
1st Qu.:11.0)5 1st Qu.	:72 1st (Qu.:19.40		
Median :12.9	0 Median	:76 Media	an :24.20		
Mean :13.2	25 Mean	:76 Mean	:30.17		
3rd Qu.:15.2	25 3rd Qu.	:80 3rd (Qu.:37.30		
Max. :20.6	30 Max.	:87 Max.	:77.00		
> summary(mod	lel <-lm(Vol	ume ~ Girtl	h*Height, data=	trees))	
Residuals:			-		
Min	1Q Median	ЗQ	Max		
-6.5821 -1.06	0.3026	1.5641 4	.6649		
Coefficients:					
	Estimate St	d. Error t	value Pr(> t)		
(Intercept)	69.39632	23.83575	2.911 0.00713	**	
Girth	-5.85585	1.92134 ·	-3.048 0.00511	**	
Height	-1.29708	0.30984 ·	-4.186 0.00027	***	
Girth:Height	0.13465	0.02438	5.524 7.48e-06	***	
Residual star	dard error:	2.709 on 2	27 degrees of f	reedom	
Multiple R-so	uared: 0.9	756,Adjust	ed R-squared:	0.9728	
F-statistic:	359.3 on 3	and 27 DF,	p-value: < 2.	2e-16	
			•	コントロント・ロント	∃ 2000

Matrices

Various operators and functions are provided for **matrices** (similar to Matlab):

- +, -, *, / etc. work element-wise
- matrix multiplication: %*%, inner and outer vector products
- transposition: t function
- inversion: solve function
- spectral decomposition: eigen function
- Singular Value Decomposition: svd function
- principal components: princomp, prcomp

イロト 不得 トイヨト イヨト 二日

Random numbers

Support for random selection, **simulation**. Density d, distribution function p, quantile function q, random generation r:

*norm *unif *binom *beta *chisq *pois

. . .

э

Graphics









Getis-Ord Gi, Syracuse leukemia incidence



SCS

ъ.

Outline

- 1 Reproducible Research
- 2 The R Project
- 3 Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem
- 6 Final words

3

・ロン ・四 と ・ ヨ と ・ ヨ と

Literate Data Analysis

- Idea: mix text and computations in one document
- text: explain assumptions, choices made in computation, comments on results
 - can include computational results in text automatically
- computation: code that generates results, including text, tables, or graphics

-

Example – input in literate programming environment

🔊 POLA	RIS.Rmd 🔀				
←→	🛲 🚃 💑 🔍 🚀 Knit 👻 🌣 🗸	€C Insert +	† ↓	-	
47 - 48	# Import POLARIS NetCDF files				Introduction Import POLA
49 50	Open a NetCDF file:				Find the pr Series code
51 ▼ 52 53 54 55 56 57 58	<pre>```{r} try chains for the second second</pre>		*	≤ ►	Extracting se Raster object Value at a p Series classifi Soil Taxono VRT files
59 60	Some information about the file:				
61 ↓ 62 63 64 65 66 67 68 69	<pre>```{r} library(ncdf4.helpers) length(v.list <- nc.get.variable.list(nc)) head(v.list) tail(v.list) tail(v.list) we see there are paired varibles: the class and its probability, ' series), and a final coverage 'original'.</pre>	for 50 class	🌣 Ses (so	z) il	
61:7	C Chunk 3 0				R Markdown 🗧
			<	× 3	

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 ののの

1 Introduction

2 Import POLARIS NetCDF files

2.1 Projection

2.2 Find the probable class

2.3 Series codes and names

3 Extracting series names at a point

4 Series classification

5 Soil Taxonomy for a set of points

6 VRT files

2 Import POLARIS NetCDF files

Open a NetCDF file:

library(ncdf4)

ncname <- "plat4243_lon-77-76_series" ncname <- "plat4243_lon-77-76" ncfname <- paste(ncname, ".nc", sep="") nc <- nc_open(file.path(ncfname))</pre>

Some information about the file:

library(ncdf4.helpers)
length(v.list <- nc.get.variable.list(nc))</pre>

[1] 101

head(v.list)

[1] "maxcl01_class" "maxcl01_prob" "maxcl02_class" "maxcl02_prob"

[5] "maxcl03_class" "maxcl03_prob"

Example – input – with LATEX formulas

- 00	🐨 - 🕞 💮 🥻 🤌 Go to file/function		
Ex_Sp	atialSampling_ModelAnswers ×		_
	🖻 📄 🥙 🔍 🐠 Knit + 🔿 + 🥙		
666 - 667 668 670 671 672 673 674	<pre># Sampling for mapping: K-means sampling using covariates Now we use information from "covariates" that we think are correlated with the target variable. Our aim is to "create" strata based on the values of the covariates. Presumably these combinations of covariates are better linked to the target than using just one covariate. This is a form of "cluster analysis". A popular method is k-means. This is explained in many references, e.g. Ch. 14 of: Nexter, T., Thishironi, R., & Friedman, J. M. (2009). The element of statistical learning data mining, inference, and production (2nd ed). New York: Springer. Retrieved from http://link.springer.com.proxyllbrary.cornell.edu/book/10.1007/251732-0-387-54858-Z The concept is to find clusters that minimize the mithin-cluster variance and maxinize the between-cluster variance. I.e., minimize: \$\$M(C) = \frac{1}{1}{2}.sum_{cl=1}+k \sum_{cl(i)}+k \sum_{c</pre>		Design-based sampling Simple random sampl Stratified random sam Sampling to fit a linear Fully random sample Compare to the know Challenge Sampling for mapping: Example dataset Select a square grid Select square grid witi Sampling for mapping Gasmet data Select square grid witi Sample dataset Select a square grid witi Sample dataset Select square grid witi Sample dataset Select square grid witi Sample data
	$W(C) = \frac{1}{2} \sum_{k=1}^{K} \sum_{C(t)=k} \sum_{C' \leftarrow k} \ x_t - x_t'\ ^2 - \sum_{k=1}^{K} N_k \sum_{C(t)=k} \ x_t - \bar{x}_k\ ^2$	l	Visualizing the sampli Sampling for mapping: Estimation of the vario Construction of the sa Computing the maxim Exploring the sensitiv Sensitivity for nuose
675 676 677 678 -	This is minimized by assigning the SMS observations to the SKS clusters so that theoverope dissimilarity_ of the observations within each cluster from that cluster's mean is minimized, over all clusters. Although this could be solved by exhaustive search, that is not feasible for any reasonably-sized dataset, so there are various "greedy" algorithms, e.g., iterative descent from a storting allocation.		Sensitivity for range Sampling for mapping: Reading the data Optimizing the sample The simulated annea Optimization Result
680 681	We use a study area from the Hunter Valley of New South Wales (AU), famous for its wines. The soil geography of this area has been extensively studied by soil scientists from the University of Sydney. See for example:		
682	Hang, J., McBratney, A. B., Molone, B. P., & Field, D. J. (2018). Mapping the transition from pre-European settlement to contemporary soil conditions in the Lower Hanter Valley, Australia. Geoderma, 329, 27-42. https://doi.org/10.1016/j.gooderma.2018.05.016		
683 684 685	Load the file with discretisation of the study area and summarize it:		
686 +	····		

1 Design-based sampling for population statistics

2 Sampling to fit a linear regression

3 Sampling for mapping; grid sampling

4 Sampling for mapping: K-means sampling using covariates

4.1 Sample data

4.2 Construct clusters

4.3 Select sample

4.4 Visualizing the sampling plan in Google Earth.

5 Sampling for mapping: Modelbased sampling for OK and KED

6 Sampling for mapping: Conditional Latin Hypercube sampling

4 Sampling for mapping: K-means sampling using covariates

Now we use information from covariates that we think are correlated with the target variable. Our aim is to create strata based on the values of the covariates. Presumably these combinations of covariates are better linked to the target than using just one covariate. This is a form of *cluster analysis*. A popular method is k-means. This is explained in many references, e.g. Ch. 14 of:

Hastie, T., Tibshirani, R., & Friedman, J. H. (2009). The elements of statistical learning data mining, inference, and prediction (2nd ed), New York: Springer, Retrieved from http://link.springer.com.proxy.library.cornell.edu/book/10.1007%2F978-0-387-84858-7

The concept is to find clusters that minimize the within-cluster variance and maximize the between-cluster variance. I.e., minimize:

$$W(C) = \frac{1}{2} \sum_{k=1}^{K} \sum_{C(i)=k} \sum_{Ci'=k} ||x_i - x_i'||^2 - \sum_{k=1}^{K} N_k \sum_{C(i)=k} ||x_i - \bar{x}_k||^2$$

This is minimized by assigning the N observations to the K clusters so that the *average dissimilarity* of the observations within each cluster from that cluster's mean is minimized, over all clusters. Although this could be solved by exhaustive search, that is not feasible for any reasonably-sized dataset, so there are various "greedy" algorithms, e.g., iterative descent from a starting allocation.

4.1 Sample data

We use a study area from the Hunter Valley of New South Wales (AU), famous for its wines. The soil geography of this area has been extensively studied by soil scientists from the University of Sydney. See for example:

Huang, J., McBratney, A. B., Malone, B. P., & Field, D. J. (2018). Mapping the transition from pre-European settlement to contemporary soil conditions in the Lower Hunter Valley, Australia. Geoderma, 329, 27-42. https://doi.org/10.1016/j.geoderma.2018.05.016

Load the file with discretisation of the study area and summarize it:

```
load("HunterValley4Practicals.RData")
# loaded object is `grdHunterValley'
summary(grdHunterValley)
```

Easting Northing elevation m slope deg

э

R Markdown

- A simple formatting language to mix code and text
- can be **compiled** to HTML, PDF, Word, websites
 - code is executed, output including graphics are written to the document, regular text is formatted
- \blacksquare code "chunks" are written inside a { ' ' ' r} ... { ' ' '} block
- outside of these blocks are regular text
- computations can be included in the regular text with the 'r ...' syntax
- graphics commands automatically produce graphs in the output document

D G Rossiter

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 ののの

e simple R

Literate Data Analysis

Running code chunks and compiling



Running code interactively



Compiling to a document

D G Rossiter

Reproducible Research with R

R Markdown example

(see demonstration)

▲口> ▲□> ▲目> ▲目> 三日 のへの

D G Rossiter

Reproducible Research with R

Outline

- 1 Reproducible Research
- 2 The R Project
- 3 Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem
- 6 Final words

D G Rossiter

Reproducible Research with R

3

イロン イロン イヨン イヨン

R ecosystem

- the R project⁵ manuals, FAQ, tutorials
 - CRAN⁶ base R and packages
 - task views⁷ compare packages for different applications
 - These are **mirrored** at many sites over the world.
- R Studio⁸ environment to use R and R Markdown
- R Markdown⁹ literate data analysis

```
<sup>5</sup>https://www.r-project.org
<sup>6</sup>https://mirrors.tuna.tsinghua.edu.cn/CRAN/index.html
<sup>7</sup>https://mirrors.tuna.tsinghua.edu.cn/CRAN/web/views/
<sup>8</sup>https://www.rstudio.com
<sup>9</sup>https://rmarkdown.rstudio.com
<sup>9</sup>https://rmarkdown.rstu
```

Learning R

- Introductions and tutorials
 - my R tutorials and applications¹⁰
- On-line help (within R and on the Internet)
- Contributed documentation
- Textbooks
- Task views
- R Journal, Mailing lists, user's conference

¹⁰http:

Tutorials



RStu	dio Education R for Data Science
R for Data Science	E Q A C R for Data Science
Welcome	
1 Introduction	
I Explore	R for Data Science
2 Introduction	Garrett Grolemund
3 Data visualisation	Hadley Wickham
4 Workflow: basics	
5 Data transformation	Welcome
6 Workflow: scripts	Welcome
7 Exploratory Data Analysis	This is the website for "R for Data Science". This book will teach
Workflow: projects	you how to do data science with R: You'll learn how to get your
II Wrangle	data into R, get it into the most useful structure, transform it, visualise it and model it. In this book, you will find a practicum of
Introduction	skills for data science. Just as a chemist learns how to clean test
10 Tibbles	tubes and stock a lab, you'll learn how to clean data and draw
11 Data import	plots-and many other things besides. These are the skills that allow data science to bappen, and here you will find the best
12 Tidy data	practices for doing each of these things with R. You'll learn how to
13 Relational data	use the grammar of graphics, literate programming, and
14 Strings	reproducible research to save time. You'll also learn how to manage cognitive resources to facilitate discoveries when
15 Factors	wrangling, visualising, and exploring data.
16 Dates and times	This website is (and will always be) free to use and is licensed

DG

Reference sheets



Textbooks

- Dalgaard, P. 2008. Introductory Statistics with R. Springer Verlag.
- Venables, W. N. & Ripley, B. D. 2002. Modern applied statistics with S. New York: Springer-Verlag, 4th edition¹¹
- Fox, J. 2011 An R and S-PLUS Companion to Applied Regression. Newbury Park: Sage.
- James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. An introduction to statistical learning: with applications in R, Springer texts in statistics. Springer, New York.

¹¹http://www.stats.ox.ac.uk/pub/MASS4/

D G Rossiter

-

・ロト ・同ト ・ヨト ・ヨト

UseR! series

Springer is publishing a series 12 of **practical introductions with R code** to topics such as:

- data manipulation
- Bayesian analysis
- spatial data anlysis
- time-series, space-time statistics
- interactive graphics
- machine learning
- applications, e.g., ecology, biomedicine, chemometrics, forestry . . .

・ロン ・回 と ・ ヨ ・ ・ ヨ ・ ・



Roger S. Bivand - Edzer J. Pebesma Virgilio Gómez-Rubio	Authors (view affiliations)
Applied Spatial	Roger S. Bivand, Edzer J. Pebesma, Virgilio Gómez-Rubio
Data Analysis with R	Book
	8 14 1.7k 4 125k
	Citations Mentions Readers Reviews Downloads
	Part of the <u>Use R!</u> book series (USE R)
Download book PDF	<u>+</u>

D G Rossiter

Reproducible Research with R

SCS

ъ.

Task views

Task Views are a summary by a task maintainer of which packages are best suited for certain tasks. Full list at https://cran.r-project.org/web/views/index.html. Examples:

- Analysis of Spatial Data¹³
- Multivariate Statistics¹⁴
- Environmetrics (Ecological & Environmental Data)¹⁵
- Hydrology¹⁶
- Clustering¹⁷

```
<sup>13</sup>https://cran.r-project.org/web/views/Spatial.html
<sup>14</sup>https://cran.r-project.org/web/views/Multivariate.html
<sup>15</sup>https://cran.r-project.org/web/views/Environmetrics.html
<sup>16</sup>https://cran.r-project.org/web/views/Hydrology.html
<sup>17</sup>https://cran.r-project.org/web/views/Cluster.html < > < > >
                                                                          3
```

iterate Data Anal.

Analysis <u>The</u>

Final words

Bayesian	Bayesian Inference
ChemPhys	Chemometrics and Computational Physics
Clinical Trials	Clinical Trial Design, Monitoring, and Analysis
Cluster	Cluster Analysis & Finite Mixture Models
Databases	Databases with R
DifferentialEquations	Differential Equations
Distributions	Probability Distributions
Econometrics	Econometrics
Environmetrics	Analysis of Ecological and Environmental Data
ExperimentalDesign	Design of Experiments (DoE) & Analysis of Experimental Data
ExtremeValue	Extreme Value Analysis
Finance	Empirical Finance
FunctionalData	Functional Data Analysis
Genetics	Statistical Genetics
Graphics	Graphic Displays & Dynamic Graphics & Graphic Devices & Visualization
HighPerformanceComputing	High-Performance and Parallel Computing with R
Hydrology	Hydrological Data and Modeling
MachineLearning	Machine Learning & Statistical Learning
MedicalImaging	Medical Image Analysis
MetaAnalysis	Meta-Analysis
MissingData	Missing Data
ModelDeployment	Model Deployment with R
Multivariate	Multivariate Statistics
NaturalLanguageProcessing	Natural Language Processing
Numerical Mathematics	Numerical Mathematics
OfficialStatistics	Official Statistics & Survey Methodology
Optimization	Optimization and Mathematical Programming
Pharmacokinetics 8 1 1	Analysis of Pharmacokinetic Data
Phylogenetics 1997	Phylogenetics, Especially Comparative Methods
Psychometrics	Psychometric Models and Methods
ReproducibleResearch	Reproducible Research
Robust	Robust Statistical Methods
SocialSciences	Statistics for the Social Sciences
Spatial	Analysis of Spatial Data
SpatioTemporal	Handling and Analyzing Spatio-Temporal Data
Survival	Survival Analysis
TeachingStatistics	Teaching Statistics
TimeSeries	Time Series Analysis
WebTechnologies	Web Technologies and Services
gR	gRaphical Models in R

Outline

- 1 Reproducible Research
- 2 The R Project
- 3 Some simple R
- 4 Literate Data Analysis
- 5 The R ecosystem



3

・ロン ・四 と ・ ヨ と ・ ヨ と …

イロト 不得 トイヨト イヨト 二日

Never ever do these things!

- Modify an item in a primary database (e.g., Excel sheet) unless it is a clear data entry error, referring to the original field or lab paper sheets
 - modify (with explanation!) in the analysis
- **Remove** any items from a primary database
 - subset appropriately (with explanation!) in the analysis
- **Compute** any composite item in a primary database (e.g., sum of base cations)
 - compute the composite item in the analysis the code shows exactly what was done
- Transform, project, or resample a GIS or raster coverage in an interactive GIS
 - use raster, rgdal, rgeos, sp, sf ... R packages

Do these things!

- Do all your data manipulation and analysis in a reproducible research environment;
- Explain the processing steps, choices and and assumptions in text;
 - note that the code itself gives explanation of exactly how a result was obtained, but not assumptions or choiuces
- Produce graphics within the program these change automatically if the processing changes

End



Amsterdam National (NL) Maritime Museum 荷兰国家海洋博物馆 Verbiest world map 1674 坤舆 全图

<ロ> <同> <同> < 回> < 回>