

BASIC STATISTICS FOR BIOLOGISTS



UNIVERSITÄT
LEIPZIG

Erik Kusch

erik.kusch@i-solution.de

Section for Ecoinformatics & Biodiversity

Center for Biodiversity and Dynamics in a Changing World (BIOCHANGE)
Aarhus University

- 1** What To Expect
 - The Seminars
 - Course Resources and Reading
- 2** The Importance Of Proper Statistics
 - The Consequences Of Bad Statistics
 - What Are Bad Statistics?
 - Statistical Concern On The Rise
 - Further benefits of a statistical background
- 3** Terminology
 - Classifying Statistics
 - Basic Vocabulary
- 4** Introduction To R
 - Why Use R?
 - The R landscape
 - Layouts
 - Coding

Course Dates & Outline I

Block I - Theory and Basics of R

Date	Time	Topic	Location
I.) Introduction			
Date	Time	(1) An Introduction to Basic Statistics for Biologists	Location
Date	Time	(2) Introduction to R	Location
II.) Basic statistical terminology			
Date	Time	(3) A Primer for Statistical Tests	Location
Date	Time	(4) Descriptive Statistics	Location
Date	Time	(5) Data Visualisation	Location
Date	Time	(6) Inferential Statistics, Hypotheses and our Research Project	Location

Course Dates & Outline II

Block II - Basic Statistics in R

Date	Time	Topic	Location
III.) Handling Data			
Date	Time	(7) Data Handling and Data Mining	Location
IV.) Non-parametric tests			
Date	Time	(8) Nominal Tests	Location
Date	Time	(9) Correlation Tests	Location
Date	Time	(10) Ordinal and Metric Tests for two-sample situations	Location
Date	Time	(11) Ordinal and Metric Tests for more than two-sample	Location
V.) Parametric tests			
Date	Time	(12) Simple Parametric Tests	Location
VI.) Closing			
Date	Time	(13) Summary and an Outlook on Advanced Statistics	Location

Learning Goals

1 A solid grasp of basic biostatistics

- Have an overview of available methods
- Be able to judge the applicability of individual methods

2 Basic proficiency in using R

- Know base commands and how they function
- Be able to prepare biologically relevant data sets for further analysis
- Be able to apply basic statistical methods to biologically relevant data sets

3 Research Design

- Understand how to formulate testable hypotheses
- Know the importance of proper statistical approaches in research
- Being able to critically assess statistical methods in research publications

Learning Methods

We will:

- Cover useful theory of biostatistics (lecture style)
- Run biostatistical analyses in R (seminar style)
- Work through basic biostatistical methods in a research project using simulated data
- Fully reproducible analyses (<https://github.com/ErikKusch/An-Introduction-to-Biostatistics-Using-R>)

When prof shows you how to do analysis on SPSS/R/MATLAB but first shows you by-hand theory so you can get a "conceptual understanding" first



We will focus heavily on actually doing the statistics!

Let Me Introduce Myself

Erik Kusch

Studies:

PhD @ Aarhus University (currently enrolled)

M.Sc. @ University of Bergen

B.Sc. @ Technical University of Dresden

Experience:

Biostatistics Tutor @ University of Leipzig

Biostatistics Research Assistant @ University of Leipzig

Biostatistics Research Assistant @ University of Kyoto

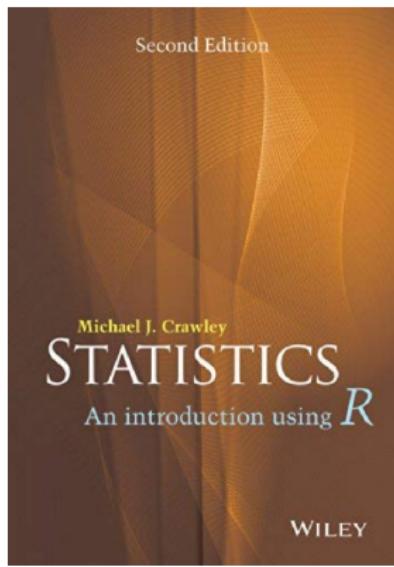
Research:

- Dryland vegetation memory analyses
- Large-scale vegetation-climate modelling
- Remote sensing approaches in macroecology
- Biostatistical approaches in behavioural ecology
- Statistical downscaling of climate reanalysis data for use in biological analyses

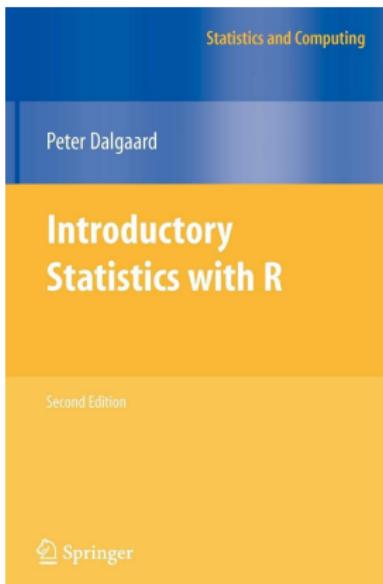


Useful Reading

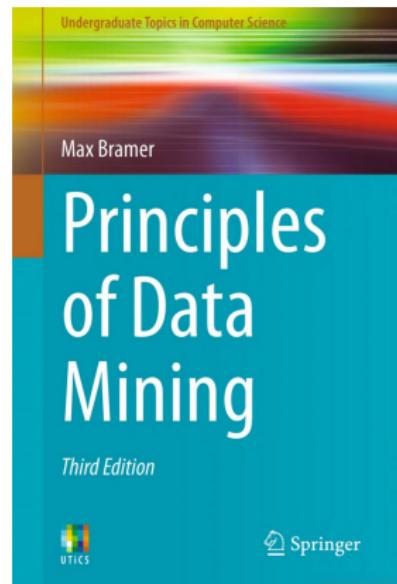
You are **NOT** required to read these!



ISBN: 978-1-118-94109-6



ISBN: 978-0-387-79053-4



ISBN: 978-1-4471-4883-8

But these books are seriously good.

When Mistakes Happen

Even the rigorous peer-review system might miss some minor flaws.

An example:

- Birkenmeyer et. al published a flawed paper in 2016.

The mistake in the data set was spotted by Dr. B. M. Weiß. in early 2017

- A corrigendum was put online
- A corrected version of the paper was uploaded

None of the results of the paper changed.

→ No big deal so long as you offer corrections to your flawed work.

c

Original Article

Sampling t₁ Swabs Sam Volatiles

Claudia S. Birkenmeyer¹,
Markus Küchlein¹

¹Institute of Analytical Chemistry,
University of Leipzig, Linie
Brüder Fuchs 4, 04109 Leipzig,
Germany; www.chemie.uni-leipzig.de/~birkenmeyer
E-mail: Birkenmeyer@chemie.uni-leipzig.de
Fax: +49 341 9442551; Tel:
Correspondence to be sent to:
Leipzig, Germany; e-mail:
Birkenmeyer@chemie.uni-leipzig.de

Abstract

The most common technique for collecting animal breath volatile organic compounds (VOCs) is passive sampling with porous volatile labeled (PVL) analysis. Most critical is the choice of solvent used for extraction. In a recent study, we found that benzene from these materials interfered with the gas chromatography-mass spectrometry (GC-MS) analysis of breath VOCs. Thus, replacing the t₁ swab with a non-porous semipermeable membrane (NPM) did not remove the influence of these solvents on the results.

Key words: body odor, GC-MS, volatile organic compounds

Introduction

In recent years, it has become increasingly evident that substances present in the human breath, either in the excretion or absorption form, can provide important information about the individual's environment (e.g., Polgar et al., 2003; Mennemeyer-Möller et al., 2007).

© The Author(s) 2016. Published by Oxford University Press, on behalf of the International Society for Biostatistics. All rights reserved.

bioRxiv preprint doi: <https://doi.org/10.1101/042260>; this version posted December 20, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

bioRxiv preprint doi: <https://doi.org/10.1101/042260>; this version posted December 20, 2015. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

Original Article

Sampling the Body Odor of Primates: Cotton Swabs Sample Semivolatiles Rather Than Volatiles

Claudia S. Birkenmeyer¹, Ruth Thomassen², Florence Jérin³,
Markus Küchlein¹, Annes Stalmø², Brigitte M. Weiß^{1,4} and Anja Widdowson⁵

¹Research Group of Mass Spectrometry, Institute of Analytical Chemistry, Faculty of Chemistry and Biochemistry, University of Leipzig, Linie Brüder Fuchs 4, 04109 Leipzig, Germany; ²Research Group of Behavioral Ecology, Institute of Biology, University of Leipzig, Linie Brüder Fuchs 4, 04109 Leipzig, Germany; ³Department of Archaeology, University College London, Gower Street, London WC1E 6BT, UK; ⁴Yves-Rémy Research Group of Private Art Selection, Department of Primatology, Max Planck Institute for Evolutionary Anthropology, Döhmenweg 11, 04342 Leipzig, Germany; and ⁵Center for Integrative Bioinformatics (IBZ), Institute for Discrete Mathematics and Geometry, Vienna University of Technology, Vienna, Austria

Correspondence to be sent to: Claudia Birkenmeyer, Institute of Analytical Chemistry, University of Leipzig, Linie Brüder Fuchs 4, 04109 Leipzig, Germany; e-mail: birkenmeyer@chemie.uni-leipzig.de
Accepted 23 March 2016

Abstract

We compared the risks of collecting animal breath volatile organic compounds (VOCs) using porous volatile labeled (PVL) analysis with a non-porous semipermeable membrane (NPM) and found that PVL analysis provides better results. However, given the increasing interest in the use of breath VOCs for medical diagnostics, it is of great importance to understand the influence of the solvent used for extraction on the quality of these measures. We thus compared the t₁ swab with a cotton swab and found that the cotton swab provided more semivolatile compounds, mainly chloro-1,2-diene, than the PVL swab.

Key words: body odor, GC-MS profiling, mass spectrometry, semivolatile organic compounds (VOCs), validation of mass sampling methods

Chemical Science, 2016, Vol. 7(12), 278–285
© The Royal Society of Chemistry 2016
DOI: 10.1039/C6SC00096A
Published online 22 November 2016
Original Article
Advance Access publication 04 October 2016



Fraudulent Practices - The Case Of Andrew Wakefield

Probably one of the most reviled doctors of the 21st century

- Claimed to have found a link for vaccines and autism (Paper from 1998)
- Paper retracted by the publisher
- General Medical Council of Britain revoked his medical license

His academic career is over despite his large community of followers in the U.S., Australia and Brazil.



DAILY REPORT

Early report

Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children

A.J. Wakefield, B.H. Murch, A. Anthony, J. Linton, D.M. Caspary, M. Scott, M. Everett, A.P. Davis, M.A. Thomson, P. Harvey, A. Valentejo, E.E. Davies, J.A. Walker-Smith

Summary

Background. We identified a cluster of cases of **ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children** (ILC-PDD) in children who received oral polio vaccine (OPV).

Introduction

Background. We identified a cluster of cases of ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children (ILC-PDD) in children who received oral polio vaccine (OPV).

Fraudulent Practices - The Case Of Diederik Stapel

Former star in academia, now a laughing stock

- Manipulated data and completely fabricated entire studies
- Fired from his position as professor at Tilburg University
- 58 retracted papers
- Papers of other authors needed to be retracted as well



→ Knowingly fraudulent practices can cost you your career, discredit your institution and your field of research, and even seriously impede the careers of unknowing co-workers.

Wrong/Mal informed Use

Lack of statistical knowledge

- Applying statistics to data which they aren't meant for
→ *Methods can “break”*
- Flawed understanding of the methodology
→ *Incorrect conclusions*

Pure biologists lack knowledge on statistics.

Uninformed Use

Lack of biological knowledge

- Delineation of nonsensical but statistically significant relationships
→ *p-hacking*
- No sense of how to establish testable, feasible hypotheses
→ *Waste of time*

Pure statisticians lack knowledge on biology.

Caveat

- Biologists often have preformed ideas of what to expect
→ data-tweaking to match expectations?
- Researchers also have a vested interest in uncovering extraordinary things
→ The more astounding a paper the better?



FOOLING OURSELVES

HUMANS ARE REMARKABLY GOOD AT SELF-DECEPTION.
BUT GROWING CONCERN ABOUT REPRODUCIBILITY IS DRIVING MANY
RESEARCHERS TO SEEK WAYS TO FIGHT THEIR OWN WORST INSTINCTS.

BY RICHARD NESEK

In 2013, three years after he co-authored a paper showing that humans are remarkably good at self-deception, Richard H. Thaler, a Nobel laureate in economics, was asked to review a paper that had been accepted for publication in a top journal. The paper, which had been submitted by a researcher at the University of California, Berkeley, claimed that people's desire to be right was so strong that it could lead them to ignore evidence that contradicted their beliefs. Thaler, who had previously argued that people's desire to be right was so strong that it could lead them to ignore evidence that contradicted their beliefs, was asked to consider whether the paper's claims were consistent with his own research.

© 2015 Scientific American

ATTENTION!

Don't let a personal bias inform your analysis!

The Recent Debate

- p-values are a cause of concern
 - More on this in seminar 6
(Inferential Statistics and Hypotheses)
- Pre-p-value statistics and data handling increasingly subject of scrutiny
 - More on this in seminar 7 (Data Handling and Data Mining)

more certain of success, Liu promises his students 13 for odds, around the same time giving them a p-value of 0.05, which is a simple process to move on from the p-value to a more serious work for a real-life project if it was, however, generalised if it had been done in a different context. Cai, who says that the Tsinghua students he has taught have a spirit grounded in strict rules and methods and that they probably "there is no way to improve it," adds, "I think both Wu Yiu-Lung and Liu are right. I am not sure how to do the concrete design phase of the years. Tsing students need to know how to do the p-value and Tsing Yiu focus on the core of the research. In my opinion, such discovery could increase their success rate. However, I still think the importance of p-value is not so large. You should stick up to this one," he says. Tsing students have now almost finished a draft of all of previous work, and Liu's own paper's "new generation of a modern statistical methodology for preprint" and "a step-out of the traditional p-value era" will be uploaded with the journal.

(Continued on page 10)

particular could be found to be safe, safe, and effective, and the pharmaceutical company based given the results to the FDA. The FDA then approves the drug and it can be used in the U.S. population. Clinical practice has been based mainly on the experience of the physician, and the technology is, therefore, the physician's responsibility. This helps often finance to a particular drug, and the manufacturer can then claim that the U.S. population will benefit from the drug.

A decision is needed once if there is a problem with a drug. Then, "now is the time to look to change drugs around," he says, to start making decisions.



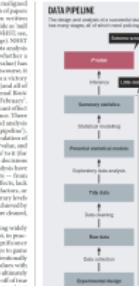
P values,
no

COMMENT

P values are just the tip of the iceberg

Ridding science of shoddy statistics will require scrutiny of every step, not merely the last one, say Jeffrey T Leek and Roger D Peng.

There is no statistic more maligned than P values. P-values are the most common way of assessing whether an intervention or treatment is effective. In a recent study of 10,000 medical papers, we found that 60% of authors used P-values to make claims about the efficacy of their interventions. Yet P-values are often misinterpreted by scientists and misused by journal editors. Given the disconnect, it is important that we move away from P-values and toward a more complete understanding of the science of statistics.



designed to address this crisis. For example, the new course "Statistical Methods for Evidence-Based Medicine" offered by Johns Hopkins University in Baltimore, Maryland, and Duke University in North Carolina, aims to train the relevant software and introduce the concepts of evidence-based medicine. This is a welcome addition to the P-values in the literature on a single topic.

The problem is that the people who do DNA sequencing or proteomics analysis have to learn how to use a machine, not how to interpret the results. They have to learn how to use the relevant software and code, and how to interpret the results. These students should be responsible for their funding and the funding for the course.

There are other courses specifically designed to address this crisis.

The ultimate goal is evidence-based data analysis. This is what evidence-based medicine is, which physicians are encouraged to use only treatments for which there is evidence. The evidence-based medicine and the people they teach and collaborate with are the ones who can interpret P-values, and prevent the rest of the field from doing the same.

Jeffrey T Leek and Roger D Peng are associate professors of biostatistics at the Institute for Computational Medicine at Johns Hopkins University, Baltimore, Maryland, USA.

Editorial: The case for less emphasis on P-values 10 | NATURE | VOL 501 | 20 APRIL 2013

© 2013 Macmillan Publishers Limited. All rights reserved.

Practices in statistics are constantly subject to change.

Why Keep Up With It?

- Journals might enact bans on studies containing p values
→ Counter-productive according to Andrew Vickers (Memorial Sloan Kettering Cancer Center)
- Statistically robust studies hold up to scrutiny much better
→ Statistical prowess enhances your research massively
- Staying up-to-date can help advance one's understanding and career

COMMENT

EXHIBITION Adolf Fleischmann, pathology sculptor and abstract artist. 

INTERVIEW I Can architecturale catalyse creativity at the Click. 

CONVERSATION Legal loophole allows mango farmers to kill their bats in Mauritius. 

BRIEFING Deep-drilling pioneers what were they drinking? 



A tragedy of errors

Mistakes in peer-reviewed papers are easy to find but hard to fix, report David B. Allison and colleagues.

Just how error-prone and self-correcting are we? We have analyzed the past 18 months' papers published in *Nature*.

We are a group of researchers working on obesity, nutrition and energy. In the summer of 2013, we were asked to read a research paper in a well-respected journal estimating how a change in food intake (either weight gain or weight loss) would affect energy expenditure. And we noted that the analysis applied a mathematical model that assumed all energy expenditure was used for metabolism, and others submitted a letter to the editor explaining the problem. Months later, we

were gratified to learn that the authors had decided to withdraw their paper. The fact that they had acknowledged the mistake is startling; this episode was an affirmation that science is self-correcting.

Now, however, in the case is not representative. In the course of assembling weekly lists of articles in our field, we began to notice that many of the peer-reviewed articles containing what we believe were serious errors were not being invalidated, errors,

mistakes or verbiage substantially from clearly erroneous, might also prove problematic.

After attempting to address more than 22 of these errors with letters to authors and editors, we gave up. As time passed and a dozen more, we had to stop — the work took too much of our time. — Our efforts have been rewarded, however, because separately (see 'Three common errors') and showed how journals and authors can be faced with mistakes that need correction.

We learned that post-publication peer review is a useful way to catch errors before they become entrenched in the literature. We hope that journals will take advantage of this opportunity to improve the quality of their publications.

David B. Allison is a professor of nutritional sciences at the University of Tennessee, Knoxville, Tennessee, USA.

✉ nature.com/nature/journal/vaop/npages/1.html

4 FEBRUARY 2016 | VOL 536 | NATURE | 27

Advancing In Statistics

"Treat statistics as a science, and not a recipe!"

~ *Andrew Vickers*

Teaching statistics



Doing Statistics



The Lack Of Biostatisticians

- Biological studies without rigorous statistical analyses are almost unpublishable
- Biostatisticians are rare
- Almost every biological research group requires at least one capable statistician
 - **Biostatisticians are sought-after**

Statistics As An Apphrodisiac

Her: I'm a stats major

Me: [trying to think of something to impress her] yea I'm bad at math too



Frequently Used Classifications

- According to how they are done:
 - Theoretical Statistics
 - Applied Statistics
- According to topic:
 - Biostatistics
 - Economic Statistics
 - Statistical Physics
 - ...
- According to what the goal is and what kind of data is available
 - Regression
 - Classification
- According to how the analyses makes use of the data
 - Supervised
 - Unsupervised

According to the kind of information returned by the methods

- Descriptive Statistics
- Inference/Inferential Statistics

Unsupervised Approaches

Unsupervised methods are often used to select the most informative X input variables for supervised approaches.

Pre-requisites:

- Only *input variables* are observed.
- *No solution/feedback (output)* is given.

Aims:

- *Divide* the observations into relatively distinct groups.
- *Model* the underlying structure or distribution in the data.

→ "Pre-processing" before a supervised learning analysis and exploratory analyses

Supervised Approaches

Supervised methods are often *informed by unsupervised approaches* and used to *gain validated information* about the data.

Pre-requisites:

- Both *predictors X*, and *responses Y* are observed (there is one y_i for each x_i).
- Data is split into *Training* and *Test Data Sets*.

Aims:

- Learn a *mapping function f* from X to Y .
- *Validate* established function/model.
- Further *prediction* and *inference*.

→ **Mostly inferential analyses**

Population vs. Sample

Population: describes the sum total of all *existing* values of a variable given a certain research question. This includes non-measured data.

Sample: describes the sum total of all *available* values of a variable for any given analysis. This can only include measured data.

An example:

In an experimental set-up, you rear an ant colony of exactly 10,000 individuals. You are interested in the average mandible strength of ants within the colony.

The problem: You cannot possibly take measurements of all 10,000 individuals.

The solution: Taking measurements on a **Sample** (e.g. 1,000 individuals) from within the **Population** (10,000 individuals).

Training Data vs. Test Data

This differentiation is only applicable when concerned with *modelling*, which we won't cover in these seminars.

Training Data: describes the subset of the total data which is used to *establish/train* the model.

Test Data: describes the subset of the total data which is used to *test* the performance of the model.

The problem: You have identified a way to model how mandible strength and ant size are interconnected but don't know how to assess the quality of your model (a model will always fit the data it was built on extremely well).

The solution: Split the available data into two non-overlapping subsets of data (**Training** and **Test Data**) and use these separately to build your model and assess its performance.

What Makes Data Truly Random?

Randomisation is one of the **most important** practices in biological studies.

A **sampling** procedure is **random** when any member of the *population* has an equal chance of being selected into the *sample*.

Training and *Test Data Sets* are established from the population with the same sense of randomness although there may be exceptions depending on the modelling procedure at hand.

Data collection: Number all units contained within the set-up and sample those units corresponding to random numbers.

In R: Use the `sample()` function to create truly random subsets. Remember to use `set.seed()` to make this step reproducible!

Random Sampling in R

```
# Making it reproducible
set.seed(42)

# Establishing a population
pop <- c(1:15)
pop

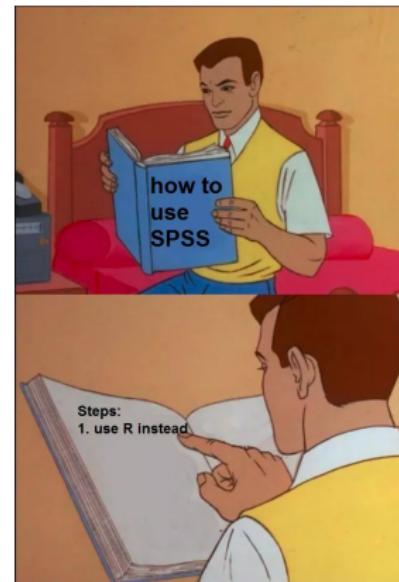
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

# Establishing a random sample
sam <- sample(pop, 5, replace = FALSE)
sam

## [1] 1 5 15 9 10
```

The Power Of R

- 1 **R** is a powerful **statistical** and **graphical tool**
- 2 Available for almost every platform (Windows, Linux, Mac, FreeBSD, etc.)
- 3 It is **completely free**
- 4 **Open source**
 - It can be modified heavily to suit individual demands
 - Constant, moderated user input to widen functionality
 - Dedicated, heavily frequented forums online
 - Allows for reproducible coding



R is the rising star of statistical applications in biological sciences!

Obtaining R

R is a free statistical environment that is used by many researchers all around the globe.

How to get it?

- R is available at
<https://www.r-project.org/>
- A host of editors is available freely on the internet. I recommend RStudio (available at
<https://www.rstudio.com/>).

What if I need help?

- Multiple dedicated forums online:
- <https://stackoverflow.com/>
 - <https://stackexchange.com/>

Layouts - The Console

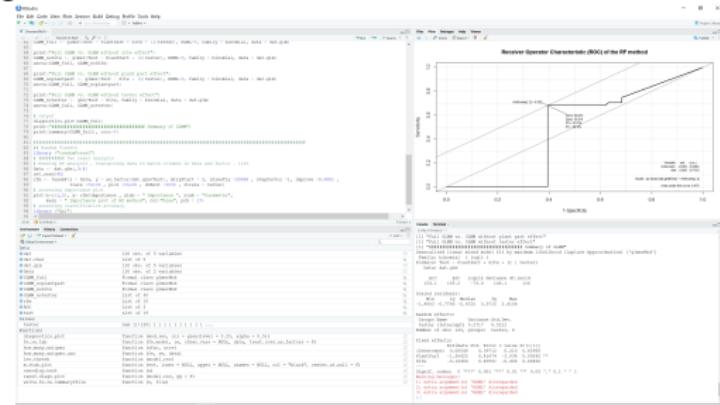
Running R through the console . . .

... is a **bad idea**.

But you will have access to it anyway as it comes with R (we will use version 3.4.2. <https://cran.r-project.org/bin/windows/base.old/3.4.2/>).

Layouts - The Editor

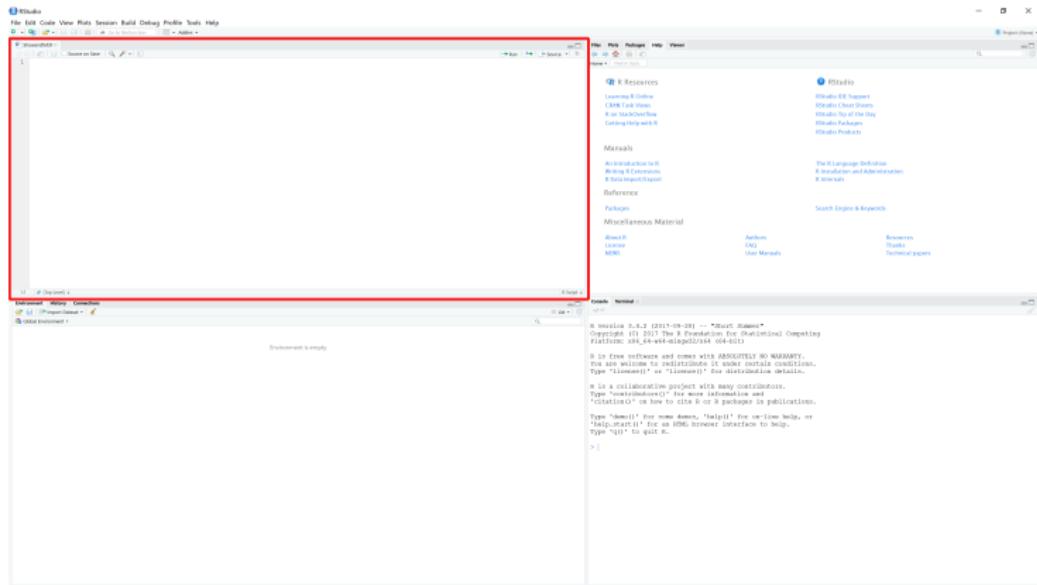
Running R through an editor...



... is a much better idea!

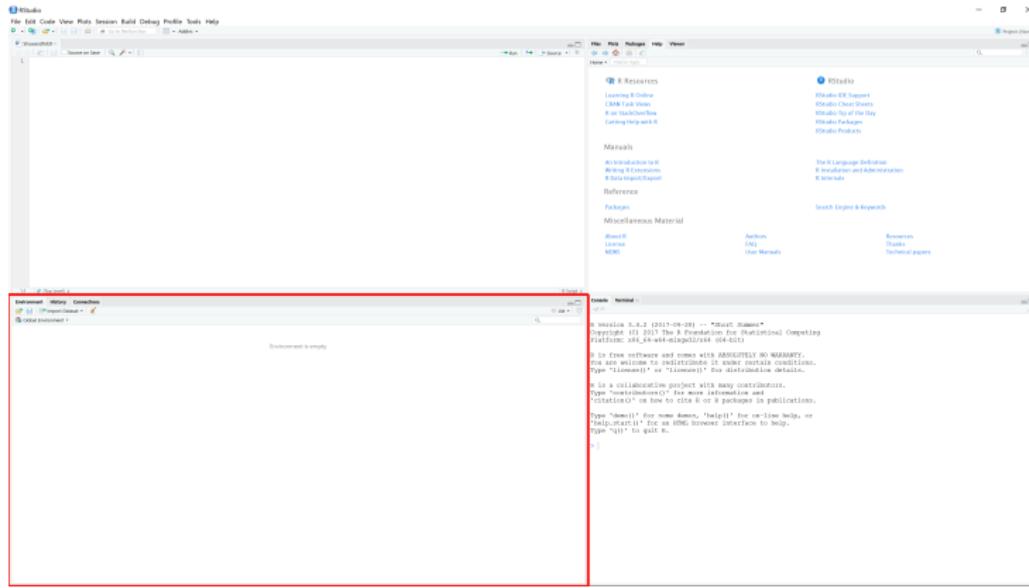
I recommend RStudio (<https://www.rstudio.com/>). If you use it a lot, I also recommend changing the appearance to 'Vibrant Ink' (setting located in the 'Global Options' window nested within the 'Tools' tab).

Layouts - The Editor Explained



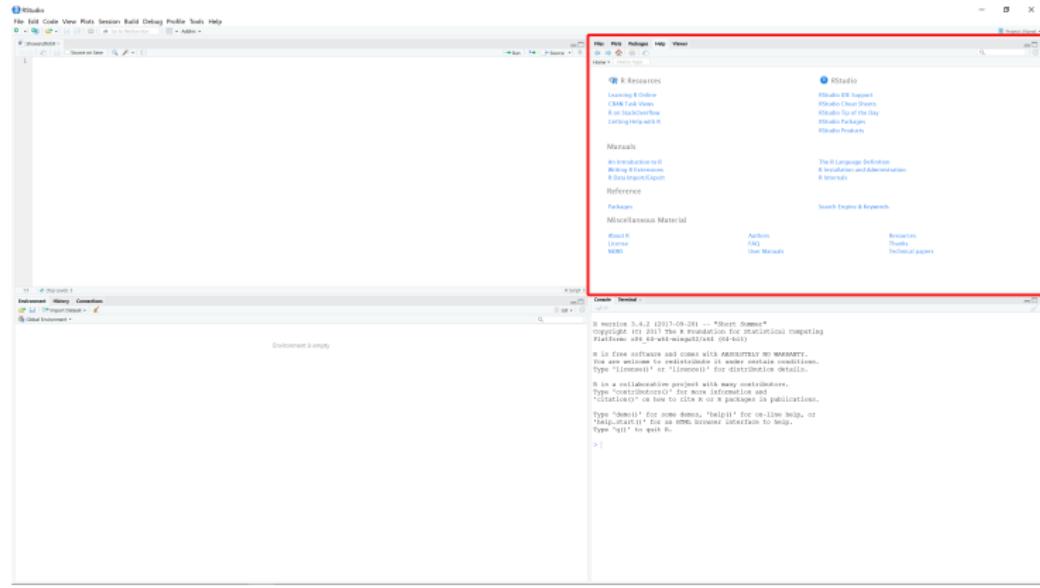
The **Source** is where you load scripts and write most of your coding document.

Layouts - The Editor Explained



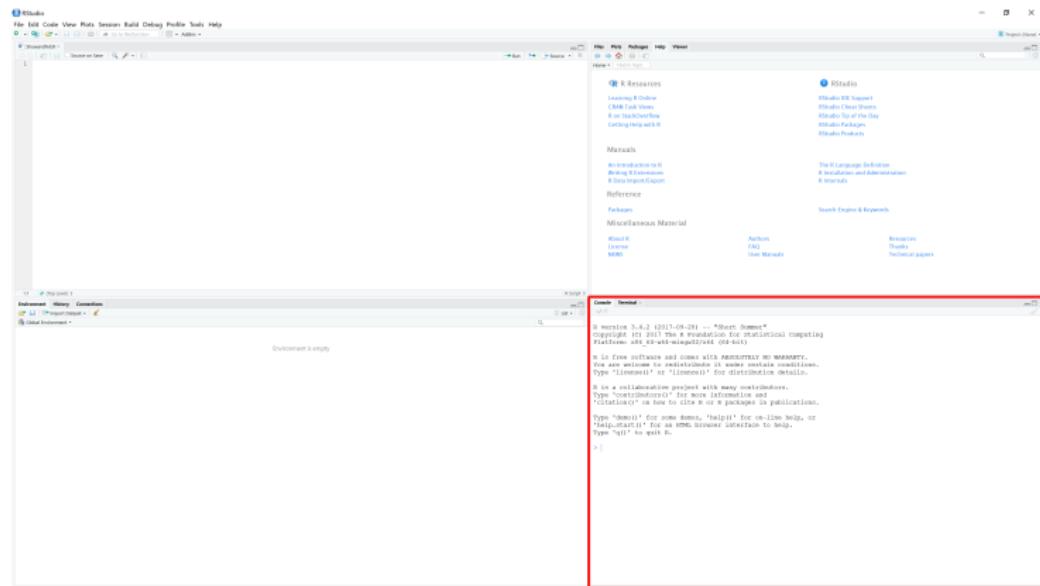
The **Environment, History, Connections** is where you will be able to quickly access all objects of your current R session.

Layouts - The Editor Explained



Files, Plots, Packages, Help Viewer are especially useful for document navigation, data visualisation and to get information on certain functions in R.

Layouts - The Editor Explained



The **Console** is where you execute short commands, and warning and error messages are displayed.

The Evolution Of Code

- Your code and coding practices evolve
- Comment **every line** of code
- Elegant code makes an analysis easier to reproduce
- **Avoid hard-coding!**

"If it looks stupid but it works, it isn't stupid."

