

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

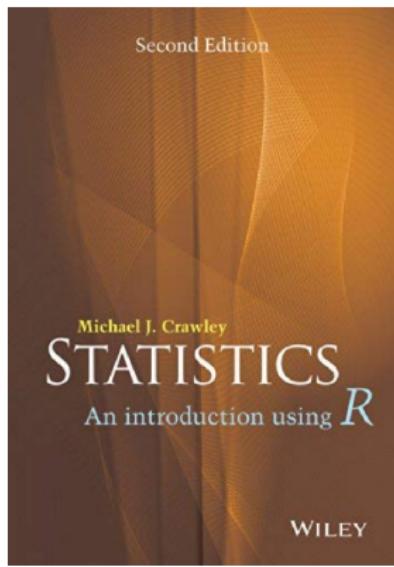
Research:

- 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

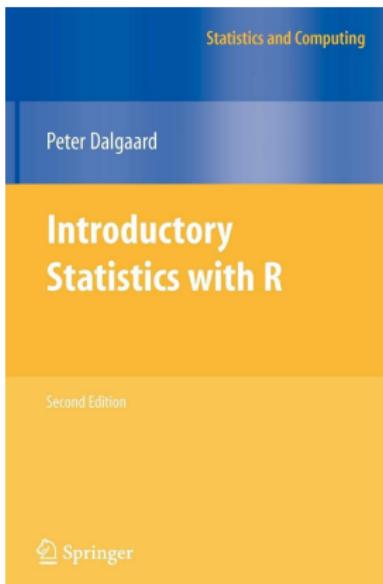


Useful Reading

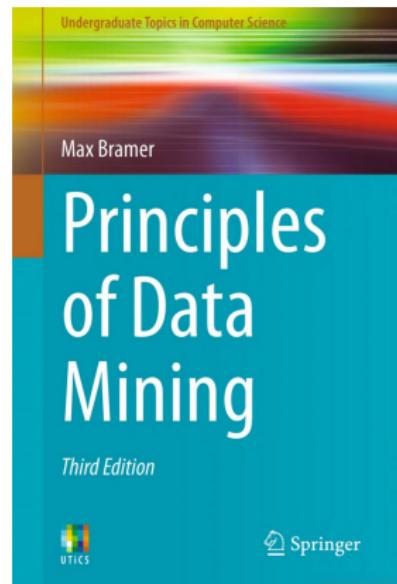
You are **NOT** required to read these!



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. Weiβ. 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.

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. Sharp, M. Biavenda; A.P. Dhillan, M.A. Thomson, P. Harvey, A. Valentejo, E.E. Davies, J.A. Walker-Smith

Summary

Background. We identified a case-control study of 120 one-year-old children who, after a course of sequential

Introduction

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 RUIWEI NEEB

In 2013, three years after he co-authored a paper showing that humans are remarkably good at self-deception, Daniel Kahneman, a psychologist at Princeton University, was asked to sign off on a manuscript that had been submitted to a top scientific journal. The paper, which had been written by a colleague, claimed that Kahneman had given his blessing to one of its results. In response, Kahneman wrote a note to the editor, declining to sign off, adding that everything in the paper's critical section should be considered suspect.

© 2015 Scientific American, a division of Springer Nature America, Inc.

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)

or a brief moment in 2010, McLean basked in the brief glow of her scientific glory, covered that year's events equaling the world in black and white. The media was "plain sailing," including PhD students at the time in Charlottesville, Virginia. Data analysis by 2,000 people convinced her that she had made a breakthrough that would revolutionize our understanding of what old-field fire effects are. The *Hepaticas*, says, "and the data provided a P value, a common index for evidence, near 0.01 — usually if $p < 0.05$ is significant." Publication in a journal seemed within reach. More likely, though, was reality intruding. Seven years on now, reproducibility, a term Brian Nosek, decided to coin, with it came the P value — not even close to the core significance, 0.05. The effect had it with it. McLean's dream of a year



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.

```

graph TD
    A[Problem Statement] --> B[Review of literature]
    B --> C[Formulate hypothesis]
    C --> D[Design research]
    D --> E[Collect data]
    E --> F[Analyze data]
    F --> G[Interpret results]
    G --> H[Report findings]
    H --> I[Final report]
    
```

The diagram illustrates the sequential steps of a research project:

- Problem Statement**: The initial phase where the research question or hypothesis is defined.
- Review of literature**: A comprehensive search and analysis of existing studies related to the topic.
- Formulate hypothesis**: Develop a clear, testable statement based on the review.
- Design research**: Plan the methodology, including sampling, data collection methods, and statistical analysis.
- Collect data**: Implement the research design to gather the necessary information.
- Analyze data**: Use statistical methods to process and interpret the collected data.
- Interpret results**: Draw conclusions from the data analysis, considering the context and implications.
- Report findings**: Communicate the results through a formal report.
- Final report**: The completed document summarizing the entire research process and findings.

Physicians are in every way best suited to decide what, if any, treatment is appropriate for their patients. This is particularly true when the evidence is inconclusive or ambiguous. In such cases, physicians can make informed decisions based on the best available information, including the results of controlled trials and other studies, as well as their own clinical judgment and experience.

collaborate with crew to stop sailing along P-Exiles, and prevent the rest of the iceberg from sinking science. ■

designed to address this crisis. For more information, visit www.cdc.gov/niosh/epidemic-preparedness/.

¹For example, the Data Science Specialization offered by Johns Hopkins University in Baltimore, Maryland, and DataCamp.com.

so few people who analyze data must be trained in the relevant software and concepts. Thus, investigating what supervision data analysis should be required by their funding agency.

2. **Emmett, J. P., Makinson, D. G., & Sverdrup, U.** *Physical Oceanography*, 2nd edn. (Elsevier, 2011).
3. **Krogh, A., & Hertz, R.** *Neuroscience* 20(9), 4028–4030 (1989).
4. **Krogh, A., & Hertz, R.** *Neuroscience* 4(4), 381–392 (1989).

⁵ See also the discussion in *Statistical Methods in Biostatistics* (Corcoran, 1990) and *Practical Data Science* (Lohr, 2014).

412 | RAYBEE | VOL 328 | 29 APRIL 2013

© 2019 Pearson Education, Inc.

Digitized by srujanika@gmail.com

state changes

ct to change.

• 100

© 2013 Pearson Education, Inc.

Digitized by srujanika@gmail.com

© 2013 Pearson Education, Inc.

17

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



A tragedy of errors

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

Jut how error-prone and self-correcting is science? We have examined 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-regarded journal estimating how a change in food availability (such as a reduction in food weight), and he noted that the analysis applied a mathematical model that assumed all individuals would eat the same amount and others submitted a letter to the editor explaining the problem. Months later, we

were gratified to learn that the authors had responded to their letter, and the face of the paper had been changed. The finding is startling; this episode was an affirmation that science is self-correcting.

Not so fast. As it turns out, 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 believed to be errors were not being corrected, or were being invalidated, errors. These involve factual

mistakes or verbiage substantially from clearly wrong to clearly right in ways that are not easily apparent to nonspecialists.

After attempting to address more than 25 of these errors with letters to authors and editors, we gave up. After all, as it is done more, we had to stop — the work took too much of our time. — Our efforts did, however, lead us to a paper we discuss separately (see 'Three common errors') and showed how journals and authors can be faced with mistakes that need correction.

We learned that post-publication ▶

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

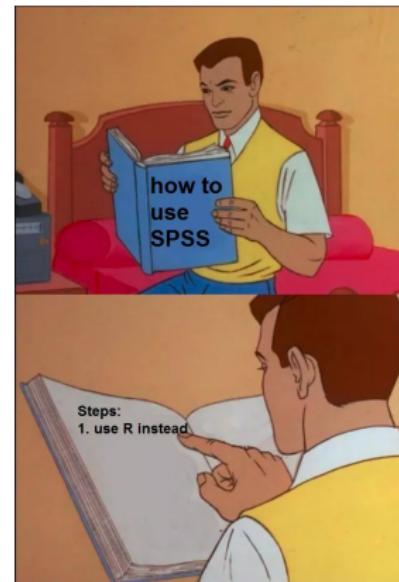
## [1] 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

# 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 ...

A screenshot of a Windows Command Prompt window titled "Windows Task Manager". The window contains the following R console output:

```
Microsoft Windows [Version 10.0.14393]
[Administrator]

C:\Users\AARHUS\Documents>R
R version 3.4.2 (2017-09-28) -- "Very Merry X-Mas Edition"
Copyright (C) 2017 The R Foundation for Statistical Computing
ISBN 3-900051-07-0

Platform: x86_64-w64-mingw32/x64 (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.

R is copyrighted free software; a complete distribution is
available at http://www.R-project.org.
All contributions are made available under the same license as R itself.

Type 'citation()' on R command-line for distribution details.

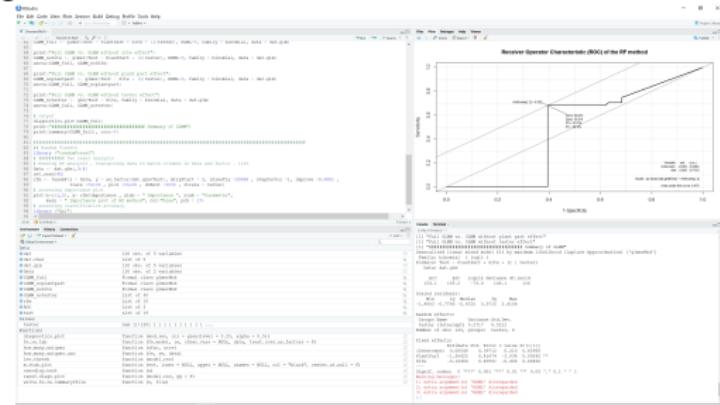
R is a language and environment for statistical computing and graphics
and is based on the S language and environment which was developed at
the Department of Statistics, University of Kent, Kent, UK. It now runs
on a wide variety of operating systems, and can be
interfaced with over 1,000 other packages that add functionality
to it. An extensive array of software packages are available
from the R Development Team (http://www.R-project.org).
```

... 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

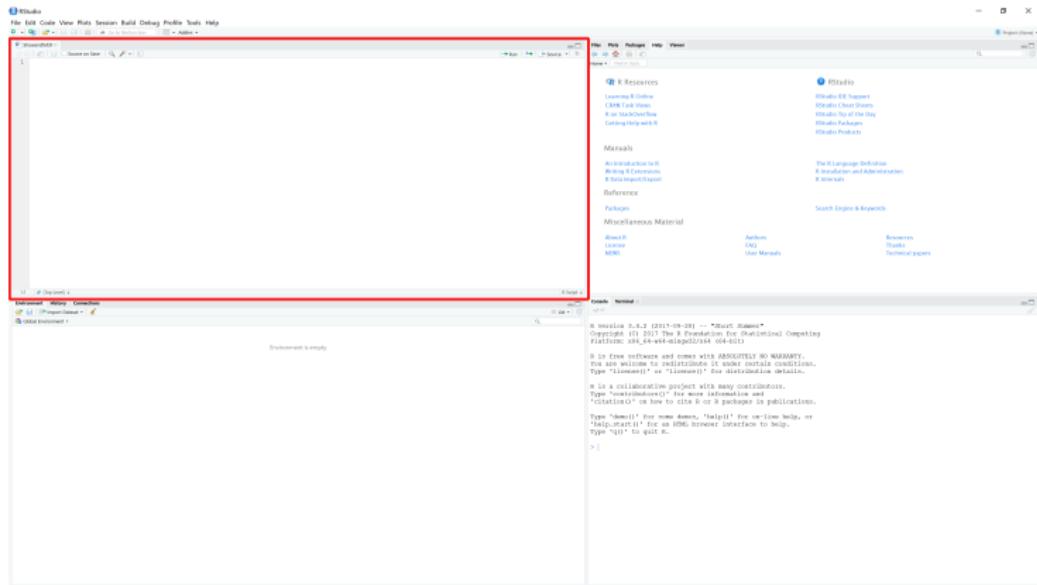
Running R through an editor...



... is a **much better idea!**

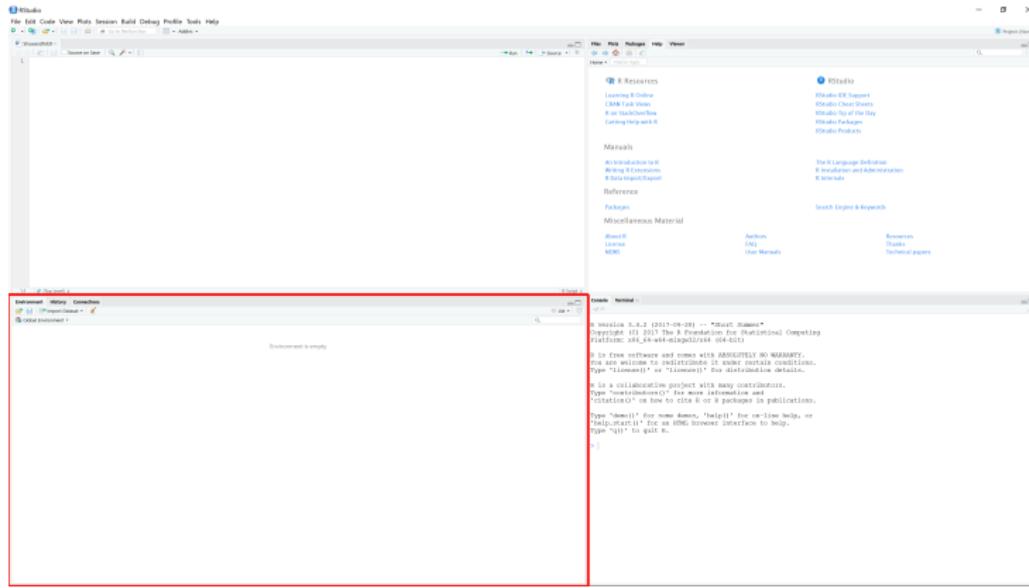
I recommend RStudio (<https://www.rstudio.com/>). If you use it a lot, I also recommend changing the appearance to 'Vibrant Ink' (setting located in the 'Global Options' window nested within the 'Tools' tab).

Layouts - The Editor Explained



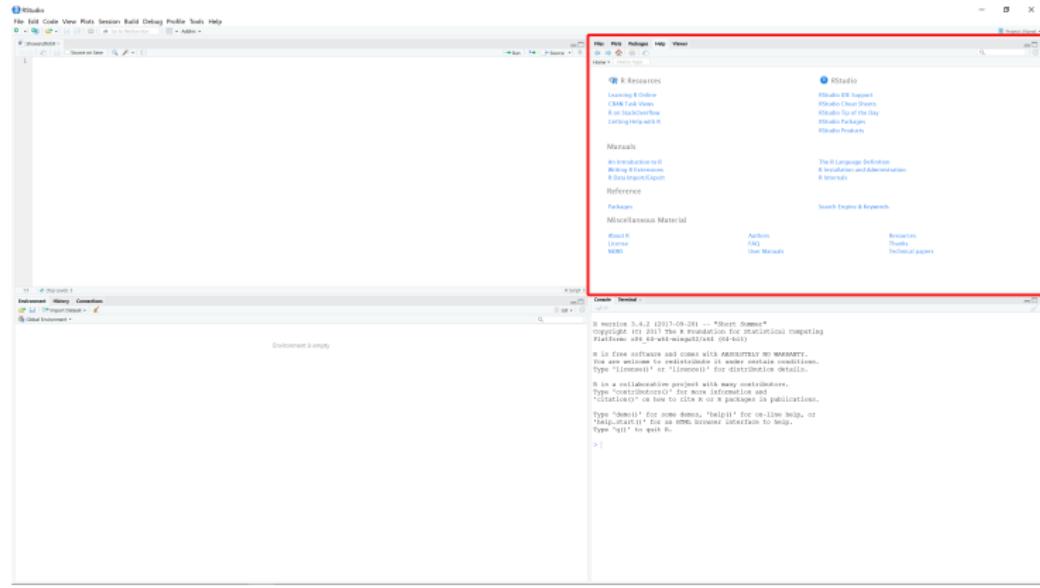
The **Source** is where you load scripts and write most of your coding document.

Layouts - The Editor Explained



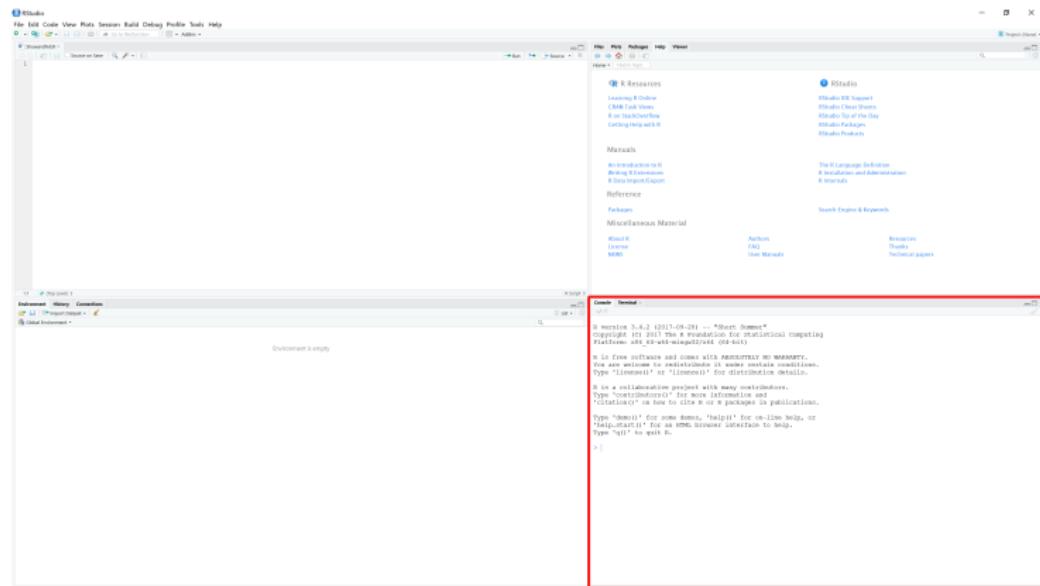
The **Environment, History, Connections** is where you will be able to quickly access all objects of your current R session.

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."

