

Experimental Design for Laboratory Biologists

Maximising Information and Improving Reproducibility

Specifically intended for lab-based biomedical researchers, this practical guide shows how to design experiments that are reproducible, with low bias, high precision, and results that are widely applicable. With specific examples from research, using both cell cultures and model organisms, it explores key ideas in experimental design, assesses common designs, and shows how to plan a successful experiment. It demonstrates how to control biological and technical factors that can introduce bias or add noise, and covers rarely discussed topics such as graphical data exploration, choosing outcome variables, data quality control checks, and data preprocessing. It also shows how to use R for analysis, and is designed for those with no prior experience. An accompanying website (<https://stanlazic.github.io/EDLB.html>) includes all R code, data sets, and the labstats R package.

This is an ideal guide for anyone conducting lab-based biological research, from students to principal investigators working either in academia or industry.

Stanley E. Lazic holds a PhD in neuroscience and a Masters in computational biology from the University of Cambridge and has conducted research at Oxford, Cambridge, and Harvard. He has written several papers on reproducible research and on the design and analysis of biological experiments and has published in *Science* and *Nature*. He is currently a Team Leader in Quantitative Biology (Statistics) at AstraZeneca.

Experimental Design for Laboratory Biologists

Maximising Information and Improving
Reproducibility

STANLEY E. LAZIC



CAMBRIDGE
UNIVERSITY PRESS

CAMBRIDGE
UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107074293

© Stanley E. Lazic 2016

This publication is in copyright. Subject to statutory exception and to the provisions of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published 2016

Printed in the United Kingdom by TJ International Ltd., Padstow, Cornwall

A catalogue record for this publication is available from the British Library

ISBN 978-1-107-07429-3 Hardback

ISBN 978-1-107-42488-3 Paperback

Additional resources for this publication at <https://stanlazic.github.io/EDLB.html>

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this publication, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

To my teachers and mentors

Contents

<i>Preface</i>	<i>page xi</i>
<i>Abbreviations</i>	<i>xiv</i>
1 Introduction	1
1.1 What is reproducibility?	1
1.2 The psychology of scientific discovery	3
1.2.1 Seeing patterns in randomness	4
1.2.2 Not wanting to miss anything	5
1.2.3 Psychological cliff at $p = 0.05$	6
1.2.4 Neglect of sampling variability	8
1.2.5 Independence bias	12
1.2.6 Confirmation bias	15
1.2.7 Expectancy effects	17
1.2.8 Hindsight bias	17
1.2.9 Herding effect	18
1.2.10 How the biases combine	19
1.3 Are most published results wrong?	21
1.3.1 What statisticians say	22
1.3.2 What scientists say	24
1.3.3 Empirical evidence I: questionable research practices	25
1.3.4 Empirical evidence II: quality of studies	26
1.3.5 Empirical evidence III: reproducibility of studies	28
1.3.6 Empirical evidence IV: publication bias	29
1.3.7 Scientific culture not conducive to ‘truth-finding’	30
1.3.8 Low prior probability of true effects	32
1.3.9 Main statistical sources of bias in experimental biology	34
1.4 Frequentist statistical inference	37
1.5 Which statistics software to use?	44
Further reading	46
2 Key Ideas in Experimental Design	48
2.1 Learning versus confirming experiments	49
2.2 The fundamental experimental design equation	52
2.3 Randomisation	59
2.4 Blocking	60
2.5 Blinding	62

2.6	Effect type: fixed versus random	65
2.7	Factor arrangement: crossed versus nested	66
2.8	Interactions between variables	68
2.9	Sampling	72
2.10	Use of controls	74
2.11	Front-aligned versus end-aligned designs	76
2.12	Heterogeneity and confounding	78
2.12.1	Batches	82
2.12.2	Plates, arrays, chips, and gels	84
2.12.3	Cages, pens, and tanks	84
2.12.4	Subject/sample characteristics	85
2.12.5	Litters	85
2.12.6	Experimenter characteristics	86
2.12.7	Time effects	86
2.12.8	Spatial effects	89
2.12.9	Useful confounding	91
	Further reading	93
3	Replication (what is 'N'?)	94
3.1	Biological units	95
3.2	Experimental units	96
3.3	Observational units	99
3.4	Relationship between units	100
3.4.1	Randomisation at the top of the hierarchy	103
3.4.2	Randomisation at the bottom of the hierarchy	109
3.4.3	Randomisation at multiple levels	118
3.5	How is the experimental unit defined in other disciplines?	121
4	Analysis of Common Designs	123
4.1	Preliminary concepts	124
4.1.1	Partitioning the sum of squares	124
4.1.2	Counting degrees of freedom	132
4.1.3	Multiple comparisons	135
4.2	Background to the designs	144
4.3	Completely randomised designs	144
4.3.1	One factor, two groups	144
4.3.2	One factor, multiple groups	145
4.3.3	Two factors, crossed	149
4.3.4	One factor with subsamples (pseudoreplication)	157
4.3.5	One factor with a covariate	166
4.4	Randomised block designs	170
4.4.1	With no replication	171
4.4.2	With genuine replication	173
4.4.3	With pseudoreplication	175

4.5	Split-unit designs	175
4.6	Repeated measures designs	181
	Further reading	191
5	Planning for Success	192
5.1	Choosing a good outcome variable	192
5.1.1	Qualitative criteria	193
5.1.2	Statistical criteria	194
5.2	Power analysis and sample size calculations	206
5.2.1	Calculating the sample size	207
5.2.2	Calculating power	210
5.2.3	Calculating the minimum detectable effect	210
5.2.4	Power curves	211
5.2.5	Simulation-based power analysis	212
5.3	Optimal experimental designs (rules of thumb)	220
5.3.1	Use equal n with two groups	223
5.3.2	Use more controls when comparing multiple groups to the control	225
5.3.3	Use fewer factor levels	227
5.3.4	Increase the variance of predictor variables	229
5.3.5	Ensure predictor variables are uncorrelated	235
5.3.6	Space observations out temporally and spatially	238
5.3.7	Sample more intensively where change is faster	240
5.3.8	Make use of blocking and covariates	245
5.3.9	Crossed factors are better than nested	251
5.3.10	Add more samples instead of subsamples	252
5.3.11	Have 10 to 20 samples to estimate the error variance	253
5.4	When to stop collecting data?	256
5.5	Putting it all together	259
5.6	How to get lucky	266
5.7	The statistical analysis plan	267
5.7.1	Why bother?	267
5.7.2	What to include in the SAP	269
	Further reading	271
6	Exploratory Data Analysis	272
6.1	Quality control checks	273
6.1.1	Data layout	274
6.1.2	Possible and plausible values	276
6.1.3	Uniqueness	281
6.1.4	Missing values	289
6.1.5	Factor arrangement	294
6.2	Preprocessing	296
6.2.1	Aggregating and summarising	296
6.2.2	Normalising and standardising	297

6.2.3	Correcting and adjusting	297
6.2.4	Transforming	297
6.2.5	Filtering	298
6.2.6	Combining	300
6.2.7	Pitfalls of preprocessing	300
6.3	Understanding the structure of the data	307
6.3.1	Shapes of distributions	307
6.3.2	Effects of interest	313
6.3.3	Spatial artefacts	326
6.3.4	Individual profiles	335
	Further reading	340
Appendix A Introduction to R		341
A.1	Installing R	341
A.2	Writing and editing code	342
A.3	Basic commands	343
A.4	Obtaining help	346
A.5	Setting options	347
A.6	Loading and saving data	347
A.7	Objects, classes, and special values	349
A.8	Conditional evaluation	353
A.9	Creating functions	355
A.10	Subsetting and indexing	357
A.11	Looping and applying	361
A.12	Graphing data	364
A.13	Distributions	371
A.14	Fitting models	375
Appendix B Glossary		381
<i>References</i>		390
<i>Index</i>		411

Preface

Everything of importance has been said before by somebody who did not discover it.

Alfred North Whitehead

Everything that needs to be said has already been said. But since no one was listening, everything must be said again.

André Gide

True to the above quotes, most of this book's contents have appeared in print before, but often where biologists are unlikely to look – statistics journals and books, and methods papers in other fields. My task is to translate ideas known to statisticians into the language of experimental biology.¹ With a background in both biology (BSc, PhD, postdoc) and data analysis (MPhil in Computational Biology and over seven years working as a preclinical statistician in the pharmaceutical industry), hopefully I am fluent enough in both languages to perform a successful translation.²

The contents of this book have little overlap with other statistics-for-biologists books because they mostly focus on statistical analysis. Analysis is but one step of the scientific workflow (Figure 0.1), and before you can analyse data you need to do an experiment. This requires planning, good execution, and quality control checks. These critical topics are rarely taught to biologists, who are expected to learn them on their own. The consequence of this approach is predictable; some biologists obtain the necessary skills, but many do not. This book focuses on the first three steps of the scientific workflow, and data analysis is briefly discussed in Chapter 4.

This book was written to improve the quality of research conducted in academic, government, and industrial labs and institutions. Scientists and funders now recognise that bias and irreproducibility are undermining preclinical biomedical research [2, 5, 28, 30, 42, 80, 83, 84, 123, 172, 240, 251, 305, 316, 342]. There are many reasons why experiments cannot be reproduced (discussed in Chapter 1) and this book focuses on the role that experimental design and data analysis have on making results reproducible.

¹ The term *biology* refers to laboratory-based experimental biology throughout. 'Field biologists' also conduct experiments, and most statistics-for-biologists books target this audience.

² There are some novel ideas here, such as the distinction between front-aligned and end-aligned designs (Section 2.11) and the distinction between biological, experimental, and observational units, to replace the biological versus technical replicate distinction (all of Chapter 3).

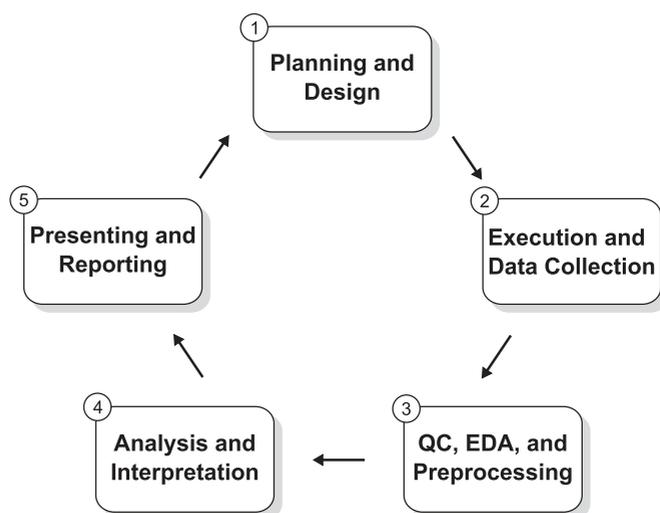


Fig. 0.1

The scientific workflow. This book focuses on steps 1–3. QC = quality control; EDA = exploratory data analysis.

Prerequisites

This book is for experimental biologists, at any level, conducting basic research or with an applied, clinical, or translational focus. Knowledge covered in an introductory statistics-for-biologists course is assumed, and concepts like the standard deviation and common statistical tests such as the *t*-test, analysis of variance (ANOVA), regression, and correlation should be familiar. It is fine if some time has passed since you formally covered these topics. Mathematical proofs are not included and equations are kept to a minimum, but given the subject, are unavoidable. The emphasis is on the ideas, concepts, and principles, and how to implement them. Hand calculations are unnecessary because statistical software is available.

Quantitative researchers who analyse biological data such as statisticians, bioinformaticians, and computational biologists might also find this book useful. Topics of interest include sources of heterogeneity and confounding in biological experiments (Section 2.12), quality control checks for biological data (Section 6.1), and understanding which types of replication address biologically interesting questions (Chapter 3).

The freely available R statistics language is used for data analysis and graphs.³ Prior knowledge is useful, but not required. The Appendix gives a brief introduction to R and the examples in the main text assume familiarity with this material. The topics however can be followed without learning or using R. The data sets can be found in the `labstats` package on CRAN⁴ and R code can be downloaded from GitHub.⁵

³ Available at www.r-project.com

⁴ <https://cran.r-project.org/web/packages/labstats/>

⁵ <https://stanlazic.github.io/EDLB.html>

The key prerequisite to derive maximum value from this book is experience conducting biological experiments and analysing the subsequent data – and the more experience the better!

How to read this book

Chapters 1–5 should be read in order as later material depends on earlier ideas, but Chapter 6 on Exploratory Data Analysis can be read at any time. Chapters 1–3 contain no R code, but for Chapters 4–6 sitting in front of a computer and running the code will reinforce the ideas.

Ideas or concepts discussed in detail later in the book will inevitably have to be mentioned earlier. To avoid excessive cross-referencing, the glossary lists the page where the main discussion of the entry is located (if there is one). For example, the term *experimental unit* is mentioned for the first time in this preface, but is discussed extensively in Section 3.2. The glossary entry for this term provides a short definition and indicates that further information can be found on page 96.

Typographical conventions

Constant width font is used for R code, R output, and when referring to R functions or objects. Lines of code entered by the user start with ‘>’ or ‘+’. These symbols do not need to be entered, only the code that follows them. A sign like the one in the margin draws attention to a warning, a key point, a subtlety with R, or a concept that is often misunderstood.



Acknowledgements

This book has benefited greatly from comments by Maarten van Dijk, Irmgard Amrein, and especially Lutz Slomianka. Pierre Farmer and Miguel Camargo also provided constructive feedback on earlier drafts. My wife, Brynn, has read every word in this book, which is beyond the call of duty, and her comments have improved it immensely. I also thank her for her support, well, at least until page 305, at which point she declared, ‘You should stop now; no one wants to read that much about statistics.’ I didn’t always follow everyone’s good advice, but I am grateful for their input.

Katrina Halliday and Jade Scard at Cambridge University Press were a pleasure to work with and made the whole process easy and enjoyable. I also thank Judith Shaw for her expert copy-editing. Finally, I would like to thank the developers and contributors of the free software R, Emacs, LaTeX, JabRef, knitr, and Inkscape, which I used to write this book.

S.E. Lazic
Cambridge, 2016

Abbreviations

AIPE	Accuracy in parameter estimation
ALS	Amyotrophic lateral sclerosis
ANCOVA	Analysis of covariance
ANOVA	Analysis of variance
AUC	Area under the curve
BMI	Body mass index
BU	Biological unit
CCC	Concordance correlation coefficient
CCLE	Cancer Cell Line Encyclopedia
CI	Confidence (frequentist) or Credible (Bayesian) interval
CRAN	Comprehensive R archive network
CRD	Completely randomised design
CSF	Cerebrospinal fluid
CSR	Complete spatial randomness
CV	Coefficient of variation
DAMP	Damage-associated molecular pattern
df	Degrees of freedom
DoE	Design of experiments
DS	Diallyl sulfide
ED50	Median (half) effective dose
EDA	Exploratory data analysis
ES	Effect size
ESS	Emacs Speaks Statistics
EU	Experimental unit
FORE-SCI	Facilities of Research Excellence – Spinal Cord Injury
FOV	Field of view
GI	Gastrointestinal
GLM	Generalised linear model
GUI	Graphical user interface
Gst	Glutathione-S-transferase
HARKing	Hypothesising after the results are known
HSD	Honestly significant difference
ICC	Intraclass correlation coefficient
i.p.	Intraperitoneally
IQR	Interquartile range

KO	Knock out
LME	Linear mixed-effects model
LOD	Limit of detection
LSD	Least significant difference
MAD	Median absolute deviation
MAR	Missing at random
MCAR	Missing completely at random
MED	Minimum effective dose
MNAR	Missing not a random
NGS	Next generation sequencing
NHST	Null hypothesis significance testing
NIH	National Institutes of Health (USA)
NINDS	American National Institute of Neurological Disorders and Stroke
OU	Observational unit
PCA	Principal components analysis
PI	Principal investigator
PK	Pharmokinetic
QC	Quality control
QRP	Questionable research practice
qPCR	Quantitative polymerase chain reaction
RE	Relative efficiency
RIN	RNA integrity number
RM-ANOVA	Repeated measures analysis of variance
SAP	Statistical analysis plan
SD	Standard deviation
SEM	Standard error of the mean
siRNA	small interfering RNA
SNP	Single nucleotide polymorphism
SOD1	Superoxide dismutase 1 (gene)
SS	Sum of squares
RSS	Residual sum of squares
SUTVA	Stable unit-treatment value assumption
TSS	Total sum of squares
VPA	Valproic acid
WT	Wild type