A Computational Approach to Statistical Arguments in Ecology and Evolution

Scientists need statistics. Increasingly this is accomplished using computational approaches. Freeing readers from the constraints, mysterious formulas, and sophisticated mathematics of classical statistics, this book is ideal for researchers who want to take control of their own statistical arguments.

It demonstrates how to use spreadsheet macros to calculate the probability distribution predicted for any statistic by any hypothesis. This enables readers to use anything that can be calculated (or observed) from their data as a test statistic, and hypothesize any probabilistic mechanism that can generate data sets similar in structure to the one observed.

A wide range of natural examples drawn from ecology, evolution, anthropology, paleontology, and related fields give valuable insights into the application of the described techniques, while complete example macros and useful procedures demonstrate the methods in action and provide starting points for readers to use or modify in their own research.

George F. Estabrook is a Professor of Botany in the Department of Ecology and Evolutionary Biology at the University of Michigan, Ann Arbor. He is interested in the application of mathematics and computing to biology, and has taught graduate courses on the subject for more than 30 years.

A Computational Approach to Statistical Arguments in Ecology and Evolution

GEORGE F. ESTABROOK University of Michigan, Ann Arbor



CAMBRIDGE UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314-321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre, New Delhi - 110025, India

103 Penang Road, #05-06/07, Visioncrest Commercial, Singapore 238467

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/9781107004306

© George F. Estabrook 2011

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 2011

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

ISBN 978-1-107-00430-6 Hardback

Additional resources for this publication at www.cambridge.org/9781107004306

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.

Contents

Acknowledgments

1	Introduction	1
1.1	About the book	1
1.2	Basic principles	10
1.3	Scientific argument	14
2	Programming and statistical concepts	20
2.1	Computer programming	20
2.2	You start programming	31
2.3	Completing the service berry example	36
2.4	Sub CARPEL	49
2.5	You practice	53
3	Choosing a test statistic	59
3.1	Significance of what	59
3.2	Implement the program	63
3.3	Sub PERIOD	71
4	Random variables and distributions	77
4.1	Random variables	77
4.2	Distributions	81
4.3	Arithmetic with random variables	88
4.4	Expected value and variance	94
5	More programming and statistical concepts	101
5.1	Re-sampling data	101
5.2	Procedures	110
5.3	Testing procedures	115

6	Parametric distributions	12
6.1	Basic concepts	12
6.2	Poisson distribution	12
6.3	Normal distribution	13
6.4	Negative binomial, Chi Square, and F distributions	13
6.5	Percentiles	13
7	Linear model	14
7.1	Linear model	14
7.2	Quantifying error	14
7.3	Linear model in matrix form	14
7.4	Using a linear model	1
7.5	Hypotheses of random for a linear model	1
7.6	Two-way analysis of variance	1
8	Fitting distributions	1
8.1	Estimation of parameters	1
8.2	Goodness of fit	1
9	Dependencies	1
9.1	Interpreting mixtures	1
9.2	Series of dependent random variables	1
9.3	Analysis of covariance	1
9.4	Confounding dependencies	2
9.5	Sub SEXDIMO	2
10	How to get away with peeking at data	2
11	Contingency	2
11.1	What is contingency	2
11.2	ACTUS2	2
11.3	Spreadsheet ACTUS	2
11.4	Sub ACTUS	2
	References	2
	Index	2

Acknowledgments

In the mid 1980s, Neal Oden quit his boring job as a computer programmer and joined the Ph.D. program of the Department of Biology at the University of Michigan to become my eighth Ph.D. student. Neal had read a journal article in which the authors claimed that when data sampled a given probability distribution, the predicted distribution of a statistic they had invented would asymptotically converge to a known pre-calculated distribution of classical statistics. From the authors' explanation, it was not clear to Neal why. So he wrote a computer program that used a random number generator to sample the given distribution to create a large data set, and calculate the statistic; his program repeated this 1000 times to estimate the predicted distribution. The middle 80% was pretty close to the classical distribution, but upper and lower 10% differed. By doing this, Neal showed me how easy it is to make accurate estimations of the probability distribution for any statistic predicted by any hypothesis of random. Thank you, Neal. I began to apply Neal's idea in collaboration with other colleagues, and eventually offered a course to graduate students.

The many students who took my course, BIO 480 and later EEB 480, have contributed much to the development of this book, which has gone through a half dozen unpublished editions over the past decade in response to their helpful suggestions and requests.

One of the major stumbling blocks for students trying to use these concepts and techniques has been the need to learn enough computer programming to implement them. I taught my course originally using PASCAL, a programming language designed explicitly to teach students good programming practices. Even with PASCAL, the INPUT, OUTPUT, and MEMORY MANAGEMENT details were especially challenging for students who had never programmed before. A few years ago, my 18th Ph.D. student, Don Schoolmaster, showed me how to program spreadsheet macros that read from, and write to,

viii Acknowledgments

spreadsheets. This eliminates problems with INPUT, OUTPUT, and MEMORY MANAGEMENT, which makes the programming much more accessible to students. Virtually all students are familiar with spreadsheets, even if they have never programmed a macro. Thank you, Don.

Keith Hunley, a former EEB 480 student and now a professor at the University of New Mexico, read much of the pre-final draft of the book. He made hundreds of suggestions that improved the clarity and organization. Thank you, Keith.

My wife, Virginia Hutton Estabrook Ph.D., helped me with drawing the figures, organizing the literature cited, and grappling with the WINDOWS operating system. She also read some of the pre-final draft and made helpful suggestions. Thank you, Virginia.