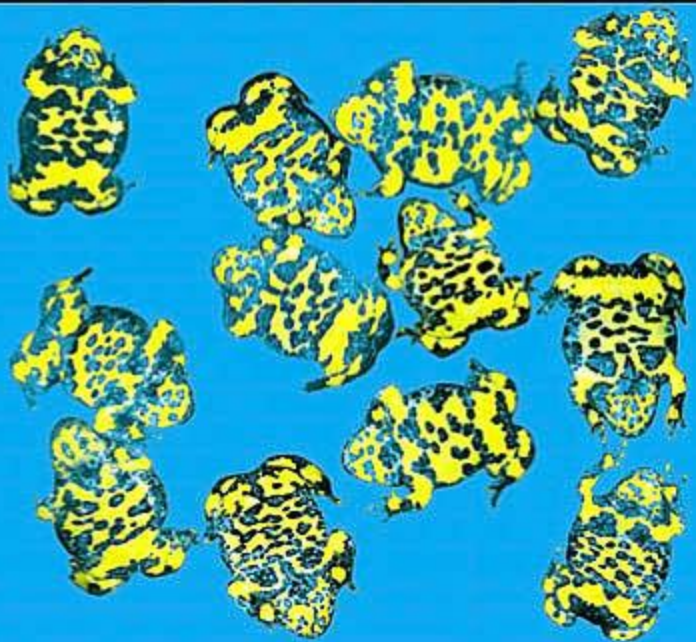




MODEL SELECTION AND MULTIMODEL INFERENCE
A Practical Information-Theoretic Approach

SECOND EDITION

KENNETH P. BURNHAM • DAVID R. ANDERSON



Springer

New York

Berlin

Heidelberg

Barcelona

Hong Kong

London

Milan

Paris

Singapore

Tokyo

Model Selection and Multimodel Inference

A Practical Information-Theoretic
Second Edition

With 31 Illustrations



Springer

and Wildlife Research Unit
Colorado State University
Fort Collins, CO 80523-1484
USA

Cover Illustration: The cover was assembled from photos of the yellow-bellied sapsucker (*Sphyrapicus varius*) taken by Jonas Barandum as part of his Ph.D. program at the University of Wisconsin. The photos show the individual patterns on their abdomen from a few weeks before they moult to new patterns that remain unchanged until death. Two pairs are duplicates—but which two? Cover photographs by Dr. Jonas Barandum, St. Gallen, Switzerland. Cover design by Dr. Jonas Barandum.

Library of Congress Cataloging-in-Publication Data

Burnham, Kenneth P.

Model selection and multimodel inference : a practical information-theoretic approach / Kenneth P. Burnham, David R. Anderson.—2nd ed.
p. cm.

Rev. ed. of: Model selection and inference. © 1998.

Includes bibliographical references (p.).

ISBN 0-387-95364-7 (alk. paper)

1. Biology—Mathematical models. 2. Mathematical statistics. I. Burnham, Kenneth P. II. Anderson, David R. III. Title.

Model selection and inference. II. Title.

QH323.5 B87 2002

570'.1'51—dc21

ISBN 0-387-95364-7

Printed on acid-free paper.

© 2002, 1998 Springer-Verlag New York, Inc.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer-Verlag New York, Inc., 175 Fifth Avenue, New York, NY 10011, USA), except for brief excerpts in connection with reviews or scholarly articles. This work may be registered with a copyright office. This work may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or by any information storage and retrieval system, now known or hereafter developed, for personal or internal use, or the personal or internal use of specific clients. The use in this publication of trade names, trademarks, service marks, and similar designations, even if not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed in the United States of America.

9 8 7 6 5 4 3 2

SPIN 10853081

www.springer-ny.com

Springer-Verlag New York Berlin Heidelberg

A member of BertelsmannSpringer Science+Business Media GmbH

*To my mother and father, Lucille R. (deceased) and J
(deceased), and my son and daughter, Shawn P. and*

*To my parents, Charles R. (deceased) and Leta M. A
Dalene F. Anderson; and my daughters, Tama
Adrienne M. Anderson*

This page intentionally left blank

Preface

We wrote this book to introduce graduate students and various scientific disciplines to the use of information-theoretic methods in the analysis of empirical data. These methods allow the selection of a “best” model and a ranking and weighting of the models from a pre-defined set. Traditional statistical inference can be based on the selected best model. However, we now emphasize that information-theoretic approaches allow formal inference to be based on more than one model (multimodel inference). Such procedures lead to more robust inferences in complex cases, and we advocate these approaches throughout the book.

The second edition was prepared with three goals in mind. First, we tried to improve the presentation of the material. Boxes now contain technical expressions and points. Some reorganization has been done to improve the flow of concepts, and a new chapter has been added. Chapters 4 and 5 have been streamlined in view of the detailed theory provided in Chapter 6. Second, concepts related to making formal inferences from multiple models (multimodel inference) have been emphasized throughout the book, particularly in Chapters 4, 5, and 6. Third, new technical material has been added to Chapters 5 and 6. Well over 100 new references to the literature are given. These changes result primarily from our experience with several seminars, workshops, and graduate courses on machine learning. In addition, we have done substantially more thinking about the problem of reading the literature since writing the first edition, and we hope this has led to further insights.

Information theory includes the celebrated Kullback–Leibler divergence between two models (actually, probability distributions), and

maximized log likelihood. This measure, now called Akaike's information criterion (AIC), provided a new paradigm for model selection in empirical data. His approach, with a fundamental link to the empirical data, is relatively simple and easy to use in practice, but little understood in the applied sciences that use it.

We do not accept the notion that there is a simple "true model" in the empirical sciences. Instead, we view modeling as an exercise in understanding the explainable information in the empirical data, in the sense of being a sample from some well-defined population or process. We view modeling as a fabric in the tapestry of science. Selecting a model approximating the inference from the data, and the "effects" (represented by parameters) can be supported by the data, on Akaike's information criterion (and various extensions), is a parsimonious model as a basis for statistical inference. Modeling on information theory represents a quite different approach. In the sciences, and the resulting selected model may differ substantially from the selection based on some form of statistical null hypothesis.

We recommend the information-theoretic approach for inference from observational studies. In this broad class of studies, where various hypothesis-testing approaches have no theoretical justification, they often perform poorly. For classic experiments (control–treatment, randomization and replication) we generally support the traditional approach (e.g., analysis of variance); there is a very large literature on this. However, for complex experiments we suggest considering alternative models, hence on estimation of the size and precision of the effects and on parsimony, with far less emphasis on "testing hypotheses, leading to the arbitrary classification "significant" versus "not significant". Instead, a strength of evidence approach is advocated.

We do not claim that the information-theoretic methods are the best for a particular situation. They do represent a unified approach, an extension of likelihood theory, an important application of information theory, and they are objective and practical to employ across a wide range of empirical problems. Inference from multiple models, or the selection of the "best" model, by methods based on the Kullback–Leibler divergence, is certainly better than other methods commonly in use now (e.g., testing of various sorts, the use of R^2 , or merely the use of a single model). In particular, subjective data dredging leads to overfitting and the attendant problems in inference, and is to be strongly discouraged in more confirmatory studies.

Parameter estimation has been viewed as an optimization problem for at least eight decades (e.g., maximize the log-likelihood or minimize the sum of squared deviations). Akaike viewed his AIC and BIC as "... a natural extension of the classical maximum likelihood

beyond merely the computation and interpretation of AIC. This is a new paradigm for inference from empirical data; it refocuses on a variety of considerations and modeling prior to the actual model selection, under the information-theoretic approach. Model selection, under the information-theoretic approach, attempts to identify the (likely) best model, orders the models from best to worst, and produces a weight of evidence that each model has for an inference.

Several methods are given that allow model selection incorporated into estimates of precision (i.e., multimodel inference). This book is to present and illustrate a consistent methodology that includes model selection, estimation of model parameters and model prediction in a unified manner, under a compelling common framework. The book also explains other information criteria (e.g., AIC_c , $QAIC_c$, and $WAIC$), and several examples to illustrate various technical issues, in comparison with BIC, a type of dimension consistent criterion. The book provides many references to the technical literature for those who wish to further on these topics.

This is an applied book written primarily for biologists using models for making inferences from empirical data. This is a science book; we say relatively little about decision making in management science. Research biologists working either in the field or in the laboratory will find simple methods that are likely to be useful. Biologists in other life sciences, econometrics, and medicine might also find the material useful but with some caveats. Examples that have been taken largely from ecological studies of vertebrates, as these are our interests. Applied statisticians and information-theoretic methods presented here quite useful. This book is native to the null hypothesis testing approach that has become standard but is uninformative. We hope material such as this will find its way into where applied data analysis and associated science philosophy. This book might be useful as a text for a course for students with experience and education in statistics and applied data analysis. The audience includes honors or graduate students in the biological and statistical sciences. Those interested in the empirical sciences will find the material useful because it offers an effective alternative to (1) the null hypothesis test, yet often both complex and uninformative, null hypothesis testing, and (2) the far less taught, but potentially very useful, Bayesian methods.

Readers should ideally have some maturity in the quantitative sciences and experience in data analysis. Several courses in contemporary statistics and methods as well as some philosophy of science would be helpful in understanding the material. Some exposure to likelihood inference is essential, but those with experience only in least squares regression will gain some useful insights. Biologists working in a

conceptual issues may preclude beginners. Chapters 1–4 are for all readers because they provide the essential material, in particular, for multimodel inference. Chapters 5 and 6 present more details and some new research results. Few readers will be able to read all of the material presented here after just one reading of the material; some additional consideration will often be necessary to understand the material. Underlying theory is presented in Chapter 7, and this material is more mathematical. A high-level summary of the material is provided in Chapter 8.

We intend to remain active in this subject area after this book has been published, and we invite comments from colleagues to help us learn more and understand differing points of view. We hope the book will not appear too dogmatic or idealized. We have tried to synthesize what we believe are important and incorporate these as recommendations in several of the chapters. This book is an effort to explore multimodel inference in some depth. We realize that there are many different attitudes toward data dredging than we advocate here. We also discuss other model selection methods, such as cross-validation, and the use of Bayesian methods. Indeed, we just learned (March, 2015) that the result derived as a Bayesian result and have added a note on this in the final page proofs (see Section 6.4.5). However, in the social sciences, we are compelled by the a priori approach of building models to represent research hypotheses, the use of information as a basis for selecting a best approximating model; model selection using multimodel inference methods, when truth is surely very uncertain; likelihood theory for deriving parameter estimators; and the use of model selection uncertainty into statistical inferences. In particular, moving beyond mere selection of a single best model by using the methods of multimodel inference.

Several people have helped us as we prepared the two editions. In particular, we acknowledge C. Chatfield, C. Hurvich, J. Rotella, R. Shibata, and K. Wilson for comments on the original manuscript. We are grateful to three anonymous reviewers for comments that allowed us to improve the first edition. D. Otiwola served as the reviewers for the second edition and offered many suggestions that were helpful; we greatly appreciate their excellent suggestions. Discussions with S. Buckland, R. Davis, R. Shibata, and C. Wilson were very useful. S. Beck, K. Bestgen, D. Beyers, L. Ellison, A. Fraley, B. Lubow, C. McCarty, M. Miller, and T. Shenk provide helpful comments and insights as part of a graduate course on model selection methods from the authors. C. Flather allowed us to use his data

by M. Zaslau under the supervision of G. White. C. Sosa used his kangaroo data from Wallaby Creek. P. Lukacs contributed analysis and some of the Monte Carlo studies of the body mass. J. Kullback allowed us to use a photo of his father, and H. R. Shibata, and K. Takeuchi kindly sent us photos and biographies that appear in the book. Chelsea Publishing Company allowed us to use a photo of L. Boltzmann from the book *Wissenschaftliche Biographie Ludwig Boltzmann*, and the International Biometric Society allowed use of a photo of R. Fisher (from *Biometrics* 1964, taken from the University of Cambridge). J. Barandun provided the toad photos for the cover design, and B. Schmidt helped in coordination. S. Kane, B. Klein, A. Lyman, and T. Sundlov helped obtain references. J. Kimmel and L. Farkas helped in countless ways as we prepared this book.

We are happy to acknowledge the long-term cooperators in our operative Fish and Wildlife Research Unit: the Colorado Natural History Museum, Colorado State University, the Biological Resources Division of the Colorado Geological Survey, and the Wildlife Management Institute. Graduate students and faculty within the Department of Fisheries and Wildlife Management at Colorado State University provided a forum for our interests in the use of biological data. We extend our appreciation to several federal agencies, particularly the Department of the Interior, particularly the U.S. Geological Survey, for support of our long-term research interests.

Fort Collins, Colorado

This page intentionally left blank

Contents

Preface

About the Authors

Glossary

1 Introduction

- 1.1 Objectives of the Book
- 1.2 Background Material
 - 1.2.1 Inference from Data, Given a Model
 - 1.2.2 Likelihood and Least Squares Theory
 - 1.2.3 The Critical Issue: “What Is the Best to Use?”
 - 1.2.4 Science Inputs: Formulation of the S Candidate Models
 - 1.2.5 Models Versus Full Reality
 - 1.2.6 An Ideal Approximating Model
- 1.3 Model Fundamentals and Notation
 - 1.3.1 Truth or Full Reality f
 - 1.3.2 Approximating Models $g_i(x|\theta)$
 - 1.3.3 The Kullback–Leibler Best Model $g_i(\cdot)$
 - 1.3.4 Estimated Models $g_i(x|\hat{\theta})$
 - 1.3.5 Generating Models
 - 1.3.6 Global Model

1.4	Inference and the Principle of Parsimony
1.4.1	Avoid Overfitting to Achieve a Good Model
1.4.2	The Principle of Parsimony
1.4.3	Model Selection Methods
1.5	Data Dredging, Overanalysis of Data, and Spurious Effects
1.5.1	Overanalysis of Data
1.5.2	Some Trends
1.6	Model Selection Bias
1.7	Model Selection Uncertainty
1.8	Summary

2 Information and Likelihood Theory: A Basis for Selection and Inference

2.1	Kullback–Leibler Information or Distance Between Two Models
2.1.1	Examples of Kullback–Leibler Distance
2.1.2	Truth, f , Drops Out as a Constant
2.2	Akaike’s Information Criterion: 1973
2.3	Takeuchi’s Information Criterion: 1976
2.4	Second-Order Information Criterion: 1978
2.5	Modification of Information Criterion for Overdispersed Count Data
2.6	AIC Differences, Δ_i
2.7	A Useful Analogy
2.8	Likelihood of a Model, $\mathcal{L}(g_i data)$
2.9	Akaike Weights, w_i
2.9.1	Basic Formula
2.9.2	An Extension
2.10	Evidence Ratios
2.11	Important Analysis Details
2.11.1	AIC Cannot Be Used to Compare Models with Different Data Sets
2.11.2	Order Not Important in Computing AIC
2.11.3	Transformations of the Response Variable
2.11.4	Regression Models with Differing Error Structures
2.11.5	Do Not Mix Null Hypothesis Testing with Information-Theoretic Criteria
2.11.6	Null Hypothesis Testing Is Still Important in Strict Experiments
2.11.7	Information-Theoretic Criteria Are Not for Strict Experiments
2.11.8	Exploratory Data Analysis

2.12.2	A Heuristic Interpretation
2.12.3	More on Interpreting Information-Theoretic Criteria
2.12.4	Nonnested Models
2.12.5	Further Insights
2.13	Bootstrap Methods and Model Selection Frequencies
2.13.1	Introduction
2.13.2	The Bootstrap in Model Selection: The Basic Idea
2.14	Return to Flather's Models
2.15	Summary

3 Basic Use of the Information-Theoretic Approach

3.1	Introduction
3.2	Example 1: Cement Hardening Data
3.2.1	Set of Candidate Models
3.2.2	Some Results and Comparisons
3.2.3	A Summary
3.3	Example 2: Time Distribution of an Insecticide in a Simulated Ecosystem
3.3.1	Set of Candidate Models
3.3.2	Some Results
3.4	Example 3: Nestling Starlings
3.4.1	Experimental Scenario
3.4.2	Monte Carlo Data
3.4.3	Set of Candidate Models
3.4.4	Data Analysis Results
3.4.5	Further Insights into the First Fourteen Nested Models
3.4.6	Hypothesis Testing and Information-Theoretic Approaches Have Different Selection Frequencies
3.4.7	Further Insights Following Final Model Selection
3.4.8	Why Not Always Use the Global Model for Inference?
3.5	Example 4: Sage Grouse Survival
3.5.1	Introduction
3.5.2	Set of Candidate Models
3.5.3	Model Selection
3.5.4	Hypothesis Tests for Year-Dependent Survival Probabilities

3.5.6	A Class of Intermediate Models
3.6	Example 5: Resource Utilization of <i>Anolis</i> Lizards
3.6.1	Set of Candidate Models
3.6.2	Comments on Analytic Method
3.6.3	Some Tentative Results
3.7	Example 6: Sakamoto et al.'s (1986) Simulated Data
3.8	Example 7: Models of Fish Growth
3.9	Summary

4 Formal Inference From More Than One Model: Multimodel Inference (MMI)

4.1	Introduction to Multimodel Inference
4.2	Model Averaging
4.2.1	Prediction
4.2.2	Averaging Across Model Parameters
4.3	Model Selection Uncertainty
4.3.1	Concepts of Parameter Estimation and Model Selection Uncertainty
4.3.2	Including Model Selection Uncertainty in Estimator Sampling Variance
4.3.3	Unconditional Confidence Intervals
4.4	Estimating the Relative Importance of Variables
4.5	Confidence Set for the K-L Best Model
4.5.1	Introduction
4.5.2	Δ_i , Model Selection Probabilities, and the Bootstrap
4.6	Model Redundancy
4.7	Recommendations
4.8	Cement Data
4.9	Pine Wood Data
4.10	The Durban Storm Data
4.10.1	Models Considered
4.10.2	Consideration of Model Fit
4.10.3	Confidence Intervals on Predicted Storm Probability
4.10.4	Comparisons of Estimator Precision
4.11	Flour Beetle Mortality: A Logistic Regression
4.12	Publication of Research Results
4.13	Summary

5 Monte Carlo Insights and Extended Examples

5.1	Introduction
5.2	Survival Models

5.2.3	An Extended Survival Model
5.2.4	Model Selection if Sample Size Is Huge or Truth Known
5.2.5	A Further Chain Binomial Model
5.3	Examples and Ideas Illustrated with Linear Regression
5.3.1	All-Subsets Selection: A GPA Example
5.3.2	A Monte Carlo Extension of the GPA Example
5.3.3	An Improved Set of GPA Prediction Models
5.3.4	More Monte Carlo Results
5.3.5	Linear Regression and Variable Selection
5.3.6	Discussion
5.4	Estimation of Density from Line Transect Sampling
5.4.1	Density Estimation Background
5.4.2	Line Transect Sampling of Kangaroos in Wallaby Creek
5.4.3	Analysis of Wallaby Creek Data
5.4.4	Bootstrap Analysis
5.4.5	Confidence Interval on D
5.4.6	Bootstrap Samples: 1,000 Versus 10,000
5.4.7	Bootstrap Versus Akaike Weights: A Comparison on QAIC _c
5.5	Summary

6 Advanced Issues and Deeper Insights

6.1	Introduction
6.2	An Example with 13 Predictor Variables and 8,191 Models
6.2.1	Body Fat Data
6.2.2	The Global Model
6.2.3	Classical Stepwise Selection
6.2.4	Model Selection Uncertainty for AIC _c
6.2.5	An A Priori Approach
6.2.6	Bootstrap Evaluation of Model Uncertainty
6.2.7	Monte Carlo Simulations
6.2.8	Summary Messages
6.3	Overview of Model Selection Criteria
6.3.1	Criteria That Are Estimates of K-L Information
6.3.2	Criteria That Are Consistent for K
6.3.3	Contrasts
6.3.4	Consistent Selection in Practice: Quasi-true Models
6.4	Contrasting AIC and BIC
6.4.1	A Heuristic Derivation of BIC

6.4.3	Performance Comparison
6.4.4	Exact Bayesian Model Selection Form
6.4.5	Akaike Weights as Bayesian Posterior Model Probabilities
6.5	Goodness-of-Fit and Overdispersion Revisited
6.5.1	Overdispersion \hat{c} and Goodness-of-Fit A General Strategy
6.5.2	Overdispersion Modeling: More Than One Model
6.5.3	Model Goodness-of-Fit After Selection
6.6	AIC and Random Coefficient Models
6.6.1	Basic Concepts and Marginal Likelihood Approach
6.6.2	A Shrinkage Approach to AIC and Random Effects
6.6.3	On Extensions
6.7	Selection When Probability Distributions Differ by Model
6.7.1	Keep All the Parts
6.7.2	A Normal Versus Log-Normal Example
6.7.3	Comparing Across Several Distributions An Example
6.8	Lessons from the Literature and Other Matters
6.8.1	Use AIC_c , Not AIC, with Small Samples
6.8.2	Use AIC_c , Not AIC, When K Is Large
6.8.3	When Is AIC_c Suitable: A Gamma Distribution Example
6.8.4	Inference from a Less Than Best Model
6.8.5	Are Parameters Real?
6.8.6	Sample Size Is Often Not a Simple Issue
6.8.7	Judgment Has a Role
6.9	Tidbits About AIC
6.9.1	Irrelevance of Between-Sample Variance of AIC
6.9.2	The G-Statistic and K-L Information
6.9.3	AIC Versus Hypothesis Testing: Results Very Different
6.9.4	A Subtle Model Selection Bias Issue
6.9.5	The Dimensional Unit of AIC
6.9.6	AIC and Finite Mixture Models
6.9.7	Unconditional Variance
6.9.8	A Baseline for $w_+(i)$
6.10	Summary

7.2	A General Derivation of AIC
7.3	General K-L–Based Model Selection: TIC
7.3.1	Analytical Computation of TIC
7.3.2	Bootstrap Estimation of TIC
7.4	AIC _c : A Second-Order Improvement
7.4.1	Derivation of AIC _c
7.4.2	Lack of Uniqueness of AIC _c
7.5	Derivation of AIC for the Exponential Family of Distributions
7.6	Evaluation of $\text{tr}(J(\theta_o)[I(\theta_o)]^{-1})$ and Its Estimation
7.6.1	Comparison of AIC Versus TIC in a Very Simple Setting
7.6.2	Evaluation Under Logistic Regression
7.6.3	Evaluation Under Multinomially Distributed Count Data
7.6.4	Evaluation Under Poisson-Distributed Counts
7.6.5	Evaluation for Fixed-Effects Normal Linear Models
7.7	Additional Results and Considerations
7.7.1	Selection Simulation for Nested Models
7.7.2	Simulation of the Distribution of Δ_p
7.7.3	Does AIC Overfit?
7.7.4	Can Selection Be Improved Based on All the Δ_i ?
7.7.5	Linear Regression, AIC, and Mean Squared Error
7.7.6	AIC _c and Models for Multivariate Data
7.7.7	There Is No True TIC _c
7.7.8	Kullback–Leibler Information Relation to Fisher Information Matrix
7.7.9	Entropy and Jaynes Maxent Principle
7.7.10	Akaike Weights w_i Versus Selection Probabilities π_i
7.8	Kullback–Leibler Information Is Always ≥ 0
7.9	Summary

8 Summary

8.1	The Scientific Question and the Collection of Data
8.2	Actual Thinking and A Priori Modeling
8.3	The Basis for Objective Model Selection
8.4	The Principle of Parsimony
8.5	Information Criteria as Estimates of Expected Kullback–Leibler Information
8.6	Ranking Alternative Models

8.9	MMI: Model Selection Uncertainty
8.10	MMI: Relative Importance of Predictor Variables
8.11	More on Inferences
8.12	Final Thoughts

References

Index

About the Authors

Drs. Kenneth P. Burnham and David R. Anderson have worked together for the past 28 years and have jointly published 9 books, 100 abstracts, 100 graphs and 66 journal papers on a variety of scientific issues. They are both in the Colorado Cooperative Fish and Wildlife Research Unit at Colorado State University, where they conduct research, teach, and mentor graduate students.

Ken Burnham has a B.S. in biology and M.S. and Ph.D. in statistics. For 29 years post-Ph.D. he has worked as a statistician, developing statistical theory in several areas of life sciences, including ecology and wildlife, most often in collaboration with subjects in ecology. Ken has worked (and lived) in Oregon, Alaska, Maryland (Patuxent Wildlife Research Center), North Carolina (U.S. Department of Agriculture, Wildlife Research Center), North Carolina (U.S. Department of Agriculture, Wildlife Research Center), North Carolina (U.S. Department of Agriculture, Wildlife Research Center), and North Carolina (U.S. Department of Agriculture, Wildlife Research Center). He is the recipient of numerous professional awards, including the Distinguished Achievement Medal from the American Statistical Association and Distinguished Statistical Ecologist Award from the International Congress of Ecology). Ken is a Fellow of the American Statistical Association.

David Anderson received B.S. and M.S. degrees in wildlife management and a Ph.D. in theoretical ecology. He is currently a Senior Scientist in the Biological Resources Division within the U.S. Geological Survey, where he is in the Department of Fishery and Wildlife Biology. He has worked at the Patuxent Wildlife Research Center in Maryland and 9 years at the Utah Cooperative Wildlife Research Unit and professor in the Department of Biology at Utah State University. He has been at Colorado State University since 1984. He is the recipient of numerous professional awards and academic contributions, including the Meritorious Service Award from the U.S. Department of the Interior.

This page intentionally left blank

Glossary

Notation and abbreviations generally used are given below. For specific examples can be found in those sections.

AIC	Akaike's information criterion.
AIC_{min}	The estimate of relative, expected K-L information for the best model in the set, given the data. For example, if the models are g_1, g_2, \dots, g_R and the data x , if the information is minimized for model g_6 , then $min = 6$, where min is the minimum over AIC_1, \dots, AIC_R . This is a random variable over samples. This notation also applies to $AIC_c, QAIC_c$.
AIC_{best}	In any set of models, one will be the best model, hence the actual best AIC model. The model for which $E_f(AIC)$ is minimized is denoted by the min is a random variable (like $\hat{\theta}$), $best$ is a value can be determined using Monte Carlo. The “best” model is the same model over a large number of samples (of which we have only a single sample). This notation also applies to $AIC_c, QAIC_c$, and TIC.
AIC_c	A second-order AIC, necessary for small samples.
Akaike weights	The relative likelihood of the model, given the data, normalized to sum to 1, are denoted by w_i probabilities.

	the data. Such a best model can be found by the methods and represents a statistical expectation. We consider the set $E(AIC_i)$, where $i = 1, 2, \dots, R$ model where $E(AIC_i)$ is minimized is AIC , AIC_c , $QAIC_c$, or TIC could be used.
Bias	(of an estimator) $Bias = E(\hat{\theta}) - \theta$.
BIC	Bayesian information criterion (Akaike 1978), also termed SIC in some literature.
c	A simple variance inflation factor used in GLS methods where there is overdispersion (extra binomial variation).
Δ_i	AIC differences, relative to the smallest AIC of the R models. Hence, AIC values are rescaled by a constant such that the model with the minimum AIC has $\Delta_i = 0$. Formally, $\Delta_i = AIC_i - AIC_{min}$ are estimates of the expected K-L information difference between the selected (best) model and the other models. The differences apply to AIC, AIC_c , $QAIC_c$, and TIC .
Δ_p	A “pivotal” value, analogous to $(\theta - \theta_0) / \sqrt{\text{var}(\hat{\theta})}$. $\Delta_p = (AIC_i - AIC_{best}) / (AIC_{best} - AIC_{min})$.
df	Degrees of freedom as associated with hypothesis testing. df is the difference between the number of parameters in the null and alternative hypotheses in standard tests.
$E(\hat{\theta})$	The statistical expectation of the estimator.
Estimate	The computed value of an estimator, given a set of sample data (e.g., $\hat{\theta} = 9.8$).
Estimator	A function of the sample data that is used to estimate a parameter. An estimator is a random variable, denoted by a “hat” (e.g., $\hat{\theta}$).
Evidence ratio	The relative likelihood of model i versus model j , $\mathcal{L}(g_i data)/\mathcal{L}(g_j data)$, which is identical to $e^{\Delta_i - \Delta_j}$.
$f(x)$	Used to denote “truth” or “full reality” where f produces multivariate data x . This concept is often considered to be an infinite-dimensional space.
$g_i(x)$	Used to denote the set of candidate models g_i that are sized to provide an adequate approximation of empirical data. The expression $g_i(x data)$ is necessary to clarify that the function involves the data.

K candidate models is represented as \mathbf{g} . Also, $\hat{g}_i = g_i(x|\hat{\theta})$.

Global model	A highly parameterized model containing all associated parameters thought to be important. Without an a priori consideration of the problem at hand, is a global model, all other models in the set are subsets of this global model.
K	The number of estimable parameters in the global model.
K-L	Kullback–Leibler distance (or discrepancy number).
LRT	Likelihood ratio test.
LS	Least squares method of estimation.
$\mathcal{L}(\theta x, g)$	Likelihood function of the model parameters θ given x and the model g .
$\mathcal{L}(g_i x)$	The discrete likelihood of model g_i , given x .
$\log(\cdot)$	The natural logarithm (\log_e).
$\text{logit}(\theta)$	The logit transform: $\text{logit}(\theta) = \log\left(\frac{\theta}{1-\theta}\right)$, $0 < \theta < 1$.
g_i	Shorthand notation for the candidate model g_i .
\min	An index to denote the fitted model that minimizes the information criterion, given the data. The model selected, based on minimizing the information criterion, given the data. AIC, AIC_c , QAIC_c , etc., are used in this context.
ML	Maximum likelihood method of estimation.
MLE	Maximum likelihood estimate (or estimator).
n	Sample size. In some applications there may be a relevant sample size (e.g., in random effects models).
Parsimony	The concept that a model should be as simple as possible concerning the included variables, model parameters, and number of parameters. Parsimony is a desirable property of a model used for inference, and it is usually associated with a tradeoff between squared bias and variance of the estimators. Parsimony lies between the extremes of under-fitting and over-fitting.
Precision	A property of an estimator related to the variability of the estimates among estimates from repeated samples.

π_i	Model selection probabilities (or relative quasi-likelihood adjustments are required from Monte Carlo studies or the bootstrap).
R	The number of candidate models in the set. One of these models is the estimated best model (in the sense of a specific model $g(x \hat{\theta})$, where the parameters have been estimated) for the data. Another model (possibly the same model) is the best model (g_{best}) to use as a basis for inference.
τ_i	Prior probability of model i . Also used to denote model redundancy (Section 4.6).
θ	Used to denote a generic parameter vector. Also used to denote conditional survival probabilities S_i .
$\hat{\theta}$	An estimator of the generic parameter θ .
θ_0	The optimal parameter value in a given model for a given sample size, but ignoring estimation issues. This is the value that minimizes K-L information for a given model structure.
TIC	Takeuchi's information criterion.
w_i	Akaike weights. Used with any of the information criteria (AIC, AIC _c , QAIC, TIC). The w_i are estimates of expected Kullback–Leibler information. The w_i sum to 1 and are interpreted as the probability that model i is the actual best model for the sampling situation considered.
$w_+(j)$	Sum of Akaike weights over all models containing explanatory variable j . These sums are used in model selection problems where one wants to assess the relative importance of the explanatory variables. They are also used to estimate the importance of the explanatory variables in models that are robust to model selection.
χ^2	A test statistic distributed as chi-square with df degrees of freedom. Used here primarily for a goodness-of-fit test of the global model for the data.
\approx	Approximately equal to.
\sim	Distributed as.

Introduction

1.1 Objectives of the Book

This book is about making valid inferences from scientific data. A useful analysis depends on a model of the information in the data. The objective is to provide scientists, including statisticians, with practical advice for the analysis of empirical data under a general theoretic paradigm. We first assume that an exciting scientific question has been carefully posed and relevant data have been collected, following a good experimental design or probabilistic sampling program. Alternative models and models to represent them, should be carefully considered at an early stage of the investigation. Often, little can be salvaged if the data have been seriously flawed or if the question was poorly posed. We acknowledge, of course, that these issues are never as ideal as one would like. Proper attention must be placed on the collection of data (Cox and Anderson 2001). We stress inferences concerning the structure of biological systems, relevant parameters, valid measures of effect, and formal prediction.

There are many studies where we seek an understanding of the underlying mechanism, especially causal ones. There are many studies to understand the relationships between variables that are important because of the parameters in them and relationships between and among variables. These parameters have real interpretations, even when they relate to quantities that are not directly measurable (e.g., survival probabilities, animal density in an area, genetic interaction terms). Science would be very limited without

ters. Observables and prediction are often critical, but so are these issues.

The first objective of this book is to outline a consensus surrounding the analysis of empirical data. Inductive statistical inference about a defined population or process from a sample or experimental data set. “Data analysis” leading to a model is the integrated process of careful a priori model formulation, parameter estimation, and measurement of precision (a major component due to model selection uncertainty). We do not think model selection should be treated as an activity that precedes model selection; model selection is a critical and integral aspect of scientific inference that leads to valid inference.

A philosophy of thoughtful, science-based, a priori model selection. Often, one first develops a global model (or set of models) and then several other plausible candidate (sub)models postulated to approximate information in the data at hand. This forms the basis of *models*. Science and biology play a lead role in this a priori model selection. Careful consideration of the problem. A simple example of alternative scientific hypotheses might be helpful at this point. Consider the importance of an interaction between age (a) and weight (w) in a particular animal population. A model including such an interaction would have the main effects plus the interaction; $a + w + a * w$. A model without the interaction term, $a + w$, lacks the interaction term. Information-theoretic methods provide lines of quantitative evidence concerning the importance of the interaction.

The modeling and careful thinking about the problem of model selection that have often received relatively little attention in statistics (and, for nonmajors), partly because such classes rarely consider the philosophy of data analysis. A proper a priori model-building process to avoid “data dredging,” which leads to overfitted models and “discovery” of effects that are actually spurious (Anderson 2001). There has often been a rush to “get to the data analysis” and begin interpreting the data and compute various estimates of interest or conduct hypothesis tests. We realize that these other philosophies may have their place in more exploratory investigations.

The second objective is to explain and illustrate methods of model selection at the interface of information theory and mathematical statistics. Selection of an estimated “best approximating model” from the a priori set of models. In particular, we review and explain the use of Akaike’s information criterion (AIC) in the selection of a model (or small set of models) for statistical inference. AIC provides a simple, effective, and principled method for the selection of an estimated “best approximating model” for inference. Model selection includes “variable selection”

quite different from the usual methods based on null hypothesis selection based on information theory is not the only reason why it is what we are focusing on here because of its philosophical advantages.

The practical use of information criteria, such as Akaike's information criterion is relatively recent (the major exception being in ecology where AIC has been used routinely for the past two decades). Information theory and mathematical statistics started with a book. Akaike considered AIC to be an extension of R. A. Fisher's theory. These are all complex issues, and the literature is often scattered widely throughout books and research journals. We bring this relatively new material into a readable text for (primarily) the biological and statistical sciences. We provide a number of examples, many of which are biological, to illustrate various aspects of its application.

In contrast, hypothesis testing as a means of selecting models has had a much longer exposure in science. Many seem to feel more comfortable with the hypothesis testing paradigm in model selection, and often judge the results of a test as *the* standard by which other approaches should be judged (we believe that they are wrong to do so). Bayesian model selection and inference have been the focus of much recent research. At the technical level of this material often makes these approaches less accessible to many in the biological sciences. A variety of cross-validation based methods have been proposed for model selection, and many of them are like very reasonable approaches. The computational demands of Bayesian and cross-validation methods for model selection are often (often 1–3 orders of magnitude higher than information-theoretic methods), especially if there are more than a dozen or so high-dimensional models.

The theory presented here allows estimates of “model selection” inference problems that arise in using the same data for model selection and the associated parameter estimation and inference. In model selection, uncertainty is ignored, precision is often overestimated, and interval coverage is below the nominal level, and prediction error rate is higher than expected. Another problem is the inclusion of model selection or factors, with no assessment of the reliability of their selection. Several methods for dealing with model- and variable-selection are suggested and examples provided. Incorporating model selection into estimators of precision is an active area of research, and several additional approaches developed in the coming years.

The third objective is to present a number of approaches to model selection and inference from more than one model in the set. That is, to compare inferences from only the model estimated to be the best, to

research area for Bayesians for the past several years (1990-2000). Model averaging can be easily done under an information-theoretic approach. Model averaging has several practical and theoretical advantages in prediction or in cases where a parameter of interest occurs in several models. Confidence sets on models is another useful approach. Models in the set represent a logical ordering (e.g., a set of models representing chronic treatment effects over 1, 2, . . . , t time periods). Model averaging can be used to assess the relative importance of explanatory variables in a general regression model. Model averaging can be assessed by summing certain quantities across models. Model averaging is useful in certain conflict resolution issues (Anderson et al. 2000).

Current practice often would judge a variable as important or not important based on whether that variable was in or out of the selected model (e.g., based on hypothesis testing). Such procedures are based on a dichotomy (see Breiman 2001) and are not in the spirit of a model selection approach. MMI allows us to discard simplistic dichotomies and focus on ranking models and variables as to their relative value and importance.

Modeling is an art as well as a science and is directed at finding a good approximating model of the information in empirical data for statistical inference from those data. In particular, the model parameters estimated from data should be substantially less than the number of parameters or inference is likely to remain somewhat preliminary (e.g., Anderson 1998 mentions a regression problem with 757 variables and a sample size of 1000). It is absurd to think that valid inference is likely to come from such a situation (data). In cases where there are relatively few data per estimated parameter, a small-sample version of AIC is available (termed AIC_c) and is used routinely rather than AIC. There are cases where quasi-likelihood approaches are appropriate when count data are overdispersed; this theory leads to criteria such as QAIC and $QAIC_c$, and these extensions are discussed in the following material.

Simple models with only 1-2 parameters are not the focus of this book; rather, we focus on models of more complex systems. The information has been firmly considered to be an optimization problem for decades, and AIC formulates the problem of model selection as an optimization problem across a set of candidate models. Minimizing the Akaike information operation with results that are easy to interpret. Models can be compared and scaled, allowing full consideration of other good models. The estimated “best approximating model.” Evidence ranking models by the strength of evidence for alternative hypotheses. Competing models with AIC values close to the minimum, are also useful in the presence of model selection uncertainty. Inference should often be based on the best model, unless the data clearly support only a single model. In some approaches are provided to allow inference from multiple models, including model averaging.

1.2 Background Material

Data and stochastic models of data are used in the empirical inferences concerning both processes and parameters of interest (e.g., Box and Jenkins 1976, Akaike 1973, Box and Jenkins 1981, Lunneborg 1994, and Shenk and Franklin 2001 for principles). Statistical scientists have worked with researchers in the sciences for many years to improve methods and understand complex processes. This book provides practical, omnibus methods for inference from models that are good approximations to biological data. We focus on statistical evidence and try to avoid ambiguity such as “significant or not significant.” A broad definition of a data set is given here. A single, simple data set might be the subject of analysis. Data collected from several field sites or laboratories are treated in a comprehensive analysis. The data might commonly be categorized by age, sex, species, treatment group, or within several geographic areas. In linear and nonlinear regression models, we use explanatory variables. There are often factors (variables) with small and large effects in these information-rich data sets (the effect sizes). Parameters in the model represent the effects of the factors. We focus on modeling philosophy, model selection, estimation of parameters, and valid measures of precision under the relative information-theoretic methods. Valid inference rests upon the model. In addition to the critical considerations relating to problem formulation, design, and protocol for data collection.

1.2.1 *Inference from Data, Given a Model*

R. A. Fisher (1922) discussed three aspects of the general theory of inference: (1) model specification, (2) estimation of model parameters, (3) estimation of precision. Here, we prefer to partition the problem into two components: formulation of a set of candidate models and selection of a model (or small number of models) to be used in making inferences. In much of the twentieth century, methods have been available to efficiently estimate model parameters and their precision (e.g., standard error, covariance matrix). Fisher’s *likelihood theory* has been the dominant approach to these issues, but it *assumes* that the model is correct (and correct, i.e., a true model) and that only the parameters

of a log-linear model for the analysis of count data displayed in a table. The parameters in these models can be estimated using maximum likelihood (ML) methods. That is, if one assumes or somehow knows the appropriate model, methods exist that are objective and asymptotically efficient for estimating model parameters and the sampling covariance structure of that model. A more challenging example might be to assume that a count is appropriately modeled by a 3-parameter gamma distribution. One can use the method of maximum likelihood to estimate these parameters, and the model-based 3×3 sampling covariance matrix. Confidence intervals for the model, and if the sample size is “large,” then maximum likelihood estimators of parameters that are consistent (i.e., asymptotically unbiased with variance tending to zero), fully efficient (i.e., minimum variance unbiased estimators), and normally distributed. With small samples, assuming an appropriate model, ML estimators often have bias, where $\text{bias} \equiv E(\hat{\theta}) - \theta$. Such bias is usually a trivial concern, but it can often be substantially less than the $\text{se}(\hat{\theta})$, and bias-adjusted confidence intervals can often be found if this is deemed necessary. The sampling distribution of ML estimators are often skewed with small samples, but probability density intervals or log-based intervals or bootstrap procedures can provide asymmetric confidence intervals with good coverage probability. **the maximum likelihood method provides an objective method for estimation of model parameters and the sampling covariance matrix given an appropriate model.**

1.2.2 Likelihood and Least Squares Theory

Biologists have typically been exposed to least squares regression in statistics classes in applied statistics. LS methods for linear models are easy to compute, and therefore they enjoyed an early history of use (e.g., Fisher 1935, Fisher 1985). In contrast, Fisher’s likelihood methods of estimation and numerical methods and were thus not popular prior to the availability of personal computers and the development of easy-to-use software. Least squares theory has many similarities with likelihood theory, and it is easy to estimate parameters of the structural parameters (but not σ^2) for linear models when the residuals are assumed to be independent and normally distributed. It is now easy to allow alternative error structures (i.e., nonnormal distributions such as Poisson, gamma or log-normal) for regression and other models using either a likelihood or quasi-likelihood framework (e.g., McCullagh and Nelder 1989, Heyde 1997), but more difficult in an LS framework.

The concepts underlying both estimation methods are the same. To understand (Silvey 1975). Consider the simple linear regression where the response variable (y) is modeled as a linear function of an

variance σ^2 . Under LS the estimates of p_0 and p_1 are $\sum(\epsilon_i)^2$ —hence the name *least squares*. The parameter σ^2 is estimated by minimizing the average squared error terms (ϵ_i) and define $\hat{\sigma}^2$ as the “best fit.” Hundreds of statistics books cover the theory and practice for least squares estimation in linear and nonlinear models. In this book the ϵ_i are assumed to be independent, normally distributed random variables.

Likelihood methods are much more general, far less taught in statistics courses, and slightly more difficult to understand at first. Much of this book relies on an understanding of likelihood methods. A brief introduction is given here. While likelihood theory is common to both frequentist and Bayesian statistics, there are no books of applied books solely on this important subject (good examples are Cullagh and Nelder 1989, Edwards 1992, Azzalini 1996, and Severini 2000).

The theory underlying likelihood begins with a probability distribution for the parameters (θ). Specifically, model g describes the probability of the data, given the model parameters and a specific model $g(x|\theta, model)$. A simple example is the binomial probability distribution. θ is the probability of a “success”; let this be the parameter of interest. Suppose we observe $y = 15$ successes out of $n = 40$ trials. Then, the discrete probability of getting 15 successes out of 40 trials, given parameter ($p \equiv 0.4$) and the binomial model, is

$$g(y, n|p, binomial) = \binom{n}{y} p^y (1 - p)^{n-y},$$

$$g(15, 40|p = 0.4, binomial) = \frac{40!}{15!25!} (0.4)^{15} (1 - 0.4)^{25}$$

The key point is that for this calculation, the model (here binomial) and its parameters (here $p = 0.4$) are known in advance. In very simple problems such as this, an excellent model can be considered given (such is rarely the case in the real world, where we are not sure what model might be used). Then one observes the data ($n = 40$) and can compute the probability of the data, given the model and its parameters.

In much of science, neither the model parameters nor the model itself are known in advance. However, data can be collected in a way that allows the parameters to be estimated if a good model can be found or assumed. The likelihood method is the basis for such parameter estimation and is a function of the data, y , and the model, g , given the data and the binomial model:

$$\mathcal{L}(p|y, n, binomial) = \binom{n}{y} p^y (1 - p)^{n-y}$$

Clearly, the likelihood is a function of (only) the unknown parameters (in this example); everything else is known or assumed. The probability model and the associated likelihood function differ only in terms of what is known and given. In the probability model, the parameters, the model, and the data are known, and interest lies in the probability of observing the data, y given n in this simple example). In the likelihood function, the data are given (observed) and the model is assumed (but unknown); interest lies in estimating the unknown parameters; thus, the likelihood function is a function of only the parameters. The probability model of the data and the likelihood function of the parameters are closely related; they merely differ in terms of what is known and given. The binomial probability model does not contain the unknown parameter p and is often used to estimate p (it contains any information about the unknown parameters p that can be used to compute it if $n > 50$).

The notation for the likelihood function is very helpful. We can consider the general expression $\mathcal{L}(\theta|data, model)$. If we follow the convention of letting x represent the empirical data and g a general model, then $\mathcal{L}(\theta|x, g)$ is read as “the likelihood of a particular value of the unknown parameter θ (θ is usually a vector), given the data x and the particular model g .”

A well-known example will help illustrate the concept. Suppose we flip n pennies and observe y “heads.” Assuming that the flips are independent and that each penny has an equal probability of a head, the binomial model is an obvious model choice in this simple setting. The likelihood function is $\mathcal{L}(p|y, n, \text{binomial})$, where p is the (unknown) probability of a head given the data (y and n) and the binomial model, one can compute the likelihood that p is 0.15 or 0.73 or any other value between 0 and 1. The likelihood function (a relative, not absolute, value) is a function of the unknown parameter p . Given this formalism, one might compute the likelihood function for various values of the unknown parameter p and pick the most likely one as the maximum likelihood estimate of p , given the data and the model. It seems compelling to call this the “most likely.” This is Fisher’s concept of *maximum likelihood estimation*; he published this when he was 22 years old as a third-year student at Cambridge University! He reasoned that the best estimate of the unknown parameter (given data and a model) was that which was most likely to have produced the data. The name *maximum likelihood*, ML. The ML estimate of the unknown parameter in the binomial model happens to have a closed-form expression $\hat{p} = y/n = 7/11 = 0.6363$. That is, the numerical value of \hat{p} maximizes the likelihood function. In most real-world cases, the ML estimate in closed form either does not exist or cannot be found, making the ML estimation difficult.

hypothesis tests, the basis for profile likelihood intervals, quantities (such as model selection criteria). More generally, it includes the broad concept of *support* (Edwards 1992) as the essential basis for Bayesian approaches to statistical inference. The likelihood is the backbone of statistical theory, whereas the log-likelihood is viewed as a limited special case and, while very useful in many applications, is not foundational in modern statistics.

For many purposes the natural logarithm of the likelihood is used; written as $\log(\mathcal{L}(\theta|data, model))$, or $\log(\mathcal{L}(\theta|x, model))$. It is clear, just $\log(\mathcal{L}(\theta))$ or even just $\log(\mathcal{L})$. Often, one sees $\log(\mathcal{L}(\theta|x))$, without it being clear that a particular model is intended. An advanced feature of $\log(\mathcal{L})$ is that it, by itself, is a type of *information*. The log-likelihood θ and the model (Edwards 1992:22–23). The log-likelihood model where 11 pennies are flipped and 7 heads are observed.

$$\begin{aligned}\log(\mathcal{L}(p|y, n, binomial)) &= \log \binom{n}{y} + y \cdot \log(p) + (n - y) \cdot \log(1 - p) \\ &= \log \binom{11}{7} + 7 \cdot \log(p) + (11 - 7) \cdot \log(1 - p) \\ &= 5.79909 + 7 \cdot \log(p) + (4) \cdot \log(1 - p)\end{aligned}$$

A property of logarithms for values between 0 and 1 is that the log is negative; thus, values of discrete log-likelihood functions are in the negative quadrant; thus, values of discrete log-likelihood functions are in the negative quadrant (unless some additive constants have been omitted). Figure 1.1 shows the likelihood (a) and log-likelihood (b) functions for the binomial distribution. In the example, 11 pennies were flipped, 7 heads were observed, and the binomial model with $p = 0.636$ maximizes both the likelihood and the log-likelihood function; this value is denoted by \hat{p} and is the maximum likelihood estimate (MLE). Relatively little information is contained in the log-likelihood function for $n = 11$ and this is reflected in the broad shape of the log-likelihood function. If the sample size had been 5 times larger, with $n = 55$ and 35 heads, the likelihood and log-likelihood functions would be more sharply peaked (c and d). In fact, the sampling variance is derived from the curvature of the log-likelihood function around its maximum point. In the usual case, a variance–covariance matrix can be estimated based on the second derivatives of the log-likelihood function. These procedures will be discussed here.

The value of the log-likelihood function at its maximum is an important quantity, and it is this point that defines the *maximum likelihood estimate*. In the example with 11 flips and 7 heads, the value of the log-likelihood is -1.411 (Figure 1.1b). This result is common to all maximum likelihood estimates.

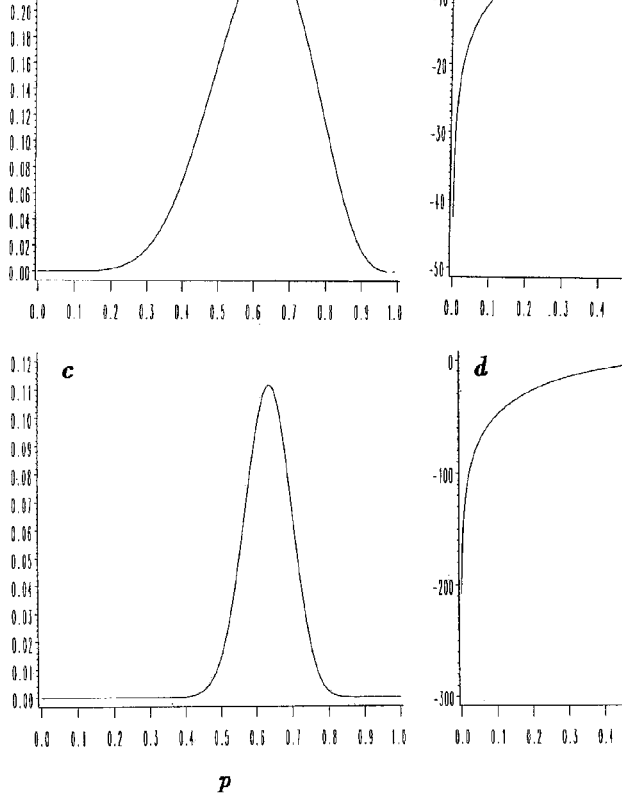


FIGURE 1.1. Plots of the binomial likelihood (a) and log-likelihood (b) function for $n = 11$ penny flips and the observation that $y = 7$ of these were heads; (c) and (d) are the same plots but for $n = 55$ penny flips and the observation that $y = 35$ of these were heads. The differing scales on the Y axis.

log-likelihood function

$$\log(\mathcal{L}(p|y, n, \text{binomial})) = \log \binom{n}{y} + y \cdot \log(p) + (n - y) \cdot \log(1 - p)$$

and substituting the MLE ($\hat{p} = 0.6363$) and the data ($y = 7$) we get

$$-1.411 = 5.79909 + 7 \cdot \log(0.6363) + (4) \cdot \log(0.3637)$$

Thus, when one sees reference to a maximized $\log(\mathcal{L}(\theta))$ that is a numerical value (e.g., -1.411).

Many do not realize that the common procedure for setting a confidence interval (i.e., $\hat{\theta} \pm 1.96 \cdot \widehat{\text{se}}(\hat{\theta})$) is merely an approximation

achieved confidence interval coverage can be much less than say, 95%). For example, if the binomial parameter is near 0, the estimator $\hat{\theta}$ will be nonnormal (asymmetric) unless n is very large. In general, rather than use the simple approximation to a 95% interval using the log-likelihood function; this procedure is called a *profile likelihood interval*. This is not a simple approximation has seen heavy use in applied data analysis. We discuss the full theory for profile likelihood intervals here, but will not do so in the binomial case where $n = 11$, $y = 7$, $\hat{p} = 0.6363$, and the log-likelihood value is -1.411 . Here, we start with 3.84, which is the chi-squared distribution with 1 degree of freedom. One minus 1.92, and this value is subtracted from the maximum point of the log-likelihood function: $-1.411 - 1.92 = -3.331$. Now, numerically, we find the values of p that are associated with the values of the log-likelihood function of -3.331 . These 2 values are the endpoints of an exact 95% likelihood interval. In this example, the 95% likelihood interval is (0.41, 0.84).

Biologists familiar with LS but lacking insight into likelihood estimation can benefit from an example. Consider a multiple linear regression model where the dependent variable y is hypothesized to be a function of r (regressor) variables x_j ($j = 1, 2, \dots, r$). Here the residuals ϵ_i are assumed to be independent, normally distributed with mean 0 and variance σ^2 , and the model structure is expressed as

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_r x_{ir} + \epsilon_i, \quad i = 1, \dots, n$$

Hence

$$E(y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_r x_{ir}, \quad i = 1, \dots, n$$

and $E(y_i)$ is a linear function of $r + 1$ parameters. The conditional distribution of y_i given x_{i1}, \dots, x_{ir} is

$$\epsilon_i = y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_r x_{ir}) = y_i - E(y_i)$$

have the joint probability distribution $g(\underline{\epsilon}|\underline{\theta})$, where $\underline{\theta}$ is a vector of parameters $(\beta_0, \beta_1, \dots, \beta_r, \text{ and } \sigma)$. Here, corresponding to the model

$$g(\epsilon_i|\underline{\theta}) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left[\frac{\epsilon_i}{\sigma}\right]^2}.$$

The likelihood is simply the product of these over the n observations. Interpreted as a function of the unknown parameters, given the data and the model structure, and the normality assumption:

$$\mathcal{L}(\underline{\theta}|\underline{x}) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left[\frac{\epsilon_i}{\sigma}\right]^2} = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^n e^{-\frac{1}{2\sigma^2} \sum_{i=1}^n \epsilon_i^2}$$

of $\underline{\hat{\theta}}$ is identical to the usual LS regression estimators (the estimator of σ^2 differs slightly). This formalism shows, *given* the data, the model, and the parameters to be estimated, using either LS or ML.

In all fitted linear models the residual sum of squares (RSS)

$$\text{RSS} = \sum_{i=1}^n \hat{\epsilon}_i^2,$$

where

$$\begin{aligned}\hat{\epsilon}_i &= y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_r x_r) \\ &= y_i - \hat{E}(y_i)\end{aligned}$$

The ML estimator is $\hat{\sigma}^2 = \text{RSS} / n$, while the estimator in the LS case is $\hat{\sigma}^2 = \text{RSS} / (n - (r + 1))$. This shows that ML and LS estimators of σ^2 differ by a factor of $n / (n - (r + 1))$; often a trivial difference if the sample size is small. The maximized likelihood is

$$\mathcal{L}(\hat{\underline{\theta}}|x) = \left[\frac{1}{\sqrt{2\pi\hat{\sigma}^2}} \right]^n e^{-\frac{1}{2}n},$$

or

$$\log(\mathcal{L}(\hat{\underline{\theta}})) = -\frac{1}{2}n \log(\hat{\sigma}^2) - \frac{n}{2} \log(2\pi) - \frac{1}{2}n.$$

The additive constants can often be discarded from the log-likelihood, since they are constants that do not influence likelihood-based inference. In all standard linear models, we can take

$$\log(\mathcal{L}(\hat{\underline{\theta}})) \approx -\frac{1}{2}n \log(\hat{\sigma}^2)^2.$$

This result is important in model selection theory because it allows mapping from LS analysis results (e.g., the RSS or the ML estimator) to the maximized value of the log-likelihood function for comparing different linear models with normal residuals. Note that the log-likelihood is only defined to an arbitrary additive constant in this usual case. If the residuals are not normally distributed (e.g., gamma, and log-normal), then all the terms in the log-likelihood must be retained, without omitting any constants. Most uses of the log-likelihood are relative to its maximum, or to other likelihoods at their maximum, or to the curvature of the log-likelihood function at the maximum.

The number of parameters $K = r + 2$ in these linear models includes the intercept (say, β_0), the r regression coefficients (β_1, \dots, β_r), and the variance (σ^2). Often, one (erroneously) considers only the $r + 1$ parameters being estimated as the intercept and the slope parameters.



Sir Ronald Aylmer Fisher was born in 1890 in East Finchley, London, and died in Australia, in 1962. This photo was taken when he was approximately 70. He was one of the foremost scientists of his time, making incredible contributions to pure and applied statistics and genetics. Details of his life and many scientific achievements are found in Box (1978). He published 7 books (one of these has been translated and printed in 7 languages) and nearly 300 journal papers. Most relevant to this book is Fisher's likelihood theory and parameter estimation using the maximum likelihood.

however, in the context of model selection, the number of parameters to estimate include σ^2 and thus $K = r + 2$. If the method of LS is used for the parameter estimators, one must use the regression-based estimator of σ^2 , $\hat{\sigma}^2 = \text{RSS}/(n - (r + 1)) = \text{RSS}/(n - K + 1)$ to obtain the ML estimator of σ^2 . For estimation, we minimize $\text{RSS} = n\hat{\sigma}^2$, which for all parameter values, itself is equivalent to maximizing $-\frac{1}{2} \cdot n \log(\hat{\sigma}^2)$.

There is a close relationship between LS and ML methods for linear and nonlinear models, where the ϵ_i are assumed to be normally distributed. For example, the LS estimates of the structural model parameters are equivalent to the MLEs. Likelihood (and related Bayesian) methods have extensions to the many other classes of models and, with the advent of modern computing equipment, likelihood methods are finding increasing use among statisticians and researchers in other scientific disciplines (Box and Jenkins, 1995 for background).

1.2.3 The Critical Issue: “What Is the Best Model?”

While hundreds of books and countless journal papers deal with the estimation of model parameters and their associated precision, relatively little attention is given to the question of model selection.

In fact, Fisher believed at one time that model specification was the field of mathematical statistics, and this attitude prevailed in the statistical community until at least the early 1970s. “*What is the critical question in making valid inference from data?*” is the critical question in making valid inference from data in the natural sciences.

The likelihood function $\mathcal{L}(\theta|x, model)$ makes it clear that θ , data and the model are taken as *given*. Before one can consider that $\theta = 5.3$, one must have data and a particular statistical model. If an investigator will have empirical data for analysis, it is the model that is known or given. Rather, a number of alternative models are somehow considered as well as the specific explanatory model. This is the *problem* in multiple regression analysis. If one has data, ML theory can be used to estimate the unknown parameters. What quantities useful in making statistical inferences. However, what is the best to use for making inferences? What is the basis for the “best”?

Model selection relates to fitted models: given the data and a model, then the MLEs of the model parameters have been estimated. Inference relates to theoretical models. It is necessary to

- (1) models as structure only (θ value irrelevant),
- (2) models as structure, plus specific θ_o (this is the theoretical model),
- (3) models as structure, plus MLE $\hat{\theta}$, fitted to data,
- (4) models as structure by fitting, downplaying θ .

If a poor or inappropriate model (3, above) is used, then the model is poor on the data and this model will often be poor. Thus, it is not possible to select (i.e., infer) an appropriate model (1, above) for the data set; however, this is not the same as trying to find the best model. Model selection methods with a deep level of theoretical support, particularly, methods that are easy to use and widely applicable. The concept of “applicability” means that the methods have good operating characteristics for realistic sample sizes. As Potscher (1991) noted, asymptotic theory has little value unless they hold for realized sample sizes.

A simple example will motivate some of the concepts. Bystrak (1992 and 1996) studied patterns of avian species-accumulation curves in forested landscapes in the eastern United States using data from the Breeding Bird Survey (Bystrak 1981). He derived an a priori set of 9 models from two sources: (1) the literature on species-area relationships (the power or exponential models were suggested) and (2) the literature for functions that increased monotonically to a plateau (1.1). Which model should be used for the analysis of the data? Clearly, none of these 9 models are likely to be the “true” model.

response variable (y) is the number of accumulated species, and x is the accumulated number of samples. Nine models and their number of parameters are shown to motivate the question, “Which fitted model should be used from these data?”

Model structure	Number of parameters (K)
$E(y) = ax^b$	3
$E(y) = a + b \log(x)$	3
$E(y) = a(x/(b + x))$	3
$E(y) = a(1 - e^{-bx})$	3
$E(y) = a - bc^x$	4
$E(y) = (a + bx)/(1 + cx)$	4
$E(y) = a(1 - e^{-bx})^c$	4
$E(y) = a \left(1 - [1 + (x/c)^d]^{-b} \right)$	5
$E(y) = a[1 - e^{-(b(x-c)^d)}]$	5

^aThere are $K - 1$ structural parameters and one residual variance parameter, $c(\epsilon) = 0$, $V(\epsilon) = \sigma^2$.

the index data from the Breeding Bird Survey over the years 1966–1999. Flather wanted an approximating model that fit the data well and was useful in making inferences about bird communities on the scale of the entire continent. In this first example, the number of parameters in the candidate models varied only from 3 to 5. Which approximating model is “best” for these data? From these data is answered philosophically by the principle of parsimony (Section 1.4) and operationally by several information-theoretic methods (Chapter 2. Methods for estimating model selection uncertainty). The results of rating this into inferences are given in Chapter 2 and illustrated in Figures 2.1 and 5.

Note, in each case, that the response variable y is being modeled directly, rather than mixing models of y with $\log(y)$, or other transformation of y . The response variable (Table 1.1). These models are in the sense of 1 above. The model structure is given but the parameter values are unspecified. Given a set of data, the model can be used to obtain $\hat{\theta}$ in the sense of 3 above. In some of the cases, the model parameters are derived from theory, without the need for fitting to specific empirical data. Such cases seem to be the exception rather than the rule in sciences, where model parameters must usually be estimated from data using least squares or likelihood theory.

1.2.4 Science Inputs: Formulation of the Set of Candidate Models

Model specification or formulation, in its widest sense, is the process of identifying the difficult than estimating the model parameters and their probabilities.

jective art, that is why scientists must be trained, educated in their discipline. The published literature and experimental sciences can be used to help formulate a set of a priori criteria. The most original, innovative part of scientific work is the posing of a proper question. Good approximating models, each representing a hypothesis, in conjunction with a good set of relevant data, are fitted into the underlying biological process and structure.

Lehmann (1990) asks, "where do models come from," and gives biological examples (also see Ludwig 1989, Walters 1996, Lin 1997). They arise from questions about biology and the manner in which biological systems function. Relevant theoretical and practical questions come from a variety of sources (see Box et al. 1978, O'Connor and Spotnitz 1991). Usually, these questions come from the scientific literature, research results, experiments, personal experience, or contemporary debate within the scientific community. More practical questions stem from real-world controversies, biomonitoring programs, quasi-experiments, and public hearings.

Chatfield (1995b) suggests that there is a need for more balance (than is usually evident) and a *better balance* between the posing of a question, analysis theory, and data. This suggestion has been in the literature for decades. One must conclude that it has not been emphasized enough in applied science or statistics courses. **Our science curriculum does not do enough to expect and enforce critical thinking.** Too much emphasis is focused on the analysis theory and data analysis, with little attention about the reason for the study in the first place (see Hayne 1991 for examples).

Tukey (1980) argues for the need for deep thinking and exploration in analysis, and that the results of these activities lead to good conclusions and confirmatory data analysis. In the exploratory phase, the importance of a flexible attitude and plotting of the data. He emphasizes the computation of test statistics, *P*-values, and so forth in the data analysis. Tukey concludes that to implement the conclusions properly we need to do a lot of exploratory work.

The philosophy and theory presented here must rest on a foundation of careful planning and execution of field or laboratory projects. Many books exist giving information on these important issues (Dunn 1979, Cook and Campbell 1979, Mead 1988, Hairston 1989, Dunn 1991, Eberhardt and Thomas 1991, Manly 1992, Skalski 1992, Thompson 1992, Scheiner and Gurevitch 1993, Cox and R. 1994, and Zimmermann 2000). Chatfield (1991) reviews statistical methods that these might be avoided. Research workers are urged to pay attention to these critical issues. Methods given here should not be used on poorly designed work. In the following material we will

Development of the a priori set of candidate models is a global model: a model that has many parameters, includes all relevant effects, and reflects causal mechanisms thought to be important in the *science of the situation*. The global model should also include all signs and attributes of the system studied. Specification of the global model should not be based on a probing examination of the data. At some early point, one should investigate the fit of the global model (e.g., examine residuals and measures of fit such as R^2 , χ^2 goodness-of-fit tests) and proceed with analysis only if the global model provides an acceptable fit to the data. If the global model parameters can then be derived as special cases of the global model, the set of reduced models represents plausible alternatives based on hypotheses hypothesized about the process under study. Generally, alternatives involve differing numbers of parameters; the number of parameters may differ by at least an order of magnitude across the set of candidate models. Chatfield (1995b) writes concerning the importance of developing candidate models: "Considerations such as accepted theory, expert background knowledge, and information in addition to known constraints on both the process and the variables in the models. All these factors should be used to make the makeup of the set of candidate models, prior to actual analysis."

The more parameters used, the better the fit of the model is likely to be achieved. Large and extensive data sets are likely to support a wide range of models, and this should be considered in the development of the set of candidate models. **If a particular model (parametrization) does not make sense, this is reason to exclude it from the set of candidate models, even if the fit is good, in the case where causation is of interest.** In developing candidate models, one must recognize a certain balance between knowledge of the process and focused on plausible hypotheses, while making it better to be conservative against omitting a very good a priori model. While this is a difficult task, considered, we advise the inclusion of all models that seem plausible. There is no justification, prior to data analysis. While one must work to avoid both underfitting and overfitting, it seems that model overfitting is more damaging than underfitting (Shibata 1989). We recommend a considerable amount of careful, a priori thinking in arriving at a set of candidate models (see Peirce 1955, Burnham and Anderson 1992, Chatfield 1995b).

Freedman (1983) noted that when there are many, many variables (x_1, x_2, \dots, x_{50}) used to predict a response variable, many selection methods will provide regression equations with many "significant" F values, and many "significant" regression coefficients shown by large t values, *even if the explanatory variables are unrelated to y* . This undesirable situation occurs most frequently when the number of variables is of the same order as the number of observations. This situation, known as Freedman's paradox, was illustrated by Freedman (1983).

methods. Miller (1998) notes that estimated regression coefficients are often far away from zero in such cases; this is a type of model selection bias. The resolution of this paradox is in the a priori modeling constraint that the number of candidate models small, achieving a large reduction in the number of parameters to be estimated, and basing the selection on more than one model.

It is not uncommon to see biologists collect data on 50 or more predictor variables in the blind hope that some analysis method will “find the variables that are significant” and sort out the rest. This is a shotgun strategy (Olden and Jackson 2000). This shotgun strategy will likely find spurious correlations (Anderson et al. 2001b), and it is poor practice. The use of many of the traditional multivariate analysis methods (e.g., principal components, stepwise discriminant function analysis, canonical discriminant analysis, methods, and factor analysis) found in the biological literature is a strategy that mostly spurious results will be found using this unthinking strategy (see Flack and Chang 1987 and Miller 1990), and we encourage biologists to give very serious consideration to a well-founded set of candidate models and predictor variables (as a reduced set of possible predictors) before minimizing the inclusion of spurious variables and relationships. Biologists are not alone in collecting a small amount of data on a very large number of variables. A. J. Miller (personal communication) indicates that the number of data sets in other fields with as many as 1,500 variables is small. In such cases is less than 40 (a purely statistical search for meaningful relationships in such data is doomed to failure).

After a carefully defined set of candidate models has been identified, the biologist is left with the evidence contained in the data; the task of the biologist is to extract this evidence from analyzing the data. Questions such as “What model is supported by the data?” can be answered objectively. This approach allows a clear place for experience (i.e., prior knowledge) to be used. The results of past studies, the biological literature, and current knowledge inform the modeling process formally. Then, one turns to the question of what is “important” within a sense of parsimony. In some cases, a small set of the number and nature of the predictor variables to be used in the model will suffice in defining the candidate models. This procedure starts with an initial set of, say, 15–40 predictor variables and a consistent set of model forms. A smaller set to use in the set of candidate models. Using Akaike’s Information Criterion methods one can only hope to select the best model from the set of candidate models; if models are not in the set of candidates, they cannot be selected. (i.e., data analysis) algorithms.

We lament the practice of generating models (i.e., “models”) in the total absence of real data, and yet “inferences” are made about the structure, and functioning of the real world based on such models. We do not object to the often challenging and stimulating

of 20 candidate models for predator-prey relationships, their “credibility” and “parsimony.” However, as is often the case with empirical data available on a variety of taxa to pursue the same question in a systematic manner (also see Turchin and Batzli (2001), who suggest a system of 2–3 differential equations, for vegetation–herbivore interactions). Such exercises help us sort out ideas that are plausible and their logical consequences are explored. Modeling exercises help us test our logical and quantitative abilities. Modeling exercises help us think into how the world *might* function, and hence modeling exercises help us explore alternative hypotheses to be explored with real data. Our caution is that the confusing of presumed insights from such models with the real world (see Peters 1991, Weiner 1995). An inference about some aspect of the real world is justified only after the model has been shown to adequately fit relevant empirical data (this will be a difficult task when the model in its totality has been fit to and tested against the data). Gause (1934) had similar beliefs when he stated, “Mathematical models independent of experiments are of but small importance.”

The underlying philosophy of analysis is important here. A conservative approach to the overall issue of *strategy* in the analysis of the biological sciences with an emphasis on a priori considerations should be considered. *Careful, a priori consideration of alternative models often require a major change in emphasis among many people* is an unfamiliar concept to both biologists and statisticians, but it is a tendency to use either a traditional model or a model with a lot of software, making its use easy (Lunneborg 1994). This approach is in contrast to strategies advocated by others who view modeling as a highly iterative and interactive exercise. Such a strategy involves deliberate data dredging and should be reserved for early stages of initial investigation. Such an exploratory avenue is not the focus of this book.

Here, we advocate the deliberate exercise of carefully considering, say, 4–20 alternative models as potential approximations to the level of information in the data available and the scientific question addressed (Lytle 2002 provides an advanced example). Such a strategy might have as many as 70–100 or more models that are considered. The number of candidate models is often larger than the number of models considered. We find that people tend to include many models that are not supported by the data than the data could reasonably support (e.g., models with too many parameters). There needs to be some well-supported guidance to help analysts better define the models to be considered. Models should be developed without first deeply examining the data, called *candidate models*.” The science of the issue enters the arena of the a priori set of candidate models.

for data analysis are the “true model” that generates the bi-serve (see, for example, Bancroft and Han 1977). We believe reality) in the biological sciences has essentially infinite d full reality cannot be revealed with only finite samples of d those data. It is generally a mistake to believe that there is a in the biological sciences and that during data analysis th covered and its parameters estimated. Instead, biological s with many small effects, interactions, individual heteroge and environmental covariates (most being unknown to us to identify a model that provides a good *approximation*. The words “true model” represent an oxymoron, except i Carlo studies, whereby a model is used to generate “data” u numbers (we will use the term “generating model” for s studies). The concept of a “true model” in biology seem may even be a source of confusion about the nature of ap (e.g., see material on BIC and related criteria in Chapter

A model is a simplification or approximation of reality reflect all of reality. Taub (1993) suggests that unproductive true models can be avoided by simply recognizing that a r definition. Box (1976) noted that “all models are wrong, b While a model can never be “truth,” a model might be rank to useful, to somewhat useful to, finally, essentially usele methods try to rank models in the candidate set relative to any of the models is actually “good” depends primarily o data and the science and a priori thinking that went into the (reality) is elusive (see deLeeuw 1988). Proper modeling o what inferences the data support, not what full reality m 1982:14–15, Lindley 1986). Models, used cautiously, tell supported by the (finite) data available.” Increased sampl allows us to chase full reality, but never quite catch it.

The concept of truth and the false concept of a true surprisingly important. Often, in the literature, one see model or simply *the* model as if to be vague as to the exact Bayesians seem to say little about the subject, even as to of the prior probabilities on models. Consider the simple size (n) at time t ,

$$n_{t+1} = n_t \cdot s_t,$$

where s is the survival probability during the interval from a correct model in the sense that it is algebraically and d rect; however, it is not an exact representation or model is not explanatory; it is definitional (it is a tautology, be

be improved in average population survival probability, but not able from a beta distribution; still, this is far from a model of the truth, even in this very simple setting. Individual variation can be caused by biotic and abiotic variables in the environment. An exact model of full reality would have, at the very least, one model individual as a nonlinear function of a large number of environmental variables and their interaction terms. Even in this simple case, it is not possible to expect any mathematical model to represent full reality. The true models in the biological sciences. We will take a set of models g_i , without pretending that one represents full reality. “true.”

In using some model selection methods it is assumed that the set of models contains the “true model” that generated the data. This is a very strong assumption, unless we use a data set generated by Monte Carlo simulation as a tutorial example (e.g., Section 3.4), and then we will know the condition clear. In the analysis of real data, it seems unlikely that the “true model” is included in the set of candidate models. If the true model exists at all. Even if a “true model” did exist and we used some method, it would not be good as a fitted model. The model (i.e., understanding or prediction) about some biological system with numerous parameters would have to be estimated from the data. The precision of these estimated parameters would be quite low.

Often the investigator wants to simplify some representation of the system in order to achieve an understanding of the dominant aspects of the system under study. If we were given a nonlinear formula with 200 parameters, we could make correct predictions, but it would be difficult to understand the main dynamics of the system without some further simplification. Thus, one should tolerate some inexactness (an inflated error) in order to get a simpler and more useful understanding of the phenomenon.

In particular, we believe that there are tapering effect sizes in many biological systems; that is, there are often several large, important effects, followed by many smaller effects, and, finally, followed by a myriad of very small effects. These effects may be sequentially unveiled as sample size increases. The dominant, effects might be relatively easy to identify and estimate with fairly poor analysis methods, while the second-order effects (e.g., treatment effect or an interaction term) might be more difficult to detect. Still smaller effects can be detected only with very large samples (see Kareiva 1994 and related papers), while the smallest effects are often never being detected, even with very large samples. Rare events that may be very important but quite difficult to study. Approximate models may be related to the amount of data and information available. Models that appropriately support only simple models with few parameters. Comprehensive data sets will support, if necessary, more

dimensioned model with relatively few parameters might be a good approximation of the full truth or reality. Biologists should not believe that a model exists that generates the data observed, although some models might be of relatively low dimension and could be well approximated by a fairly simple model. The issue of a range of tapering effects is common in epidemiology, where Michael Thun notes, "... you cannot get away from a big thing. What's very hard to do is to tell a little bit of the truth at all" (Taubes 1995). *Full reality will always remain elusive in the biological sciences.*

At a more advanced conceptual level, there is a concept of a model of the population (or process or system) under study. The model is used and the goal is to express this information in a more compact form using a "model." Conceptually, this is a change in the way we think about the data to using a different "alphabet." The data have only a finite amount of information. The *goal* of model selection is to achieve a good translation so that no information is lost; in fact, we cannot lose information. The data can be ideally partitioned into *information* and *noise*. The part of the data is not information. However, noise could contain information that we cannot decode. Conceptually, the role of a good model is to separate information from noise.

Our main emphasis in modeling empirical data is to understand the empirical structure, process, or system. Sometimes prediction is the goal, but here, however, one would hopefully have an understanding of the system as a basis for making trustworthy predictions. The process of developing a set of candidate models prior to intensive comparison and selecting one that is "best," and estimating the parameters of the best model with precision (using maximum likelihood or least squares methods) is a basis for valid inferences, and there are several methods to allow additional inferences and insights. In particular, we want to allow formal inference from more than one model, and to understand the advantages (Hoeting et al. 1999). Statistical science is not just a branch of mathematics, but rather it is concerned with the development of a theory of information using what is known or postulated about the nature of the matter. In our investigations into these issues we are concerned by how much uncertainty there is in selecting a good approximation of the data. The variability in terms of what model is selected or chosen from a set of independent data sets, for example, is often large.

1.2.6 An Ideal Approximating Model

We consider some properties of an ideal model for valid inference from the analysis of data. It is important that the best model is selected from a set of models that were defined prior to data analysis and based on the

thermore, precise, unbiased estimators of parameters would be appropriately simple, based on conceptual model. The best model would be accurate estimators of precision. The best model would have confidence interval coverage close to the nominal level (95%). Achieved confidence intervals of minimum width. Achieved coverage is a convenient index to whether parameter estimates and precision are adequate. Finally, one would like as good a model of the structure of the system as the information permits. The adjusted R^2 can be computed and σ^2 estimated as a measure of explained or residual variation, respectively. Ideally, the parameters of the model would have biological interpretations. If predictions from a model having the above issues in place might warrant some tentative predictions. There are many cases where two or more models are tied for “best,” and this should be fully recognized in further model selection, especially when they produce different predictions. There might be 4–10 models that have at least some support, and a careful scrutiny in reaching conclusions from the data, based on information, rather than a single model.

1.3 Model Fundamentals and Notation

This section provides a conceptualization of some important concepts as they are used in this book. Some of these classes are partially overlapping. A general notation is introduced that is intended to be useful to readers.

1.3.1 *Truth or Full Reality* f

While there are no models that exactly represent full reality, the concept of full truth can be denoted as f . The concept of f is abstract, representing the full truth which we want to make inferences, based on data and a model. We use the notation $f(x)$ to denote that integration is over x . We do not want to convey the notion that f is a function of x . The models arise from full reality and can be used to make formal inferences about full truth, if data collection has been carefully planned and the experimental design has been achieved.

1.3.2 *Approximating Models* $g_i(x|\theta)$

We use the notation $g_i(x|\theta)$ or often, if the context is clear, g_i for the i th approximating model. We use θ to represent general parameters of the model.

These models specify the structural form (including how covariates enter), but not the parameter values (the β_i); each β_i is an independent Bernoulli random variable. The first model is the null model. The second model represents the hypothesis that z is unimportant, while the third model is like the first, but z is hypothesized to be unimportant. The fifth model implies that z is constant and not a function of concentration. Often, one hypothesizes that the compound that model g_5 is not worth exploration. Of course, a log-log, or probit function could have been hypothesized, rather than the hypothesized relationships in this example, rather than the

1.3.3 The Kullback–Leibler Best Model $g_i(x|\theta_0)$

For given full reality (f), data (x), sample size (n), and model (g_i), g_i is a best model in the sense of Kullback–Leibler information (Chapter 2). That is, given the possible data, the form of g_i , and possible parameter values, K-L information can be computed for g_i in the set and the model best approximating full reality f is selected.

The parameters that produce this conceptually best single model $g_i(x|\theta)$, are denoted by θ_0 . Of course, this model is general, but can be estimated; such estimation involves computing the MLEs of the parameters in each model ($\hat{\theta}$) and then *estimating* K-L information for model selection and inference. The MLEs converge to the true parameters and the concept of bias is with respect to θ_0 , rather than “the true parameters” associated with full reality f .

1.3.4 Estimated Models $g_i(x|\hat{\theta})$

Estimated models have specific parameter values from θ_0 based on the given data and model. If another, replicate data set is generated and based on the same sample size, the parameter estimates will be different; the amount of difference expected is related to model g_i (e.g., standard errors and confidence intervals). It is important to note that the model form $g_i(x|\theta)$ from specific estimates of this model $g_i(x|\hat{\theta})$ and the process of parameter estimation, $g_i(x|\hat{\theta})$.

In the models of mortality as a function of concentration (above), there are associated likelihoods and log-likelihoods. The log-likelihoods can be used to obtain the MLEs $\hat{\beta}_0$ and $\hat{\beta}_1$ for model g_1 . The log-likelihood function is

$$\mathcal{L}(\beta_0, \beta_1 | \text{data}, \text{model}) = \prod_{i=1}^n (\mu_c(i))^{y_i} (1 - \mu_c(i))^{1-y_i}$$

1 + exp{-($\beta_0 + \beta_1 c$)}

Thus, the only parameters in the likelihood are β_0 and β_1 , one can obtain the MLEs. The value of the maximized log-likelihood and the estimated variance–covariance matrix can also be computed. If we have only the model form $g(x|\theta)$ we have an infinite number of models where all such models have the same form but different parameter values. Of all of these models there is a unique K-L best model. Consider how to find this model, given f .

1.3.5 *Generating Models*

Monte Carlo simulation is a very useful and general approach in theoretical and applied statistics (Manly 1991). These procedures are commonly specified as the basis for generating Monte Carlo data. Such a model is called reality, and thus we call it a *generating model*. It is “truth” or computerized truth. One should not confuse a generating model based on Monte Carlo data with full reality f .

1.3.6 *Global Model*

Ideally, the global model has in it all the factors or variables that are important. Other models are often special cases of this global model. It is not always a global model. If sample size is small, it may not fit the global model. Goodness-of-fit tests and estimates of overdispersion parameter for count data should be based (only) on the global model. The concept of overdispersion is relatively model-independent; however, it must be used to compute or model any overdispersion that is not accounted for by the data. Thus, the most highly parametrized model will serve best for assessing overall fit and estimating a parameter associated with overdispersion. In the models of mortality (above), model g_1 would be the global model.

The advantage of this approach is that if the global model is fitted adequately, then a selected model that is more parsimonious will fit the data (this is an empirical result, not a theorem). Parsimony should not lead to a model that does not fit the data (this principle is the basis for the selection methods we advocate here). Thus, goodness-of-fit tests and the estimation of overdispersion parameters should be based on the global model (this could also be computed for the selected model).

In summary, we will use the word “model” to mean distribution. Fully, the context will be clear. Certainly it is important to distinguish between f and g . The general structural form is denoted $g(x|\theta)$, specifying the numerical value of the parameter θ (e.g., $\theta = 0.5$).

by the space over which θ varies. Frequently, we will refer to MLEs (the most likely, given the data and the model) have cases we will mean the best model, $g(x|\theta_0)$, which is one K-L best relative to f).

1.3.7 Overview of Stochastic Models in the Biological Sciences

Models are useful in the biological sciences for understanding of systems, estimating parameters of interest and their associated covariance matrix, predicting outcomes and responses, and testing hypotheses. Such models might be used for “relational” or “descriptive” purposes or might be used for prediction. In the following material, we will discuss the main types of models used in the biological sciences. Although not meant to be exhaustive, it will allow the reader an impression of the types of models of empirical data that we will treat under an inferential framework.

Simple linear and multiple linear regression models (Sellers and Smith 1981, Brown 1993) have seen heavy use in the biological sciences in the past four decades. These models commonly employ a small number of parameters, and the statistical theory is fully developed (based on least squares or likelihood theory). Similarly, analysis of variance models have been widely used, and the theory underlying them is closely related to regression models and is fully developed (for models of general linear models). Theory and software for this work are readily available.

Nonlinear regression models (Gallant 1987, Seber and Lee 1989, et al. 1995) have also seen abundant use in the biological sciences (e.g., nonlinear regression is a common example). Here, the underlying theory is likelihood based, and some classes of nonlinear models require specialized software. In general, nonlinear estimation is a more advanced topic, somewhat less well understood by many practicing researchers.

Other types of models used in the biological sciences include generalized linear (McCullagh and Nelder 1989, Morgan 1992, 2000), additive (Hastie and Tibshirani 1990) models (these can be thought of as regression models). These modeling techniques have seen widespread use in the past decade. Multivariate modeling approaches such as multiple regression and regression, canonical correlation, factor analysis, principal component analysis, and discriminate function analysis have had a checkered history in the biological and social sciences, but still see substantial use (e.g., Brown and McCulloch 1990). Log-linear and logistic models (Brown 1993) have become widely used for count data. Time series models (Brown 1993, 1987, 1991) are used in many biological disciplines. Variational models of an organism’s growth (Brisbin et al. 1987, Gochfeld 1987) have

Compartmental models are a type of state transition model with a continuous response and are usually based on systems of ordinary differential equations (Brown and Rothery 1993, Matisoff and Rothery 1971); these have found use in a wide variety of fields including mathematical models of disease transmission. More advanced methods of wide application include the class of models called “random walk models” (deLeeuw 1998).

Models to predict population viability (Boyce 1992), often based on a type of Leslie matrix, are much used in conservation biology. Many alternative model forms given serious evaluation. A common criticism is that these models are rarely based on empirical data; the focus is on parameter values are often merely only “very rough guesses” by the lack of empirical data (White 2000).

Biologists in several disciplines employ differential equations in their research (see Pascual and Kareiva 1996 for a reanalysis of data and Roughgarden 1979 for examples in population genetics and primary ecology). Many important applications involve exploitation of computer software (Myers et al. 1995). Computer software exists to allow models to be estimated using least squares or maximum likelihood methods (e.g., Splus). These are powerful tools in the analysis of empirical data, but they beg the issue of “what model to use.”

Open and closed capture–recapture (Lebreton et al. 1992) and mark–recapture (Brownie et al. 1985) models represent a class of models based on multinomial distributions (see issues 5 and 6 of volume 10 of *Journal of Applied Statistics*, 1995). Distance sampling theory (Borchers 2001) relies on models of the detection function and often involves metric models. Parameters in these models are nearly always estimated by maximum likelihood.

Spatial models (Cressie 1991 and Renshaw 1991) are widely used in the biological sciences, allowing the biologist to take advantage of spatial data sets (e.g., geographic information systems). Stein and Cressie (1991) have shown how Kriging (perhaps the most widely used spatial model) can be expressed as a least squares problem, and the development of Monte Carlo methods such as the Gibbs sampler (Robert and Casella 2004). Chen et al. (2000) allow other forms of spatial models, including least squares or maximum likelihood (Augustin et al. 1996). Further, the use of methods widely used on biological data has been common (e.g., Brown and Brooks (1990)). Geographic information systems potentially allow large numbers of covariates for biological models, so that models can be particularly important.

Spatiotemporal models are potentially invaluable to biologists. Most researchers model changes over space or time, and

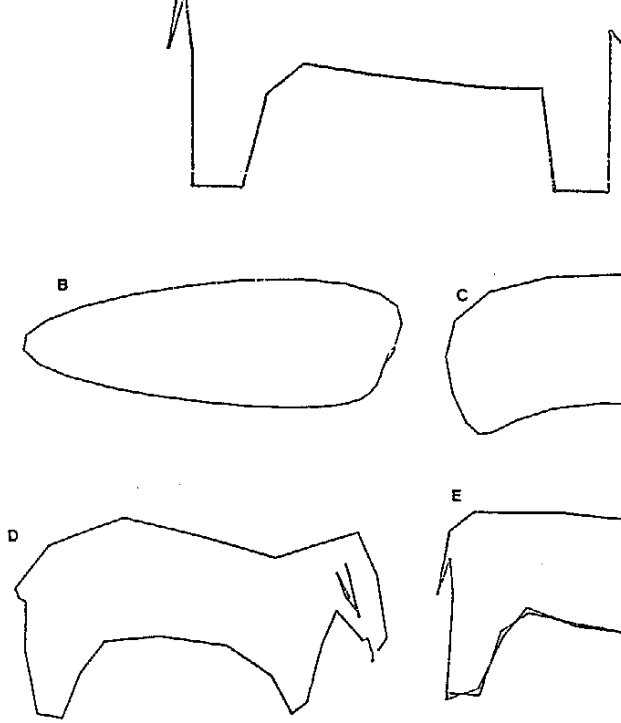


FIGURE 1.2. “How many parameters does it take to fit an elephant?” by Wel (1975). He started with an idealized drawing (A) defined by least squares Fourier sine series fits of the form $x(t) = \alpha_0 + \sum \alpha_i \beta_0 + \sum \beta_i \sin(it\pi/36)$ for $i = 1, \dots, N$. He examined fits for $K = 1, 2, \dots, 30$ in B–E) and stopped with the fit of a 30 term model. He concluded “may not satisfy the third-grade art teacher, but would carry most of the preliminary design.”

in both science and technology, and it is embodied in the principle of parsimony. Albert Einstein is supposed to have said, “Everything should be made as simple as possible, but no simpler.”

Success in the analysis of real data and the resulting inference depend importantly on the choice of a best approximating model. In the biological sciences should be based on a parsimonious model. A good, accurate approximation to the structural information in the data should not be viewed as searching for the “true model.” Model selection are essentially concerned with the “art of approximation” (Box, 1974).

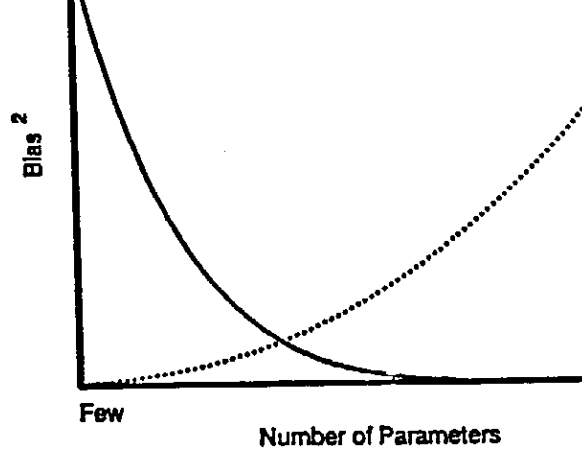


FIGURE 1.3. The *principle of parsimony*: the conceptual tradeoff between bias (solid line) and variance (dotted line) vs. the number of estimable parameters in the model. Model selection methods implicitly employ some notion of this tradeoff. The optimal model need not occur exactly where the two curves intersect. Full model fit is not attainable with finite samples and usually lies well to the right of the intersection. The region to the left of the intersection is the *tradeoff region*. Bias decreases and variance increases as the number of parameters in a model increases.

1.4.2 The Principle of Parsimony

If the fit is improved by a model with more parameters, when should we stop? Box and Jenkins (1970:17) suggested that the *principle of parsimony* should lead to a model with "... the smallest possible number of parameters for adequate representation of the data." Statisticians view model selection as a bias versus variance tradeoff. In general, bias decreases and variance increases as the dimension of the model (K) increases. When, then, we may use the number of parameters in a model as a measure of the amount of structure inferred from the data. The fit of any model is improved by increasing the number of parameters (e.g., the elephant-fit model). However, even, a tradeoff with the increasing variance must be considered. The *principle of parsimony* for model selection is a tradeoff between bias and variance. All model selection methods are based on the *principle of parsimony* (Breiman 1992, Zhang 1994).

In understanding the utility of an approximate model for inference, it is convenient to consider two undesirable possibilities: underfitting and overfitting. Here, we must avoid judging a selected model as the supposed "true model," as occurs when data are simulated from a very simple, model using Monte Carlo methods. In this case,

to parameters). This interpretation is often of little value. If one ignores the principle of parsimony and its implications and holds the conception that such a simple true model exists in biology, and yet believe that truth is essentially infinite-dimensional, then one is defined in terms of the number of parameters in the fitted model. In this use of the terms “underfitted” and “overfitted” that such a model of a low-dimensional “true model” as a “standard.”

Instead, we reserve the terms underfitted and overfitted only in relation to a “best approximating model” (Section 1.2.6). Here, a model that would ignore some important replicable (i.e., conceptual or structural) other samples) structure in the data and thus fail to identify features actually supported by the data. In this case, bias in the parameter estimates is often substantial, and the sampling variance is underestimated, resulting in poor confidence interval coverage. Underfitted models miss important treatment effects in experimental settings. Models as judged against a best approximating model, are often overfitted. Parameter estimators, but have estimated (and actual) sample variances are needlessly large (the precision of the estimators is poor). The model could have been accomplished with a more parsimonious model. The treatment effects tend to be identified, and spurious variances are avoided with overfitted models. Shibata (1989) argues that underfitting is a more serious issue in data analysis and inference than overfitting. Model assessment breaks down in many exploratory studies where the number of parameters be only 35–80 and there are 20–80 explanatory variables. The model may expect substantial overfitting and many effects that are not real (Freedman 1983, Anderson et al. 2001b).

The concept of parsimony and a bias versus variance trade-off are important. Thus we will provide some additional insights (also see Forster and Sober 1994, and Jaffe and Spirer 1987). The purpose of estimation and analysis is to make inferences from the sample to the population. The inferences relate to the *information* about the system under study as inferred from the models considered. The parameter estimated in each model. A paramount consideration is the quality of good precision, of any inference reached. When we imagine many samples, there will be some recognizable features common to all samples. Such features are the sort of inference about which we have strong inferences (from our single sample). Other features are not. For example, say, 60% of the samples yet still reflect something real about the process under study, and we would hope to make weaker inferences about these. Yet additional features appear in only a few samples. These be best included in the error term (σ^2) in modeling. If one makes an inference about these features quite unique to just the single sample, as if they applied to all (or most all) samples (hence to the

replicable over samples is underfitting. The data are not replicable over such samples (see Chatfield 1996, Collopy et al. 1996) that structure with a model form and parameter estimates is “sampling variation” that must also be estimated (inferred).

True replication is very advantageous, but this tends to be in the case of strict experiments where replication and random foundation. Such experimental replication allows a valid estimate of variation (σ^2). An understanding of these issues makes one realize that when observational studies seem possible and practical, they seem less feasible.

A best approximating model is achieved by properly balancing the trade-off of underfitting and overfitting. Stone and Brooks (1992) state, “... straddling pitfalls of underfitting and overfitting.” This balance is achieved when bias and variance are controlled to achieve a good coverage at approximately the nominal level and where the error is at a minimum. Proper model selection rejects a model that is far from the truth. Attempts to identify a model in which the error of approximation to random fluctuations are well balanced (Shibata 1983, 1980) and model selection methods are “parsimonious” (e.g., BIC, Schwarz 1978) are realistic situations, to select models that are too simple (i.e., underfitting) where bias is large, precision is overestimated, and achieved coverage is well below the nominal level. Such instances are called underfitting. One has only a highly precise, quite biased result.

Sakamoto et al. (1986) simulated data to illustrate the consequences of underfitting and the errors of underfitting and overfitting models (Figure 1). The data (each with $n = 21$) were generated from the simple model

$$y = e^{(x-0.3)^2} - 1 + \epsilon,$$

where x varied from 0 to 1 in equally spaced steps of 0.05, and ϵ was a random variable with mean 0 and variance 1. Thus, in this case, they considered the generating model with parameters: 0.3, -1 , and 0.01. They considered the set of candidate models (the approximating models) to be simple polynomials of order 0 to 5, as shown in the table below.

Order	K	Approximating Model
0	2	$E(y) = \beta_0$
1	3	$E(y) = \beta_0 + \beta_1(x)$
2	4	$E(y) = \beta_0 + \beta_1(x) + \beta_2(x^2)$
3	5	$E(y) = \beta_0 + \beta_1(x) + \beta_2(x^2) + \beta_3(x^3)$
4	6	$E(y) = \beta_0 + \beta_1(x) + \beta_2(x^2) + \beta_3(x^3) + \beta_4(x^4)$
5	7	$E(y) = \beta_0 + \beta_1(x) + \beta_2(x^2) + \beta_3(x^3) + \beta_4(x^4) + \beta_5(x^5)$

Thus, each of these 6 models was fit to each of the 10 simulated data sets.

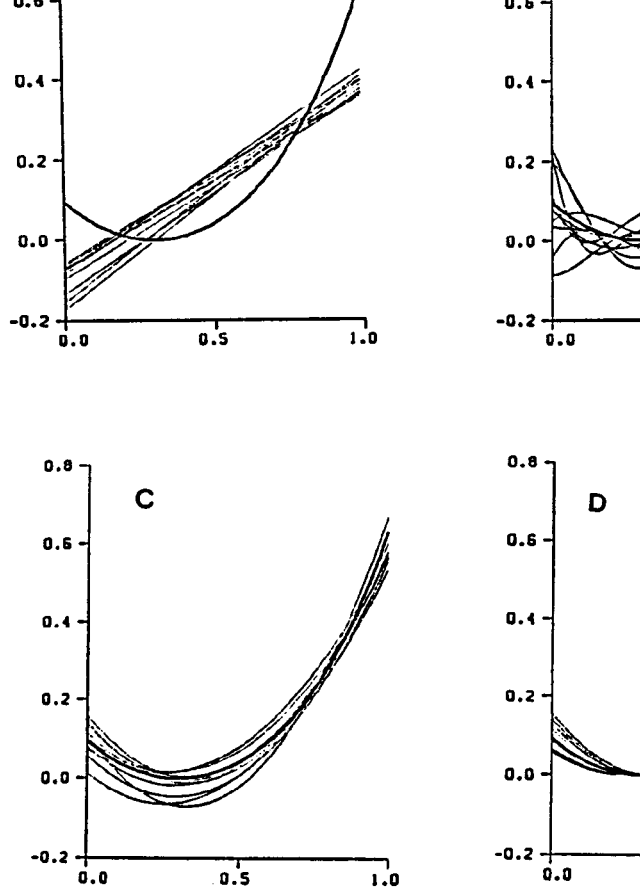


FIGURE 1.4. Ten Monte Carlo repetitions of data sets ($n = 21$) generated by $y = e^{(x-0.3)^2} - 1 + \epsilon$; $0 \leq x \leq 1$, $\epsilon \sim N(0, .01)$ (from Sakamoto et al.). A 5th-order polynomial (A) clearly misidentifies the basic nonlinear structure and is unsatisfactory. A 2nd-order polynomial (B) has too many parameters and large variance, and will have poor predictive qualities because it is overfitted. Neither A nor B is properly parsimonious, nor do they represent a best fit. A 2nd-order polynomial seems quite good as an approximating model, but since the function is nonnegative and has its minimum at $x = 0.3$, a model that enforces these conditions is improved further (D). In none of the plots one lacks the benefit of simple plots and 10 independent data sets, A–D. See Section 3.7 for a full analysis of these data.

1.4B). Here bias is obvious, the nonlinear structure of the model is poorly approximated, and confidence interval coverage for the model will be quite poor. Of course, there is *some* model bias in 5 models because they are only simple polynomial approximations. This is illustrated in Figure 1.4B, where a 5th-order polynomial is used as an approximating model. Here, there is little evidence of bias (quantity), precision is obviously poor, and it is difficult to discern the structure of the model. Prediction will be quite imprecise, and it has features that do not occur in the generating model. One extrapolates beyond the range of the data (always a bad idea). Underfitting and overfitting are undesirable in judging approximations for data analysis.

If a second-order polynomial ($K = 4$) is used as the approximating model, the fits seem quite reasonable (Figure 1.4C), and one may have more confidence from this model. Finally, *if* it were known a priori of the situation that the function was nonnegative and had a minimum at $x = 0.3$, then an improved quadratic approximating model could be obtained, information very effectively (Figure 1.4D). The form of the model is

$$E(y) = \beta_0(x + \beta_1)^2$$

with $K = 3$ (i.e., β_0 , β_1 , and σ^2), whereas the second-order model has 4 parameters. This example illustrates that valid statistical inference is partially dependent on the analysis process; the science of statistics can play an important role through modeling. This particular example provides a visual image of underfitting and overfitting in a simple case. The choice of an approximating model and various approximating models can be evaluated in high dimensions. Parsimony issues with real data in the biological sciences always defy such a simple graphical approach because true functions rarely have 10 independent data sets on exactly the same plane. In high dimensions are problematic to produce and interpret. The generating model contained no tapering effects. However, the approximating models do have tapering effects. Therefore, objective and subjective models are needed that do not rely on simple graphics and can cope with the complexities and high dimensionality.

1.4.3 Model Selection Methods

Model selection has most often been viewed, and hence taught, in the context of null hypothesis testing. Sequential testing has most often been viewed in the context of stepup (forward) or stepdown (backward) methods. Stepwise selection for variables to be added or deleted at each step. These testing schemes remain popular in many computer software packages in the biological sciences. Operating characteristics. Testing schemes are based on

tests are to be made (see Westman and Young 1995), independent. Tests between models that are not nested model is nested if it is a special case of another model; first-degree polynomial is nested within a fourth-degree polynomial. Hypothesis testing is a very poor basis for model selection (Scolove 1994b). McQuarrie and Tsai (1998) do not even treat for a short appendix on stepwise regression—the final three

Cross-validation has been suggested and well studied for model selection (Mosteller and Tukey 1968, Stone 1974, 1977; see also the data are divided into two partitions. The first partition is used for fitting; and the second is used for model validation (so that each partition has only one observation). Then a new partition is chosen, the whole process is repeated hundreds or thousands of times, and the partition then chosen, such as minimum squared prediction error, is used for model selection. There are several variations on this theme, and it has been applied to ecology (Craven and Wahba 1979, Burman 1989, Shao 1993, Hjorth 1994). These methods are quite computer intensive, but are practical if more than about 15–20 models must be evaluated. If the sample size is large, still, cross-validation offers an interesting alternative to model selection.

Some analysts favor using a very general model in all situations (the fitted model). We believe that this is generally poor practice. Others have a “favorite” model that they believe is good for nearly all situations. For example, some researchers always use a simple model (Buckland et al. 1993) with 2 parameters ($K = 2$) for model selection to the detection function in line transect sampling. This may be what is reasonable for situations where a simple model suffices (e.g., to 3), but will be poor practice in more challenging models where $10 \leq K \leq 30$ or more is required. These *ad hoc* rules ignore the principle of parsimony and data-based model selection, in which the model is chosen for model to be used for inference.

If goodness-of-fit tests can be computed for all alternative models, and some are not nested within others, then one could use the model with the fewest parameters that “fits” (i.e., $P > 0.05$ or 0.10). However, good fits can often be achieved by using models with more parameters (e.g., the elephant-fitting problem), and this can make the choice of α very critical. A large α -level leads to overfitted models and associated problems. In addition, other problems may be encountered, such as underdispersion and low power if one must pool small expected values so that the test statistic is chi-square distributed. Perhaps, more research is needed. There is no theory to suggest that this approach will lead to satisfactory model selection with good inferential properties (i.e., an adequate bias vs. variance trade-off, or achieved confidence interval coverage and width).

where R^2 is the usual coefficient of multiple determination (1981:91–92). Under this method, one selects the model in which the adjusted R^2 statistic is largest. McQuarrie and Tsai (1998) found this method to be poor (also see Rencher and Pun (1980). While adjusted R^2 is a descriptive statistic, it is not useful in model selection. Mallows's C_p (Mallows 1973, 1995) is also used in LS regression with normal errors and a constant variance and in this special case provides a ranking of candidate models that is the same as the rankings under AIC (though C_p vs. AIC, will differ, see Atilgan 1996). The selection of the model with the adjusted R^2 statistic and Mallows's C_p are related for linear regression (see Seber 1977:362–369). Hurvich and Tsai (1989) and Tsai (1998) provide some comparisons of AIC_c vs. several other criteria for regression problems.

Bayesian researchers have taken somewhat different assumptions, and have proposed several alternative methods of model selection. Methods such as CAIC, BIC (SIC), WIC, and HQ are mentioned in the literature as well as full Bayesian model selection (see especially Kass and Raftery 1995). These other Bayesian approaches to model selection and model comparison are in the current state of the art in statistics but may seem very difficult to understand and implement and are very computer intensive (e.g., Laud and Carlin and Chib 1995). Draper (1995) provides a review of some advanced methods (also see Potscher 1991). Spiegelhalter (1994) has developed a deviance information criterion (DIC) from a Bayesian perspective that is analogous to AIC. This seems to represent a blending of frequentist and Bayesian thinking, resulting in an AIC-like criterion.

The general approach that we advocate here is one derived from Akaike (1974, 1977, 1978a and b, and 1981a and b), based on information theory. It is discussed at length in this book. Akaike's information criterion has led to a number of alternative methods having desirable properties for the selection of best approximating models in practice (e.g., AIC, BIC, and TIC—Chapters 2 and 7). Our general advocacy concerning model selection criteria is somewhat stronger than that of Akaike (1986) but similar in that they also recommend objective criteria. We recommend some well-defined criterion with a strong, fundamental basis.

1.5 Data Dredging, Overanalysis of Data, and Spurious Effects

The process of analyzing data with few or no a priori questions, but rather, tentatively and iteratively searching the data for patterns and “significant” results, is called by the derogatory term “data dredging.” Other terms

large data sets. Often the problem arises when data on a large number of variables are available, but the number of observations is small. In such cases, the data are often analyzed without any prior knowledge of the problem, and the results are often misleading. No specific objectives or alternatives were in place; thus the data are submitted for analysis in the hope that a plethora of null hypothesis test results will provide information on what is significant. A model is fit, and variables not in that model are added, creating a new model, letting the data and intermediate results be investigated. Patterns seen in the data are used to create new models and variables to be investigated. Patterns seen in the analysis are "chased" as new variables, cross products, and other transformations are added to the model and alternative transformations are added. The process is iterative, and new models are clearly based on the intermediate results from the previous analyses. The final model is the result of effective dredging, and everything remaining is "significant." Under this view, Hoaglin et al. (1989:169) comment that "Model fitting is an iterative process, and it is unlikely that one can obtain the final model on the first pass through the data." It is also likely that such a final model is probably overfitted and unstable, and that its predictive performance (i.e., on new data) often will be considerably less than the actual predictive performance (i.e., on new data) often will be expected from the statistics provided by the terminal model. (Hoaglin et al. 1989:169, Wang 1993). The inferential properties of a model based on dredging are very different. For example, (traditionally) no statement of precision can be made from the model following data dredging (Hoaglin et al. 1989:169).

1.5.1 Overanalysis of Data

If data dredging is done, the resulting model is very much tailored to the data in a post hoc fashion, and the estimates of the parameters of the model are likely to be overestimated. Such tailoring overdescribes the data, and the validity of inferences made about the information in the data is likely to be compromised. Many naive applications of classical multiple regression are merely "fishing trips" hoping to find "significant" linear relationships among the many variables subjected to analysis (Rexstad et al. 1998, Reid 2000).

Computer routines (e.g., SAS INSIGHT) and associated software have made data dredging both easy and "effective." Some statistical literature has referred to this as the *iterative process of model building* (e.g., Hendry 1981). One looks for patterns in the residuals, employs various transformations, and adds variables in their decreasing order of "importance," and so on. Stepwise regression and discriminant functions, for example, are used to search for "significant" variables; such methods are especially problematic when many variables (Freedman's paradox) are available for analysis. If data are available on over 100 variables, and the sample size is small, the results of such analyses are likely to be misleading.

White (2000, 1097) notes, "It is widely acknowledged by researchers that data snooping [dredging] is a dangerous practice, but in fact it is endemic." Examples of data dredging include the selection of crossplots or a correlation matrix of the explanatory variables, the selection of the response variable, the selection of linear or nonlinear relationships and interactions *in the sample*, and the selection of the model. These data-dependent activities can lead the investigator to consider additional models. These models are avoided, because they probably lead to overfitted models with poor parameter estimates and inclusion of unimportant variables as regressors (Anderson et al. 2001b). The sample may be well fit, but the inference from the sample to the population is invalid. This type of analysis, exploratory data analysis, has a place in the earliest stages of research, to explore a possible relationship but should probably remain unpublished. Cases are not the subject of this book, and we can only hope that the results of such procedures be treated as possible hypotheses (Longford and Nelder 1999). New data should be collected to test hypotheses effectively and then submitted for a comprehensive strategy of analysis such as we advocate here.

Two types of data dredging might be distinguished. The first type is the one described above; a highly interactive, data dependent, iterative process. The second is also common and also leads to likely overfitting and spurious effects that are actually spurious. In this type, the investigator has a priori information; thus "all possible models" are considered (e.g., SAS PROC REG allows this as an option). Note that the "all possible models" approach usually does not include interaction terms, but various transformations such as $(x_1)^2$ or $1/x_3$ or $\log(x_2)$. In many problems, the number of candidate models in this approach is large (e.g., 20 variables > a million models, 30 variables > a billion models). At least this second type is not explicitly data dependent, but it is data dependent and leads to the same "sins." Also, it is usually a one-time exercise rather than taking the results of one set of analyses and incorporating them into the consideration of new models. Still, in some applications, software often can systematically search all such models and thus the strategy of trying all possible models (or at least a large number of models) continues, unfortunately, to be popular. While these situations could be substantially improved if the researcher focuses on the science of the situation before proceeding with such an approach.

Standard inferential tests and estimates of precision (e.g., standard errors of the sampling covariance matrix, given a model) are based on model results from the first type of data dredging. Results are misleading, and there is no valid basis to claim "significantly different" or "significant" when there is no way to estimate precision because of the

a levels of P values. However, if there were 1,000 models, the P value would be 0.00005, instead of the usual 0.05! Problems are often linked with the problems with hypothesis testing (Anderson et al. 2000). This approach is hardly satisfactory, ignored the issue and merely pretended that data dredging that the usual inferential methods somehow still apply. **Journal referees rarely seem to show concern for the validity of conclusions where substantial data dredging has occurred. This methodology based on data dredging has been allowed to continue in an unthinking manner.**

We certainly encourage people to understand their data and the scientific questions of interest. We advocate some examination prior to the formal analysis to detect obvious outliers and to determine a preliminary truncation point or the need for growth of distance sampling data). One might examine the residuals of the chosen global model to determine likely error distribution models (e.g., normal, lognormal, Poisson). However, if a problem is noticed while examining the residuals and this leads to a new variable, then we might suggest caution concerning data dredging. There can be a fine line between a largely a priori approach and data dredging.

Thus, this book will address primarily cases where there is prior knowledge concerning the issue at hand and where a set of good candidate models can be specified in advance of the analysis. Of course, there is some latitude where some (few) additional models are investigated as the analysis proceeds; however, results from these models should be kept clearly separate from the purely a priori results. That objective science is best served using a priori considerations, limited peeking at plots of the data, parameter estimates from the models, correlation matrices, or test statistics as the analysis proceeds, is condoned data dredging in confirmatory analyses, but allowed in more preliminary explorations. If some limited data dredging is done, careful analysis based on prior considerations, then we believe that types of results should be carefully explained in resulting publications (Anderson 1980). For this philosophy to succeed, there should be more consideration of alternative candidate models than has been in the past.

1.5.2 Some Trends

At the present time, nearly every analysis is done using computers. Biologists and researchers in other disciplines are increasingly using computer methods for more generalized analyses. Standard computer

should generally be avoided, except in (1) the early stages or (2) *after* a more confirmatory analysis has been done. In the first case, the investigator should fully admit to the process that the results and should treat them much more cautiously than the initial, a priori, approach. When done carefully, we can explore their data beyond the important a priori phase. We recommend a substantial, deliberate effort to get the models in place and try to obtain more confirmatory results. The post hoc issues that often arise after one has seen the results.

Data dredging activities form a continuum, ranging from the trivial to the grievous (mortal). There is often a fine line between the two. Our advice is to stay well toward the a priori end of the continuum to achieve a more confirmatory result.

One can always do post hoc analyses after the a priori analysis. They never go from post hoc to a priori. Why not keep one's eyes on the regard?

Grievous data dredging is endemic in the applied literature. It is taught or implied in statistics courses without the needed caution about the attendant inferential problems.

Running all possible models is a thoughtless approach and a waste of finding effects that are, in fact, spurious if only a single model is for inference. If prediction is the objective, model averaging estimates of precision should include model selection uncertainty. In this case, surely one can often rule out many models on

allow likelihood methods to be used where LS methods have been in the past. LS methods will see decreasing use, and likelihood methods will see increasing use as we proceed into the twenty-first century. Likelihood methods allow a much more general framework for addressing inference (e.g., a choice of link functions and error distributions for generalized linear logistic regression models). Another advantage in a likelihood framework is that confidence intervals with good properties can be set using the likelihood. Edwards (1976), Berger and Wolpert (1984), Azam (1997), and Morgan (2000) provide additional insights into the use of likelihood methods while Box (1978) provides the historical setting relating likelihood methods.

During the past twenty years, modern statistical science has moved away from traditional formal methodologies based on hypothesis testing (Clayton et al. 1986, Jones and Matloff 1994, Bozdogan 1994, Johnson 1995, Stewart-Oaten 1995, Nelder et al. 1999, Anderson et al. 2000). The historic emphasis on hypothesis testing

models as an integral component of data analysis. Computerized and bootstrapping approaches and other resampling schemes (Duval 1993, Efron and Tibshirani 1993), and such techniques are increasing use in the future.

The size or dimension (K) of some biological models has been increasing. This has tended to increase over the past two decades. Open and band recovery models commonly have 20–40 estimated parameters. A single data set and might have well over 200 parameters. Some of several data sets (see Burnham et al. 1987, Preface, for a review of these trends). Analysis methods for structural equations commonly have 30 parameters (Bollen and Long 1993). These are applications of model specification and selection is essential to answer the question: *inferences do the data support about the population?*

1.6 Model Selection Bias

The literature on model selection methods has increased rapidly in the past 15–25 years; much of this has been the result of Akaike's work in the mid-1970s. However, relatively little appears in the literature on the properties of the parameter estimators, given that a data selection procedure has been used (see Rencher and Purvis 1992, Tsai 1990, Miller 1990, Goutis and Casella 1995, Ye 1995). This is used to both select a parsimonious model and estimate the parameters and their precision (i.e., the conditional sampling covariance matrix of the selected model). These issues prompt a concern for both the accuracy and model selection uncertainty (Section 1.7).

Bias in estimates of model parameters often arises when model selection has been done. Miller (1990) provides a technical discussion of selection bias in the context of linear regression. He notes that in a stepwise analysis of meteorological data with large sample sizes and many date models. When selecting only about 5 variables from a pool of 10, he found t statistics as large as 6, suggesting that a parameter was very highly significant, and yet even the sign of the corresponding regression coefficient could be incorrect. Miller warns that P -value tests and selection software are totally without foundation, and large t statistics and coefficients are often caused by data-based model selection.

Consider a linear model where there is a response variable y and explanatory variables x_j , where $j = 1, \dots, 4$. Order is not important, for example, so for convenience let x_1 be, in fact, very important and somewhat important, while x_4 is barely important. Given a

data sets of the same size, from the same stochastic process, and x_2 would be included in the model in nearly all cases. Inference from a sample data set to the population would be selected that included predictors x_1 and x_2 (essentially all estimators of the regression coefficients associated with x_1 and x_2 would have good statistical properties with respect to bias). Standard theory tends to hold for the estimators $\hat{\beta}_1$ and $\hat{\beta}_2$.

Variable x_3 is somewhat marginal in its importance; a t -test that $|\beta_3|/\text{se}(\beta_3) \approx 1$, and thus its importance is somewhat marginal. It might be included in the model in only 15–30% of the 1,000 data sets where it is selected, it tends to have an estimated regression coefficient that is biased away from zero. Thus, an inference from a sample concerning the population tends to exaggerate the importance of x_3 . An inference from a data set in one of the remaining 70–80% would imply that x_3 was of no importance. Neither of these inferences is correct.

Variable x_4 is barely important at all (a tapering effect). A t -test that $|\beta_4|/\text{se}(\beta_4) \approx \frac{1}{4}$. This variable might be included in only 1–5% of the 1,000 data sets and, when it is selected, there will be a bias (away from 0) in the estimator of this regression parameter. In any particular sample where this variable is included in the model, the variable x_4 was much more important than is actually the case. The investigator has no way to know that $\hat{\beta}_4$, when selected, is in the upper 5–10% of its sampling distribution). Then, if one performs a t -test, where $t = \hat{\beta}_4 / \widehat{\text{se}}(\hat{\beta}_4)$, the likely decision will often be that x_4 is significant, and should be retained in the model. This is a false conclusion that comes from the fact that the numerator in the test is biased and the denominator is biased low. The analyst has no way to know this; the result is probably spurious.

When predictor variables x_3 and x_4 are included in the model, the estimator for a σ^2 is negatively biased and precision is exaggerated. These types of bias are called model selection bias and can occur in many types of models (Miller 1990, Ye 1998). Ye (1998) warns, "... the identification of a true structure bears little cost [i.e., including variables x_1 and x_2], but including through white noise has a heavy cost [i.e., including variables x_3 and x_4]. Of course, in the analysis of real data, the investigator typically does not know which (if any) variables are dominant versus those that are not. This is a problem of importance. Model selection bias is related to the problem of model selection, the notion of tapering effect sizes, and Freedman's (1983) paradox.

The problem of model selection bias is particularly serious in exploratory studies where it is available to guide the analysis. Many exploratory studies examine even thousands of models, based on a large number of explanatory variables. Very often the number of models exceeds the size of the population. If a model has been (somehow) selected, the analyst is usually

given the data, and estimates a conditional sampling variance for the selected model. Then estimated precision will be too small if the sampling component for model selection uncertainty is missing. Model selection uncertainty is the component of variance that reflects that model selection *estimates* which model is best, based on the single data set (in the fixed set of models considered) may be selected a different replicate data set arising from the same experiment.

Failure to allow for model selection uncertainty often leads to underestimating sampling variances and covariances that are too low, and confidence interval coverage will be below the nominal value. Current methods for coping with model selection uncertainty are at the forefront of statistical research; better methods might be expected in the coming years as the continued increases in computing power. Model selection uncertainty is problematic in making statistical inferences; if the goal is only to select a model, then perhaps selection uncertainty is a minor issue.

One must keep in mind that there is often considerable uncertainty in the selection of a particular model as the “best” approximating model. The data are conceptualized as random variables; their values are fixed only if another, independent sample were available. It is this “sampling” that results in uncertain statistical inference from the particular data set analyzed. While we would like to make inferences that would be true for other (hypothetical) data sets, our ability to do so is still limited. Even with procedures such as AIC, with its cross-validation procedure, using independent and identically distributed sample data. Various resampling methods will further improve our assessment of the reliability of our inferences, but it remains important to understand that model selection is accompanied by a substantial amount of uncertainty. This technique can effectively allow insights into model uncertainty. Similar issues are the subject of Chapter 5.

Perhaps we cannot totally overcome problems in estimating the reliability of a data-dependent selection method such as AIC (e.g., Ye 1998). This limitation certainly warrants exploration but model selection uncertainty is a quite difficult area of statistical inference. One must also consider the “cost” of *not* selecting a good parsimonious model for the analysis of a particular data set. That is, a model is just a tool, independent of the data and used to approximate the data. The model is used for inference. This procedure simply ignores both the uncertainty in model selection and the benefits of selection of a model that best approximates the data. This naive strategy certainly will incur substantial costs in terms of unreliable inferences because model selection uncertainty is ignored (Ye 1998). Alternatively, one might be tempted into an iterative, highly complex procedure of data analysis (unadulterated data dredging). Again, this approach incurs costs in terms of reliable inference using this approach. In

1.8 Summary

Truth in the biological sciences and medicine is extreme: we cannot hope to find exact truth or full reality from a finite amount of data. Thus, inference about truth must be based on an approximating model. Likelihood and least squares methods are a good inference theory if the model structure is “given.” For most scientific problems, the model is *not* “given.” Thus, the central question is “What is the best model to use.” This is the model selection problem.

The emphasis then shifts to the careful a priori definition of the set of models. This is where the science of the problem enters in. There should be a good rationale for including each parameter in the set, as well as a careful justification for why other models are excluded. The degree to which these steps can be implemented suggests a confirmatory analysis, rather than a more exploratory analysis. Critical evaluation of the scientific question and modeling alternatives, prior to looking at the data, has been underemphasized in many statistics classes in the past. This is an important issue, and one must be careful not to engage in data dredging, which weakens inferences that might be made. Information-theoretic methods provide a simple way to select a best approximating model from a set of models.

In general, the information-theoretic approach should be used instead of searching for a single best model as a basis for inference. Model uncertainty is included in estimates of precision, this is not always done in many cases. Instead, multimodel inference should be used for making valid inference. Here, models are ranked and the uncertainty in understanding of model uncertainty over the set. These methods are easy to understand and compute. Specific methodologies for this model selection are the subject of this book.

We cannot overstate the importance of the scientific method, the formulation of multiple working hypotheses, and the use of models to clearly and uniquely represent these hypotheses. The issues to be presented in the following chapters are “easy” to understand and interpret; however, they rest on both good science and good statistics to the issue. We try to emphasize a more confirmatory end to the issue. We try to emphasize a more confirmatory end to the issue, rather than exploratory work that has become common in the biological sciences, often led to so little (Anderson et al. 2000).

Data analysis is taken to mean the entire integrated process: a priori model specification, model selection, and estimation of the parameters and their precision. Scientific inference is based on this process.

The principle of parsimony is fundamental in the selection of a parsimonious model. The reward for proper model selection in terms of valid inference is substantial. There are substantial dangers in either underfitting or overfitting. If the analyst has selected a good approximating model, there are issues of model bias and model selection uncertainty. Perhaps these cannot be avoided, but their effects can be lessened. These issues will be addressed in the next section to follow.

Zhang (1994) notes that for the analyst who is less concerned with optimality it is more important to have available methods that are flexible enough to be used in a variety of practical situations. Theoretic methods fall in this broad class and, when used, provide reliable inference.

Information and Likelihood Theory Basis for Model Selection and

Full reality cannot be included in a model; thus we seek to approximate the effects or factors supported by the empirical data. The choice of an appropriate approximating model is critical to statistical inference for many types of empirical data. This chapter introduces concepts of information theory (see Guisasu 1977), which has been a discipline since the mid-1940s and covers a variety of theories and methods used in many of the sciences (see Cover and Thomas 1991 for a review). Figure 2.1 is produced from their book and shows the relationship of information theory to several other fields. The Kullback–Leibler “distance,” or “information,” between two probability distributions (Kullback and Leibler 1951) is introduced, discussed, and linked to entropy in this chapter. Akaike (1973) found a simple relationship between the Kullback–Leibler distance and Fisher’s maximized likelihood function (see deLeeuw 1992 for a brief review). This relationship is an effective, and very general methodology for selecting a model for the analysis of empirical data.

Akaike introduced his “*entropy maximization principle*” in the mid-1970s (Akaike 1973, 1974, 1977) as a theoretical basis for model selection. He followed this pivotal discovery with several papers beginning in the early 1980s (Akaike 1981a and b, 1983). This chapter introduces AIC and related criteria such as BIC and TIC. No mathematical derivations of these criteria are given; they are given in full detail in Chapter 7. We urge readers to read the derivation (given in Chapter 7), for without it, the simple

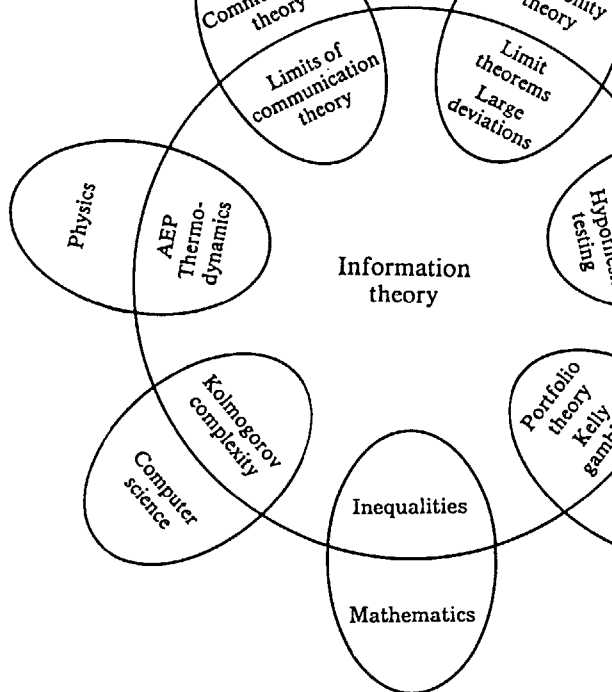


FIGURE 2.1. Information theory and its relationships to other disciplines (see Thomas 1991). Information theory began in the mid-1940s, at the time of this book, the most relevant components of information theory are information, entropy (from thermodynamics and communication theory), and Kullback–Leibler information.

underlying Kullback–Leibler information and the various other concepts cannot be fully appreciated.

2.1 Kullback–Leibler Information or Distance: Two Models

We begin without any issues of parameter estimation and confidence intervals. We start with expressions for the models f and g , assuming that they are probability distributions. In initial sections of this chapter we will let both f and g be probability distributions, since this will allow an understanding of the distance in a simple setting. However, we will soon switch to the more general case of

Kullback–Leibler Information

Kullback–Leibler information between models f and g for continuous functions as the (usually multi-dimensional)

$$I(f, g) = \int f(x) \log \left(\frac{f(x)}{g(x|\theta)} \right) dx,$$

where \log denotes the natural logarithm. The notation **“information lost when g is used to approximate f .”**

As a heuristic interpretation, **$I(f, g)$ is the distance from f to g .**

We will use both interpretations throughout this book, so far as it is useful. Of course, we seek an approximating model that loses as little information as possible; this is equivalent to minimizing $I(f, g)$, over all models g that are considered to be given (fixed), and only f varies over all models indexed by θ . Similarly, Cover and Thomas (1991) note that $I(f, g)$ is a measure of the inefficiency of assuming that the distribution is g when the true distribution is f .

Kullback–Leibler Information

The expression for the Kullback–Leibler information between two cases of discrete distributions such as the Poisson, binomial, or multinomial is

$$I(f, g) = \sum_{i=1}^k p_i \cdot \log \left(\frac{p_i}{\pi_i} \right).$$

Here, there are k possible outcomes of the underlying random process. The true probability of the i th outcome is given by p_i , while π_i constitute the approximating probability distribution (i.e., the model). In the discrete case, we have $0 < p_i < 1$, $\sum p_i = \sum \pi_i = 1$. Hence, here f and g correspond to the true and approximating distributions respectively.

As in the continuous case the notation **$I(f, g)$ denotes the information lost when g is used to approximate f or the distance from f to g .**

In the following material we will generally think of Kullback–Leibler information in the continuous case and use the notation f and g for simplicity.

Well over a century ago measures were derived for assessing the difference between two models or probability distributions. Most relevant is Boltzmann’s (1877) concept of generalized entropy (see Section 2.1) and thermodynamics (see Akaike 1985 for a brief review). Shannon (1948) employed entropy in his famous treatise on communication theory (see Cover and Thomas 2001 for an exciting review of information theory, its practical applications, and its role in evolution). Kullback and Leibler (1951) derived an information measure that is now known as the Kullback–Leibler information.



Ludwig Eduard Boltzmann, 1844–1906, one of the most famous made incredible contributions in theoretical physics. He received most of his work was done in Austria, but he spent some years in full professor of mathematical physics at the University of Graz, Austria. His mathematical expression for entropy was of fundamental importance in areas of science. The negative of Boltzmann’s entropy is a measure of information, over half a century later by Kullback and Leibler. J. Bronowski wrote of him “an irascible, extraordinary man, an early follower of Darwin, quantitative and everything that a human should be.” Several books chronicle the history of science, including Cohen and Thirring (1973) and Broda (1983). His papers appear in Hasenöhl (1909).

happened to be the negative of Boltzmann’s entropy, now called Kullback–Leibler (K-L) information or distance (but see Kullback, 1959; he preferred the term *discrimination information*). The negative of Kullback and Leibler’s work was to provide a rigorous definition of information in relation to Fisher’s “sufficient statistics.” The K-L distance is called the K-L discrepancy, divergence, information, and many other names. These terms are synonyms, but tend to use *distance* or *information* to follow.

The Kullback–Leibler distance can be conceptualized as the “distance” between two models, say f and g (Kullback 1959). This is a measure of “discrepancy”; it is not a simple distance. The measure from f to g is not the same as the measure from g to f , directed, or oriented, distance (Figure 2.2). The K-L distance

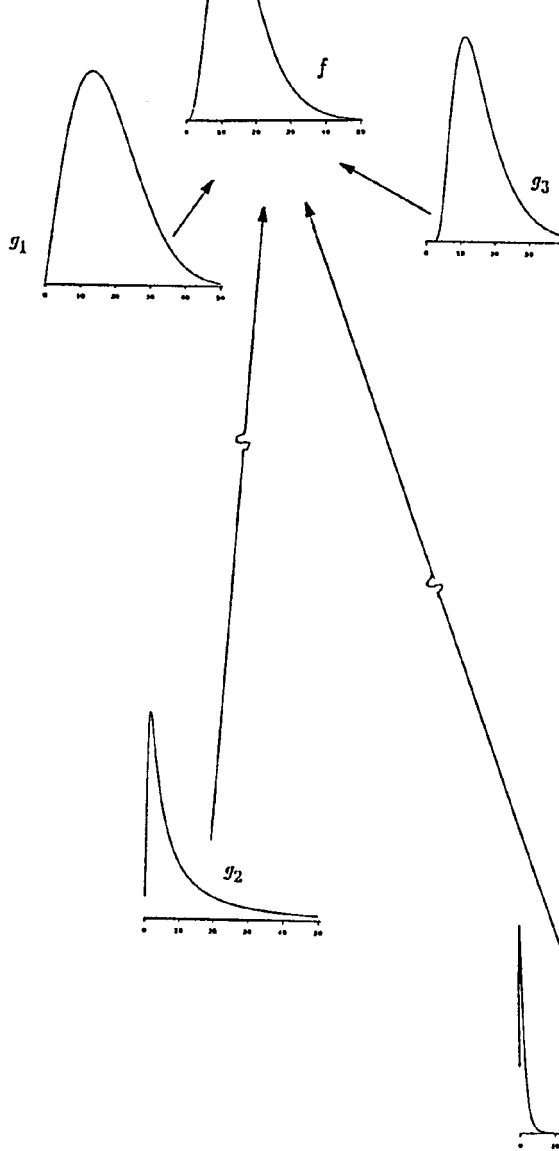


FIGURE 2.2. The Kullback–Leibler discrepancy $I(f, g_i)$ is a direct measure of the distance between the target distribution f and various candidate models g_i to f . Knowing the K-L distances with respect to f , we can determine which of the 4 approximating models is *closest* to model f . Here, the 4 approximating models are $g_1 = \text{Weibull}(2, 20)$, $g_2 = \text{lognormal}(16, 64)$, $g_3 = \text{Gaussian}(16, 64)$, and $g_4 = \text{F distribution}(4, 10)$. In each case, the K-L distance is known exactly (not estimated).

an extension of Shannon's concept of information (Foster and Soofi 1994) and is sometimes called a “relative entropy” between models is a *fundamental quantity* in science and (see Akaike 1983) and is the logical basis for model selection with likelihood inference.

At a heuristic level, “information” is defined as $- \int f(x) \log_e(p_i)$ for a continuous probability density function or $-\log_e(p_i)$ for a discrete probability mass function. Kullback–Leibler information is a type of “cross entropy” between two probability distributions. In either the continuous or discrete representation, it is an expected value (i.e., $\int f(x)(\cdot)dx$ for the continuous case and $\sum p_i(\cdot)$ for the discrete case) of the logarithm of the ratio of the two probability densities (f and g) or two discrete probabilities (p_i and π_i). In the continuous case, we can think of this as an average (with respect to f) of $\log_e(f(x)/g(x))$. In the discrete case it is an average (with respect to the p_i) of the $\log_e(p_i/\pi_i)$. The foundations of these expressions are both deep and broad (see Boltzmann 1877, Kullback and Leibler 1951, or Cover and Thomas 1991 for information theory).

The K-L distance ($I(f, g)$) is always positive, except when the two distributions f and g are identical (i.e., $I(f, g) = 0$ if and only if $f = g$ everywhere). More detail and extended notation will be in Section 2.2; here we will employ a simple notation and use it to measure the distance between f and g in the sample data (x) and the multivariate function $f(x)$.

2.1.1 Examples of Kullback–Leibler Distance

An example will illustrate the K-L distances ($I(f, g_i)$). We will consider a target distribution with 2 parameters ($\alpha = 4, \beta = 4$). Then consider four approximating models g_i , each with 2 parameters (see below): Weibull, lognormal, Gaussian, and the F distribution. Details on these simple probability distributions can be found in Johnson and Kotz (1970). The particular parameters for the four g_i are not material here, except to stress that they are not estimated. “Which of these parametrized distributions is closest to f ?” is answered by computing the K-L distance between each g_i and f (see Section 2.2). These are as follows:

	Approximating model	$I(f, g_i)$
g_1	Weibull distribution ($\alpha = 2, \beta = 20$)	0.046
g_2	lognormal distribution ($\theta = 2, \sigma^2 = 2$)	0.672
g_3	inverse Gaussian ($\alpha = 16, \beta = 64$)	0.060
g_4	F distribution ($\alpha = 4, \beta = 10$)	5.745

Here, the Weibull distribution is closest to (f loses the least information by approximating f), followed by the inverse Gaussian. The lognormal distribution is the farthest from f .

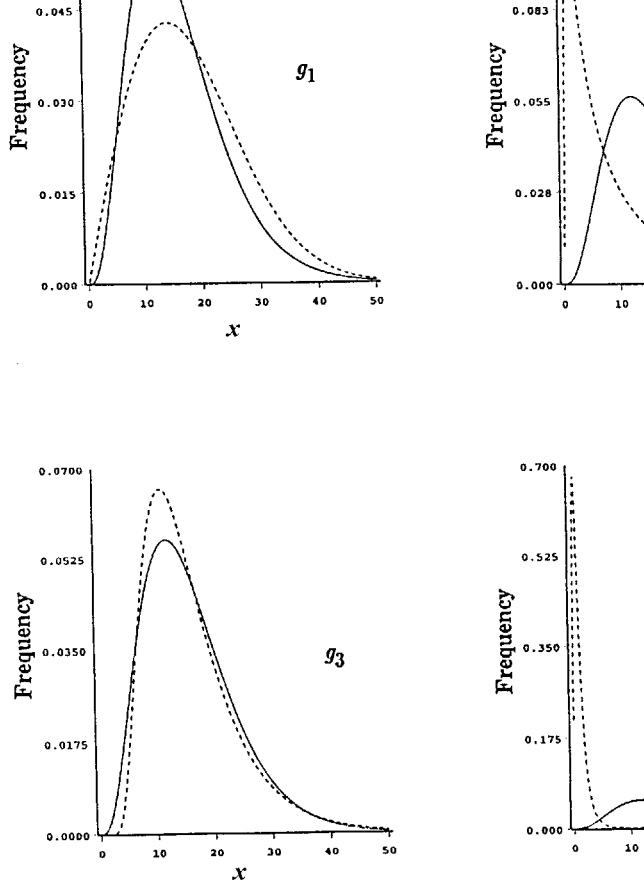


FIGURE 2.3. Plots of f ($=$ gamma (4, 4), solid line) against each models g_i (dashed lines) as a function of x . Here, $g_1 =$ Weibull (2, 2), $g_3 =$ inverse Gaussian (16, 64), and $g_4 =$ F distribution (4, 1) cases can plots such as these be used to judge closeness between same in all 4 graphs; it is merely scaled differently to allow the g_i same graph.

while the F distribution is relatively far from the gamma (Figure 2.3).

Further utility of the K-L distance can be illustrated by approximating models g_i might be closest to f when the allowed to vary (i.e., what parameter values make each g_i f ?). Following a computer search of the parameter space found that the *best* Weibull had parameters $\alpha = 2.120$ and

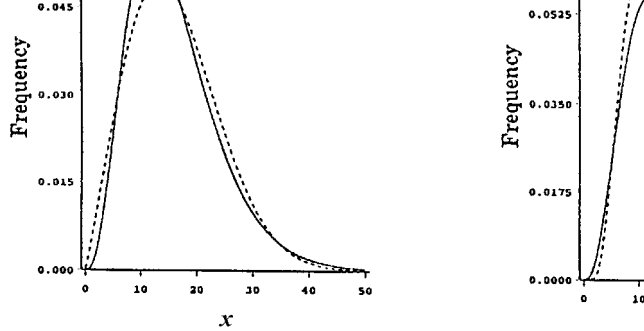


FIGURE 2.4. Plots of f (= gamma (4, 4)) against the best Weibull and lognormal models. The Weibull model that was closest to f had parameters (2, 2) with a K-L distance = 0.02009, while the best lognormal had parameters (2, 2) with a K-L distance = 0.02195. Compare these optimally parametrized models with the gamma (4, 4) distribution in Figure 2.3 (top).

K-L distance of 0.02009; this is somewhat closer than the gamma (4, 4) distribution's K-L distance of 0.04620 above. Using the same approach, the best lognormal model had parameters $\theta = 2.642$ and $\sigma^2 = 0.2838$ and a K-L distance of 0.02195. The best inverse Gaussian model had parameters $\alpha = 16.0$ and $\beta = 0.767$ and a K-L distance of 0.03726, and the approximately best Weibull model had parameters $\alpha \approx 300$, $\beta = 0.767$ and a K-L distance of approximately 0.02009. Note that the K-L distance is not sensitive to α in this case, but is quite sensitive to β (this is a numerical artifact). Thus, K-L distance indicates that the best Weibull model is closer to f than is the best lognormal (Figure 2.4). Note that the formula for the K-L distance requires knowing the true distribution f as well as the model g_i (i.e., parameter estimation has not yet been performed). Thus, the K-L distance cannot be computed for real-world problems.

These values represent *directed* distances; in the first case, $I(f, g_1) = 0.04620$, while $I(g_1, f) = 0.05552$ (in fact, $I(g_1, f) = 0.05552$ is the information loss from approximating f by g_1). The point here is that these are directed distances, and $I(f, g_1) \neq I(g_1, f)$; *nor should they be equal, because the model and the data are not interchangeable*.

These are all univariate functions; thus one could merge them on the same scale and visually compare each g_i to f ; however, this will work only in the simplest cases. In addition, if two distributions are fairly close to f , it might be difficult to decide which is closer to f by only visual inspection. Values of the K-L distance are not a measure of the mean and variance of the distributions; rather, the distributions themselves are the subject of comparison.

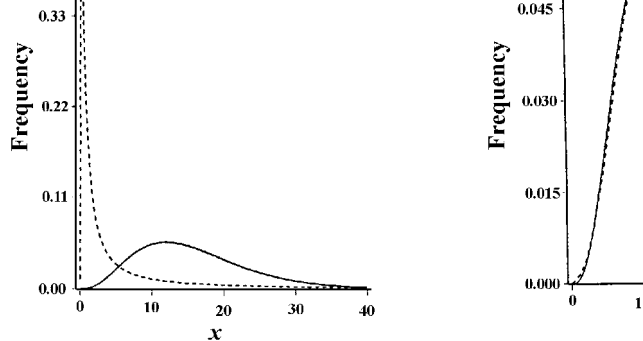


FIGURE 2.5. Plots of f ($=$ gamma (4, 4)) against the best 2-parameter and the best 3-parameter (noncentral) F distribution. The best 2-parameter is a poor approximation to f (K-L distance $=$ 1.486), while the best 3-parameter is an excellent approximation (parameters 1.322, 43.308, 18.856) with K-L distance $=$ 0.001097. Approximating models with increasing numbers of parameters typically result in better approximations than approximating models with fewer parameters.

The F distribution ($\alpha = 4$, $\beta = 10$) provided a relatively poor approximation to the gamma distribution with ($\alpha = 4$, $\beta = 4$). Even the noncentral F distribution remains a relatively poor approximation (K-L distance $=$ 1.486). However, in general, adding more parameters will result in a better approximation (e.g., the classic use of the Fourier series in the physics literature (1975) elephant-fitting problem). If we allow the addition of a third parameter (λ) in the F distribution (the noncentral F distribution), the best 3-parameter model ($\alpha = 1.322$, $\beta = 43.308$, and $\lambda = 18.856$) has a K-L distance of 0.001097; this is better than any of the other 2-parameter models (Figure 2.5). Closeness of approximation can always be improved by adding more parameters to the candidate model. When we compare models with more parameters and the associated uncertainty, then the problem of overfitting must be addressed (see Section 1.4), or overfitted models.

In the remainder of the book we will want a more general concept of a model, f , and we will use it to reflect truth or full reality. Here, reality is not a model; rather, it reflects the complex biological (and mechanical) process that generated the observed data x . For this reason, we do not want to parametrize the complex function f , because it represents full reality. It might not even have parameters in a sense that would be useful in a modeling framework. In fact, thinking that truth is parameterizable is a product of (artificial) model-based conceptualization. Sometimes we use f as full reality and let it have (conceptually) an infinite number of parameters (see Section 1.2.4). This “crutch” of infinite-dimensionality is a useful concept of reality even though it is in some unattainable place.

rise to a set of data x . Finally, we will see how this concept of $I(f)$ collapses into a nonidentifiable constant in the context

2.1.2 Truth, f , Drops Out as a Constant

The material above makes it obvious that both f and g (a model) must be known to compute the K-L distance between them. However, if only relative distance is used, this requirement drops. $I(f, g)$ can be written equivalently as

$$I(f, g) = \int f(x) \log(f(x)) dx - \int f(x) \log(g(x | \theta)) dx$$

Note that each of the two terms on the right of the above equation is a statistical expectation with respect to f (truth). Thus, the K-L distance can be expressed as a difference between two statistical expectations

$$I(f, g) = E_f [\log(f(x))] - E_f [\log(g(x | \theta))]$$

each with respect to the distribution f . This last expression provides insights into the derivation of AIC.

The first expectation $E_f [\log(f(x))]$ is a constant that depends on the unknown true distribution, and it is clearly not known (in general, f in actual data analysis). Therefore, treating this unknown constant, a measure of *relative* directed distance is possible (Kaplan and Kesavan 1992:155). Clearly, if one computed the expectation $E_f [\log(g(x | \theta))]$, one could estimate $I(f, g)$ up to a constant $E_f [\log(f(x))]$,

$$I(f, g) = C - E_f [\log(g(x | \theta))],$$

or

$$I(f, g) - C = -E_f [\log(g(x | \theta))].$$

The term $(I(f, g) - C)$ is a *relative* directed distance between f and g . $E_f [\log(g(x | \theta))]$ becomes the quantity of interest for selection. For two models g_1 and g_2 , if $I(f, g_1) < I(f, g_2)$, so g_1 is better than g_2 . $C < I(f, g_2) - C$, and hence $-E_f [\log(g_1(x | \theta))] < -E_f [\log(g_2(x | \theta))]$. Moreover, $I(f, g_2) - I(f, g_1) \equiv -E_f [\log(g_2(x | \theta))] + E_f [\log(g_1(x | \theta))]$. So we know how much better model g_1 is than model g_2 . However, we just do not know the absolute measure of how good a model is. We can identify the fact that model g_1 is better than g_2 . No estimation is involved here, but the concepts carry over to model selection. From the preceding example, where f is the true distribution, $\int f(x) \log(f(x)) dx = 3.40970$, and this term is constant.

The *relative* distances between the gamma (4, 4) and the four approximating models are shown below:

	Approximating model	Relative distance $I(f, g_i)$
g_1	Weibull distribution ($\alpha = 2, \beta = 20$)	3.455
g_2	lognormal distribution ($\theta = 2, \sigma^2 = 2$)	4.082
g_3	inverse Gaussian ($\alpha = 16, \beta = 64$)	3.469
g_4	F distribution ($\alpha = 4, \beta = 10$)	9.155

Note that the ranking of “closeness” of the four candidate models is preserved, and the relative ranking of distance between the gamma and the approximating models is unchanged, even though only relative distances are used.

Kullback-Leibler distance $I(f, g)$ is on a true ratio scale with a true zero. In contrast, $-\int f(x)(\log(g(x|\theta)))dx \equiv -E_f[\log(g(x|\theta))]$ is on an interval scale and lacks a true zero. A difference of magnitude 10 means the same thing anywhere on the scale. Thus, $D = 10 = 12 - 2$ means a difference of 10 means the same thing anywhere on the scale. Similarly, $10 = V_1 - V_2$, regardless of the size of V_1 and V_2 .

The calculation of the two components of K-L distance is affected based on a sample size of 1. If the sample size were 100, the two components would be 100 times larger, and the difference between the two components would also be 100 times larger. For example, $\int f(x) \log(f(x))dx = 3.40970 \times 100 = 340.970$ and $E_f[\log(g_1(x|\theta))]$ for Weibull = $3.45591 \times 100 = 345.591$. Thus, the difference between the two components of K-L distance would be 4.620; the *relative* difference would be the same when sample size is large. A large sample size magnifies the differences between search hypotheses and the models used to represent them. **size conveys a wide variety of advantages in making v**

Typically, as in the example above, the analyst would select a set of *a priori* candidate models $g_i(x|\theta)$ and want to select the best model as a basis for data analysis and inference. Definition of “best” is based on the principle of parsimony and the related concept of a best approximation. In data analysis, the parameters in the various candidate models are unknown and must be estimated from the empirical data. This represents a key distinction from the material above, since one usually has known or estimated parameters, denoted by $g_i(x|\hat{\theta})$. In this case, the relative of the relative directed distances between the unknown true model $f(x)$ and the various candidate models $g_i(x|\hat{\theta})$. Then, knowing the relative distance from each $g_i(x)$ to $f(x)$, we select the candidate model that is *estimated* to be closest to truth for inference (Figure 2.2). In other words, the model with the smallest estimated, *relative* distance is selected as the approximating model that loses the least information.

irrelevant for comparison. (Summary, it is interesting to note that the likelihood function also involves an additive constant that is the same for all models; this term is known, but generally ignored, since it does not affect the relative ordering of models.) In practice, we can obtain only an *estimate* of the distance from each approximating model $g_i(x | \hat{\theta})$ to f .

2.2 Akaike's Information Criterion: 1973

Akaike's (1973) seminal paper proposed the use of the K-L information or distance as a fundamental basis for model selection. The K-L distance cannot be computed without full knowledge of the truth (f) and the parameters (θ) in each of the candidate models. Akaike found a rigorous way to estimate K-L information, based on the log-likelihood function at its maximum point.

Given a parametric structural model there is a unique parameter value that minimizes K-L distance $I(f, g)$. This (unknown) minimum parameter depends on truth f , the model g through its structure, the sample space, and the sample space (i.e., the structure and nature of the data collected). In this sense there is a "true" value of θ under which the K-L information loss is minimized at θ_0 . If one somehow knew the true g was, in fact, the K-L best model, then the MLE $\hat{\theta}$ would be the best property of the model $g(x|\theta_0)$ as the minimizer of K-L, a key and important feature involved in the derivation of AIC (Chapter 2).

In data analysis the model parameters must be estimated, and there is substantial uncertainty in this estimation. Models based on estimates of parameters, hence on $\hat{\theta}$ not θ , represent a major distinction from the case where parameters would be known. This distinction affects how K-L distance as a basis for model selection. The difference between having the true θ_0 (we do not) and having the estimate $\hat{\theta}$ (we do) is quite important.

Selection Target

Akaike (1973, 1974, 1985, 1994) showed that the critical quantity for the applied K-L model selection criterion was to estimate

$$E_y E_x [\log(g(x|\hat{\theta}(y)))],$$

where x and y are independent random samples from the truth f , and both statistical expectations are taken with respect to f . The double expectation, both with respect to truth f , is the key feature of the selection approaches, based on K-L information.

over the set of K models considered.

It is tempting to just estimate $E_y E_x [\log(g(x|\hat{\theta}(y)))]$ $\log(\mathcal{L}(\hat{\theta})|data)$ for each model g_i . However, Akaike (1973) found that the maximized log-likelihood is biased upward as an estimator of the selection target (above). He also found that under certain conditions (which are important, but quite technical) this bias is of order $1/K$, the number of estimable parameters in the approximation. This is an asymptotic result of fundamental importance.

The Key Result

Thus, an approximately unbiased estimator of

$$E_y E_x [\log(g(x|\hat{\theta}(y)))]$$

for large samples and “good” models is

$$\log(\mathcal{L}(\hat{\theta}|data)) - K.$$

This result is equivalent to

$$\log(\mathcal{L}(\hat{\theta}|data)) - K = \text{constant} - \hat{E}_{\hat{\theta}}[I(f|g)]$$

where $\hat{g} = g(\cdot|\hat{\theta})$.

The bias-correction term ($K =$ the number of estimable parameters) is a special case of a more general result derived by Tibshirani (1986) described in the following section and in Chapter 7. **Akaike’s information criterion: the relation between the relative expected K-L distance and the log-likelihood has allowed major practical and theoretical advances in model selection and the analysis of complex data sets** (Akaike 1973, Bozdogan 1987, and deLeeuw 1992).

Akaike’s Information Criterion

Akaike (1973) then defined “*an information criterion*” by multiplying $\log(\mathcal{L}(\hat{\theta}|y)) - K$ by -2 (“taking historical reasons into account”)

$$\text{AIC} = -2 \log(\mathcal{L}(\hat{\theta}|y)) + 2K.$$

This has become known as “**Akaike’s information criterion**”.

Thus, rather than having a simple measure of the distance between two models (i.e., the K-L distance), one has instead the expected, relative distance between the fitted model and the true mechanism (perhaps of infinite dimension) that act on the observed data.

The expression $\log(\mathcal{L}(\hat{\theta}|y))$ is the numerical value of the log-likelihood at the maximum point (see Section 1.2.2). This maximum point is

clear as to what the correct count should be (see below for
els). In some types of models there are some parameters that
estimable from the data, and these should not be counted in
can occur in the analysis of count data where a cell has no ob
a parameter that is identifiable becomes nonestimable for
timability can also arise due to inherent confounding (e.g.
and f_i in certain band recovery models of Brownie et al. 1986).
one computes AIC for each of the candidate models and selects
the smallest value of AIC. It is this model that is estimated to be
the unknown reality that generated the data, from among the models
considered. This seems a very natural, simple concept; selecting the
imating model that is estimated, on average, to be closest to the
Basing AIC on the expectation (over $\hat{\theta}$) of $E_x[\log(g(x)|\hat{\theta})]$ has
with a cross-validation property for independent and identically
samples (see Stone 1977, Stoica et al. 1986, Tong 1994). It also
show that AIC asymptotically coincides with generalized likelihood
subset regression (also see review by Atilgan 1996).

**Of course, models not in the set remain out of consideration. This is
useful in selecting the best model in the set; however, if the set is
very poor, AIC will still select the one estimated to be the best. The
relatively best model might be poor in an absolute sense. It must
be made to ensure that the set of models is well fitted to the data.**

$I(f, g)$ can be made smaller by adding more known parameters in the approximating model g . Thus, for a fixed f , the
addition of parameters in a model g_i will allow it to be closer to f .
when these parameters must be estimated (rather than being known).
further uncertainty is added to the *estimation* of the relationship.
some point, the addition of still more estimated parameters will have a
positive from desired effect (i.e., to reduce $E_{\hat{\theta}}[I(f, \hat{g})]$ as desired).
the estimate of the relative K-L distance will increase because of the
estimated parameters that are not really needed to achieve the desired
phenomenon can be seen by examination of the information criterion
minimized,

$$\text{AIC} = -2 \log(\mathcal{L}(\hat{\theta}|\mathbf{y})) + 2K,$$

where the first term on the right-hand side tends to decrease as more
rameters are added to the approximating model, while the second term
gets larger as more parameters are added to the approximating model.
is the tradeoff between bias and variance or the tradeoff between underfitting
and overfitting that is fundamental to the principle of parsimony (see
1.4.2). Some investigators have considered K to be a measure of model
ity,” but this is unnecessary, though not irrational. We can write

expected K-L information, thus parsimony arises as a by-product of the approach. Further books and papers on the derivation of AIC are Akaike (1983, 1989), Linhart and Zucchini (1986), Bozdogan (1993), and Burnham and Anderson (1998).

Usually, AIC is positive; however, it can be shifted by a constant and some shifts can result in negative values of AIC. Computing AIC from regression statistics (see Section 1.2.2) often results in negative values. In our work, we have seen minimum AIC values that range from negative numbers to as high as 340,000. **It is not the absolute size of the AIC values that matters; it is the relative values over the set of models considered. It is the differences between AIC values (Section 2.5), that matter.**

The material to this point has been based on likelihood methods, a very general approach. In the special case of least squares regression with normally distributed errors, and apart from an arbitrary additive constant, AIC can be expressed as a simple function of the residual sum of squares.

The Least Squares Case

If all the models in the set assume normally distributed errors with constant variance, then AIC can be easily computed from least squares regression statistics as

$$\text{AIC} = n \log(\hat{\sigma}^2) + 2K,$$

where

$$\hat{\sigma}^2 = \frac{\sum \hat{\epsilon}_i^2}{n} \text{ (the MLE of } \sigma^2 \text{),}$$

and $\hat{\epsilon}_i$ are the estimated residuals for a particular candidate model. A common mistake with LS model fitting, when computing AIC, is to use σ^2 from the computer output, instead of computing the MLE of σ^2 . **Also, for LS model fitting, K is the total number of estimated parameters, including the intercept and σ^2 .**

Thus, AIC is easy to compute from the results of LS estimation for linear models and is now included in the output of many software packages for regression analysis. However, the value of K is sometimes computed incorrectly because either β_0 (the intercept) or σ^2 (or both) is not included in determining K .

The fact that AIC is an estimate only of relative expected K-L information is almost unimportant. It is the fact that AIC is only an estimate of the distances from each model g_i to f that is less than ideal. We must recognize that there is usually substantial uncertainty as to which model is best for a given data set. After all, these are stochastic biological processes with relatively high levels of uncertainty.

the basic underlying structure of the model that fits the data “speak” is of interest to both biologists and statisticians. Coming from empirical data. The data then help determine the (order or dimension) of the approximating model used for inference. To determine what effects or factors are justified. In this sense, the given data set are conditional on sample size. We must be careful. If more data were available, then further effects could probably be supported. “Truth” is elusive; model selection tells us what is best supported, not what full reality might be.

Akaike (1973) multiplied the bias-corrected log-likelihood by -2 for “technical reasons” (e.g., it is well known that -2 times the log-likelihood ratio of two maximized likelihood values is asymptotically distributed as a chi-square under certain conditions and assumptions). The term -2 occurs in many contexts, so it was not unreasonable that Akaike performed the multiplication to get his AIC. Two points frequently arise, and both are addressed here. First, the model associated with the minimum AIC is not necessarily the best model if the bias-corrected log-likelihood (i.e., $\log(\mathcal{L}) - K$) is negative, such as -34 , or -51.3 , or any other negative number. Thus, the model is not changed by the multiplication of both terms by any negative number. Akaike merely chose -2 . Second, some investigators have not recognized the link between K-L information and AIC and believed, that the -2 in the second term in AIC was somehow “arbitrary”. Other values should also be considered. This error has led to confusion in the technical literature; clearly, K is the asymptotic bias, not arbitrary. Akaike chose to work with $-2 \log(\mathcal{L})$, rather than $\log(\mathcal{L})$. The term $+2K$ is theoretically correct, for large sample sizes. If the terms (the log-likelihood and the bias correction) are multiplied by a negative constant, the model where the criterion is minimized is the same, and there is nothing arbitrary.

It might be argued that we should have merely defined $AIC = -2l + 2K$ (where l is the log-likelihood of the model); then $AIC = -2l + 2K$, making the criterion look simpler. However, if we may have advantages, we believe that the full notation $AIC = -2 \log(\mathcal{L}(\theta|x, g_i)) + 2K$ and helps in understanding exactly what is meant. The use of the abbreviations such as $\log(\mathcal{L}(\theta|x, g_i))$, makes it explicit that the log-likelihood is a function of (only) the parameters (θ), while the data (x , say multinomial) must be given (i.e., known). These distinctions are important when we introduce the concept of a likelihood of the data: $\mathcal{L}(g_i|data)$. Both concepts are fundamental and useful. The notation in this book and the notation serves an important purpose.

If the approximating models in the candidate set are not too far from the truth, then Takeuchi’s information criterion (TIC) is an alternative. If the true model is quite large. AIC is a special case of TIC, and as such, AIC is a good approach to the estimation of relative expected K-L distance.

At one point in Akaike's derivation of an estimator of $K \cdot L$, the assumption that the model set included f (full reality) was subject of attention and criticism. Akaike maintained that was asymptotically unbiased and free from any notion of model or that such a true model was required to be in the models. This section will indicate that such claims were just another insight into the concept of parsimony. The key important, little-known paper (in Japanese) by Takeuchi came just 3 years after Akaike's initial breakthrough in 1973.

Takeuchi (1976) provides a very general derivation of a criterion, without taking expectations with respect to g . His criterion, TIC (Takeuchi's information criterion) and was thought of where the candidate models were not particularly close to g . TIC has a more general bias-adjustment term to allow for a justed to be an asymptotically unbiased estimate of relative information,

$$\text{TIC} = -2 \log(\mathcal{L}) + 2 \cdot \text{tr}(J(\theta)I(\theta)^{-1})$$

The $K \times K$ matrices $J(\theta)$ and $I(\theta)$ involve first and second derivatives of the log-likelihood function, and "tr" denotes the trace function. One might consider always using TIC and worrying about the quality of the models in the set of candidates. This consideration has issues that are problematic. First, one must *always* worry about the set of approximating models being considered; this is a shortcut. Second, using the expanded bias adjustment term in estimation of the elements of the matrices $J(\theta)$ and $I(\theta)$ (see in Chapter 7). Shibata (1999) notes that estimation errors in these matrices can cause instability of the results of model selection where a candidate model has $K = 20$ parameters. The matrices $J(\theta)$ and $I(\theta)$ are of dimension 20×20 , and reliable estimation of each matrix will be difficult unless sample size is very large so that $\text{tr}(J(\theta)I(\theta)^{-1})$ itself has a very simple parsimonious expression. This is an interesting and important general result.

Thus, AIC is an approximation to TIC, where $\text{tr}(J(\theta)I(\theta)^{-1})$ approximation is excellent when the approximating model is good, comes poor when the approximating model is a poor. For models that are poor, the first term, $-2 \log(\mathcal{L})$, dominates the criterion. If the model is poor and this term will tend to be relatively large, compared to the second term. Thus, with the final approximation that $\text{tr}(J(\theta)I(\theta)^{-1}) \approx K$, we can see that AIC is an asymptotically unbiased estimator of $K \cdot L$ (K-L information, derived without assuming that full reality is in the set of candidate models). Very important contribution to the literature, it has rarely been seen ap

2.4 Second-Order Information Criterion:

While Akaike derived an estimator of K-L information poorly if there are too many parameters in relation to the sample size (Sugiura 1978, Sakamoto et al. 1986). Sugiura (1978) derived a variant of AIC that he called c-AIC.

A Small Sample AIC

Hurvich and Tsai (1989) further studied this small-sample bias adjustment, which led to a criterion that is called AIC_c .

$$AIC_c = -2 \log(\mathcal{L}(\hat{\theta})) + 2K \left(\frac{n}{n - K - 1} \right)$$

where the penalty term is multiplied by the correction factor $\frac{n}{n - K - 1}$. This can be rewritten as

$$AIC_c = -2 \log(\mathcal{L}(\hat{\theta})) + 2K + \frac{2K(K + 1)}{n - K - 1}$$

or, equivalently,

$$AIC_c = AIC + \frac{2K(K + 1)}{n - K - 1},$$

where n is sample size (also see Sugiura 1978).

Unless the sample size is large with respect to the number of parameters, use of AIC_c is recommended.

AIC_c merely has an additional bias-correction term. If n is large relative to K , then the second-order correction is negligible and AIC_c is close to AIC . Findley (1985) noted that the study of bias correction is complex; the exact small-sample bias-correction term varies with the distribution (e.g., normal, exponential, Poisson). Bedrick and Tsai (1990) provided a refinement, but it is more difficult to compute (also see Hurvich et al. 1995a and b, and Hurvich et al. 1990). While AIC_c is a good approximation for Gaussian assumptions for linear models (fixed effects), Bedrick and Tsai found this second-order approximation to the K-L distance to be poor for product multinomial models. **Generally, we advocate the use of AIC_c when the ratio n/K is small (say < 40).** In reaching a decision between AIC vs. AIC_c , one must use the value of K for the highest (global) model in the set of candidates. If the ratio n/K is large, then AIC and AIC_c are similar and will strongly tend to select the same model. One must use either AIC or AIC_c consistently in a given analysis.

2.5 Modification of Information Criterion for Overdispersed Count Data

In general, if the random variable n represents a count discrete distribution (e.g., Poisson or binomial), it has a mean $\mu(\theta)$, and a known theoretical variance function, $\sigma^2(\theta)$. In a model of overdispersed data the expectation of n is $\mu(\theta)$ and the variance model must be generalized, for example using a factor, e.g., $\gamma(\theta)\sigma^2(\theta)$. The form of the factor $\gamma(\theta)$ can be determined by theoretical considerations and can be complex (see, e.g., Nelder 1989). Overdispersion factors typically are small, usually above 1 to perhaps 3 or 4 if the model structure is correct. Overdispersion is due to small violations of assumptions such as independence and homogeneity over individuals. Hence, a first approximation to overdispersion is to use a simple constant c in place of $\gamma(\theta)$, generalized to more than one c for different partitions of the data.

Count data have been known not to conform to simple models based on binomial or multinomial distributions (e.g., Fisher 1949, Armitage 1957, and Finney 1971). There are many statistical models for count data (e.g., Poisson, binomial, multinomial). In these, the sampling variance is theoretically determined (e.g., for the Poisson model, $\text{var}(n) = E(n)$; for the binomial model, $\text{var}(\hat{p}) = p(1 - p)/n$). If the sampling variance exceeds the theoretical (model-based) variance, the situation is called “overdispersion.” There is here on a lack of independence in the data leading to “extrabinomial variation.” Eberhardt (1978) provides a collection of examples of issues in the biological sciences. For example, Canada geese frequently mate for life, and the pair behaves almost as a single unit rather than as two independent “trials.” The young of some species stay with the parents for a period of time, which can also cause dependence of individual responses. Further reasons for overdispersion in ecological systems include species whose members exist in schools or flocks. Members of such populations can be expected to have positive correlations between individuals within the group; such dependence causes overdispersion. Another type of overdispersion stems from parameter heterogeneity, where individuals have unique parameters rather than the same parameter values. This is a problem in applying probability to all individuals.

The estimators of model parameters often remain unbiased in the presence of overdispersion, but the model-based theoretical variance estimates are biased (McCullagh and Nelder 1989). To properly cope

for the analysis of overdispersed data (also see Whitman and Pregibon 1985, Moore 1987, and McCullagh and Nelder 1999a). Hurvich and Tsai (1995b) provide information on overdispersed data.

Cox and Snell (1989) discuss modeling of count data and a useful approximation is based on a single variance inflation factor. The variance inflation factor can be estimated from the goodness-of-fit chi-square statistic and its degrees of freedom,

$$\hat{c} = \chi^2/\text{df}.$$

The variance inflation factor should be estimated from the data. Cox and Snell (1989) assert that the simple approach of a single variance inflation factor should often be adequate, as opposed to the task of seeking a detailed model for the $\gamma(\theta)$. In a study of several approaches on five data sets, Liang and McCullagh (1993) found that the use of a single \hat{c} for overdispersion was clearly better than use of a single \hat{c} in each model examined.

Given \hat{c} , empirical estimates of sampling variances ($\text{cov}_e(\hat{\theta}_i, \hat{\theta}_j)$) can be computed by multiplying the theoretical (model-based) variances and covariances by \hat{c} . This approach has long been used; see, e.g., Finney 1971). These empirical estimates (i.e., $\hat{c} \cdot \widehat{\text{var}}_i(\hat{\theta}_i)$) must be treated as having the degrees of freedom to compute \hat{c} for purposes of setting confidence limits (or intervals). **The number of parameters (K) must include one for the variance inflation factor, if used.** Generally, quasi-likelihood (i.e., use of $\hat{c} > 1$) are made only if some distinct lack of fit is evident (for example, if the observed significance level $P \leq 0.05$ and the goodness-of-fit degrees of freedom ≥ 10 , as rough guide).

We might expect $c > 1$ with real data but would not expect $c > 4$ if model structure is acceptable and only overdispersion is present (Eberhardt 1978). Substantially larger values of c (say, 6–10) are due partly by a model structure that is inadequate; that is, the model does not account for an acceptable amount of variation in the data. The methods of variance inflation are most appropriate only when the structural adequacy of the model has been achieved. The variance inflation factor can be computed only for the global model; one should not multiply the estimates of this variance inflation factor for each of the parameters in the set. The issue of the structural adequacy of the model is at the heart of good data analysis (i.e., the reliable identification of signal versus residual variation in the data). Patterns in the goodness-of-fit statistics (Pearson χ^2 or G-statistics) might be an indication of structural inadequacy of the model. Of course, the biology of the organism in question is also important.

When data are overdispersed and $c > 1$, the proper model is $\mathcal{L}(\theta|\mathbf{x}, \mathbf{g}; c)$ (not just $\log(\mathcal{L})$). Principles of quasi-likelihood suggest substituting c for 1 in the formula for AIC and AIC_c ; we denote these modifications by (Lebreton et al. 1992)

$$QAIC = - \left[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c} \right] + 2K,$$

and

$$\begin{aligned} QAIC_c &= - \left[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c} \right] + 2K + \frac{2K(K+1)}{n-K-1} \\ &= QAIC + \frac{2K(K+1)}{n-K-1}. \end{aligned}$$

If an overdispersion factor is estimated, then one parameter is added to the model, so K is replaced by $K+1$. Of course, when no overdispersion exists, then $c = 1$ and the formulas for $QAIC$ and $QAIC_c$ reduce to AIC and AIC_c , respectively. Lebreton et al. (1994) found that these criteria performed well in product of capture–recapture data in the presence of differing levels of overdispersion.

One must be careful when using some standard software (e.g., SAS GENMOD), since they were developed some time ago and do not have a variance testing mode (i.e., adjusting χ^2 test statistics by \hat{c} to account for overdispersion). In some cases, a separate estimate of c is made for each model, and the test statistics and covariances are multiplied by this model-specific estimate of c . This is an inflation factor. Some software packages compute an estimate of c for each model, thus making the correct use of model selection criteria more difficult. One is careful. Instead, we recommend that the global model be used for the estimation of a single variance inflation factor c .

Overdispersed Count Data: A Review

Try to ensure that the structural part of the data is well described by a global model.

If there is biological reason to suspect overdispersion, the overdispersion parameter c can be estimated as χ^2/df , using the global model.

If overdispersion is present, the log-likelihood of the data and the model, should be computed as

$$\frac{\log(\mathcal{L}(\theta|\mathbf{x}, \mathbf{g}_i))}{\hat{c}}.$$

The number of parameters K is now the number of parameters in the model, plus one to account for the estimation of the overdispersion parameter c .

The estimated overdispersion parameter should generally be used to adjust the test statistics. Otherwise, some structural lack of fit is probably entering into the model. If $\hat{c} < 1$, just use $c = 1$.

and covariances should also be adjusted using \hat{c} from the data, so there are few degrees of freedom left.

AIC for Overdispersed Count Data

Model selection should use either

$$\text{QAIC} = -[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c}] + 2K,$$

or

$$\begin{aligned} \text{QAIC}_c &= -[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c}] + 2K + \frac{2K(K)}{n - K} \\ &= \text{QAIC} + \frac{2K(K + 1)}{n - K - 1} \end{aligned}$$

The variance–covariance matrix should be multiplied by the overdispersion parameter \hat{c} (i.e., $\hat{c}(\text{cov}(\hat{\theta}_i, \hat{\theta}_j))$).

Some commercial software computes AIC, while AIC_c and no general software package computes QAIC or QAIC_c. In cases, AIC, AIC_c, QAIC, and QAIC_c can be computed from the material that is output from standard computer packages (or least squares estimation). In general, we recommend the information-theoretic criterion for count data, and we will use it in the practical examples in Chapter 3. Of course, often the parameter is near 1, negating the need for quasi-likelihood, just as often the ratio n/K is large, negating the need for the correction term in AIC_c. AIC, AIC_c, and QAIC_c are all estimates of K-L information. We often use the generic term “AIC” to refer to all three criteria.

2.6 AIC Differences, Δ_i

AIC, AIC_c, QAIC_c, and TIC are all on a relative (or information) scale, strongly dependent on sample size. Simple differences of the estimates of $E_{\hat{\theta}}[\hat{I}(f, g_i)] - \min E_{\hat{\theta}}[\hat{I}(f, g_i)]$, where the \min is over the estimated parameters and \min is over the models.

The larger Δ_i is, the less plausible it is that the fitted model is the K-L best model, given the data x . Some rough rules of thumb for Δ_i are particularly useful for nested models:

Δ_i	Level of Empirical Support of Model
0-2	Substantial
4-7	Considerably less
> 10	Essentially none.

AIC differences,

$$\Delta_i = \text{AIC}_i - \text{AIC}_{\min},$$

over all candidate models in the set. We use the term “ Δ_i ” in a generic sense here to mean AIC, AIC_c , QAIC_c , or TIC. These Δ_i values estimate the relative expected K-L differences between f_i and the best model. Δ_i values are easy to interpret and allow a quick comparison of candidate models and are also useful in computing Akaike weights (2.9). The model estimated to be best has $\Delta_i \equiv \Delta_{\min} \equiv 0$.

Models with $\Delta_i > 10$ have either essentially no support or are omitted from further consideration, or at least those models with substantial explainable variation in the data. These guidelines are useful if the sample size is small (even as many as 100), but may break down in large samples where there may be thousands of models. The guideline values are larger for nonnested models, and more research is needed (Burnham and Linhart 1988). If observations are not independent, but are approximately independent, then these simple guidelines cannot be expected to hold. If the log-likelihood is corrected for overdispersion in count data, then the guidelines above will be useful.

As an example, candidate models g_1 , g_2 , g_3 , and g_4 have AIC values of 3,400, 3,560, 3,380, and 3,415, respectively. Then one would choose g_3 as the best single model as the basis for inference because it has the lowest AIC value. Because these values are on a relative (interval) scale, subtract, say, 3,380 (the minimum of the 4 values) from each value to have the following rescaled AIC values: 20, 180, 0, and 75. Rescaling does not change the ranks of the models, nor the relative differences in the AIC values. People are often surprised that Δ_i values are not as important, when the associated AIC values that led to the models are in the order of 97,000 or 243,000.

AIC Differences

It is not the absolute size of the AIC value, it is the relative size, and particularly the AIC differences (Δ_i), that are important.

An individual AIC value, by itself, is not interpretable on a constant (interval) scale. AIC is only comparative, relative to the best model in the model set; thus such differences Δ_i are very important.

We can say with considerable confidence that in real life, with several or more models and large sample size (say $n > 100$), a model having $\Delta_i = 20$, such as model g_4 , would be a good approximating model for the data at hand.

K-L best model. Consider Δ_i values for 7 models as 0, 1, 2, 3, 4, 5, and 7.3. An important question is, how big a difference must be asked in the sense of when a model is not to be considered selected best model as plausibly the actual K-L best model used, for the sample size and data at hand. The question has no answer; it is like asking how far away from an MLE $\hat{\theta}$ an alternative θ must be (assuming that the model is a good model) to be that an alternative θ is unlikely as “truth.” This question could be answered with a confidence (or credibility) interval on θ based on the uncertainty. A conventionally accepted answer here is that a model is not selected if it is further away than $\pm 2 \widehat{\text{se}}(\hat{\theta})$ (there is a fundamental arbitrariness in a procedure). Relative scaling of alternative models can be done using Akaike weights (Section 2.9) and evidence ratios (Section 2.10).

2.7 A Useful Analogy

In some ways, selection of a best approximating model is like a car racing or other similar contests. The goal of such a race is to select the (fastest) car/driver combination, and the data represent the results of the race (e.g., the Indianapolis 500 in the USA, the 24 Heures du Mans in France). Only a relatively few car/driver combinations “qualify,” based on the results (e.g., 33 cars at Indianapolis)—this is like the set of candidates for a model. Only certain models “qualify,” based on the science of the situation. The race can be chaotic if all car/driver combinations with an interest in winning are included, just as it makes little sense to include a very large number of candidates in a set of candidates (and risk Freedman’s paradox). Cars that do not win, even though they might indeed have been the best, are not considered. Similarly, models, either good or bad, that do not qualify are not considered. Candidates remain out of consideration.

At the end of the race the results provide a ranking of the car/driver combination, from first to last. Furthermore, if a measure of quality is available (e.g., elapsed time for each finisher), a “scaling” can be considered. Clearly, the primary interest is in the “winner” or “which was the first”; this is like the model selection problem. The AIC value. This answers the question, “Which is best in the current data set?” It could differ for another (future) race or another data set, but the current results are unavailable to us.

Some (secondary) interest exists in the question, “Which is second place?” and in particular, was second place only thousandths of a second behind the winner or 5 minutes behind? The race time results answer these questions, as do the Δ_i values in model selection. In the

first place winner.); while in the second case, the inference is based on a single car/driver combination as the clear best (with a finish). The finishing times provide insights into the third etc. In trying to understand the performance of car/driver combinations, there is considerable information from both the rankings and the AIC values (both the ranks and the Δ_i). In 2.9 and 2.10 we will see how the Δ_i can be used to estimate the relative performance of the models. Note that the winner is of little interest because of temperature differences and other variables; only the *relative* times for a given race are of interest. Similarly, the absolute values of AIC are also of little interest; they reflect sample size and some constants, among other things. The values of the maximized log-likelihood (i.e., $\log(\mathcal{L}(\hat{\theta}|x))$) vary from sample to sample. However, all comparisons of models are based on the same data, so this sample-to-sample variation is irrelevant. Comparing log-likelihood values across data sets is like comparing times across races when some races are 500 miles whereas others are 400 miles.

The winner of the race is clearly the best for the particular data set. To make a broader inference concerning races for an entire season (i.e., ranks) from several races can be pooled or weighted. In general, inferences beyond a single observed data set can sometimes be made by some type of model averaging using, for example, the bootstrap (details in Chapters 4 and 5) and the incorporation of uncertainty in estimators of precision.

The race result might not always select the best car/driver combination because the fastest qualifying car/driver may have had bad luck (e.g., engine failure) and finished well back from the leader (if the race is not a model selection one has only one realization of the stochastic process). The *estimated* relative distance as the basis for the selection of a model (a winner). If the same race is held again with the same data, the winner and order of finishers are likely to change somewhat. If a new sample of data could be obtained, the model ranks would change somewhat.

To carry the analogy a bit further, data dredging would be like watching a race as cars dropped out and others came to the front. The continually shifts the bet and predicted winner, based on the current lead at any point in time (i.e., an unfair advantage). In general, the prediction would surely be improved, but the rules of play would be altered! Alternatively, the definition of winning might not be based on the initiation of the race. Only after the race are the winners based, in part, on who they think “ought” to win). Then, one might question the applicability of this specific prediction to other races. In general, “new rules” when data dredging has been done. That is,

all the information and insights they might provide. However, leading to data dredging should be revealed, and results discussed in this light.

Many realize that there is considerable variation in cars as to race and track to track. Similarly, many are comfortable that there is often considerable sampling variation (uncertainty in estimate of a parameter from data set to data set. Similarly, if many races could be taken, the estimated best model (car/driver) could vary from sample to sample (or race to race). Both components of model selection uncertainty should ideally be incorporated into measures of precision.

2.8 Likelihood of a Model, $\mathcal{L}(g_i | data)$

While the AIC differences Δ_i are useful in ranking the models, it is not possible to quantify the plausibility of each model as being the actual best model. This can be done by extending the concept of the likelihood of a model, given both the data and model, i.e., $\mathcal{L}(\theta | x, g_i)$, to the concept of the likelihood of the model given the data, hence $\mathcal{L}(g_i | x)$. Such quantities can be used in making inferences concerning the relative strength of evidence for each model in the set.

Likelihood of a Model, Given Data

The likelihood of model g_i , given the data, is simply the likelihood of the model in the set:

$$\mathcal{L}(g_i | x) \propto \exp \left(-\frac{1}{2} \Delta_i \right),$$

where “ \propto ” means “is proportional to.” Such likelihoods can be used to make inferences concerning the relative strength of evidence for each model.

Akaike (see, e.g., Akaike 1983b) advocates the above approach to the relative likelihood of the model, given the MLEs of model parameters on the same data. Such quantities can also be expressed as

$$C \mathcal{L}(\hat{\theta} | x, g_i) e^{-K},$$

where C is an arbitrary constant.

Model Probabilities

To better interpret the relative likelihood of a model, the set of R models, we normalize the $\mathcal{L}(g_i|x)$ to be a set of weights,” w_i , adding to 1:

$$w_i = \frac{\exp(-\frac{1}{2}\Delta_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta_r)}.$$

The w_i depend on the entire set; therefore, if a model is added during a post hoc analysis, the w_i must be recomputed for the newly defined set.

This idea of the likelihood of the model given the data and model weights, has been suggested for many years by Akaike (1978b, 1979, 1980, 1981b and 1983b; also see Bozdoga et al. 1991) and has been researched some by Buckland (1984). This name will herein apply also when we use AIC_c , $QAIC_c$. **A given w_i is considered as the weight of evidence in favor of model i . The actual K-L best model for the situation at hand given the set of R models must be the K-L best model of that set of R models.** That there are only R models and one of them must be best is obvious. It is convenient to normalize the relative likelihoods to sum to 1.

For the estimated K-L best model (let this be model g_{min}) for that model $\exp(-\frac{1}{2}\Delta_{min}) \equiv 1$. The odds for the i^{th} model relative to the K-L best model are thus $\exp(-\frac{1}{2}\Delta_i)$ to 1, or just the “odds.” It is convenient to reexpress such odds as the set of Akaike weights, w_i . A Δ_i is, the smaller the w_i , and the less plausible is model i . The K-L best model for f based on the design and sample size is g_{min} . The weights provide an effective way to scale and interpret the relative likelihoods. The weights also have other important uses and interpretations in the following chapters.

In general, likelihood provides a good measure of evidence about parameter values, given a model and data (Akaike 1977). We think that this concept extends to evidence about the best model, given a set of models. That is, evidence for the best model is represented by the likelihood of a model.

of these models might be the K-L best model for the data (compelled by a certain aspect of information theory itself (Jessop 1995)). Let τ_i be the prior probability that model i is the best. Lacking any prior information, we set the τ_i all equal, and $1/R$. In fact, doing so places all R of the models on an equal footing as the K-L best model.

If there is prior information or belief, this opens the door to subjective probabilities. Ignoring any model redundancy (this subject is discussed in section 4.6), τ_i is our prior state of information or belief that model i , given the data, provides the K-L best model for the design and the data. This is a deceptively complex issue, as it relates both to idealized models and to approximations to truth and to expected model fitting trade-offs and sampling variances.

To us it seems impossible to have any real prior basis for a subjective differential assessment of the τ_i (other than on how the models are structurally interrelated or partially redundant). Using the principle of Jaynes (1957) we should take the τ_i to represent our state of uncertainty about all unknown aspects of the probability distribution of the τ_i . Thus we determine the τ_i that maximize the entropy of the τ_i subject to constraints that express whatever information (in the form of subjective probabilities) we have about the distribution. In the “no information” case, where we have is that $\sum \tau_i = 1$ (plus the essential $0 < \tau_i < 1$), the maximum entropy (hence maximum uncertainty) prior is then $\tau_i \equiv 1/R$. We do not feel to delve into the aspects of information theory underlying the maximum entropy principle. This principle is fundamentally tied to the concept of entropy and to information theory and can be used to justify the use of Bayesian priors—when they exist. The interested reader should consult Jaynes and Kesavan 1992, or the less technical Jessop 1995.]

Given any set of prior probabilities (the τ_i), generalized Akaike weights can be given by

$$w_i = \frac{\mathcal{L}(g_i|\mathbf{x})\tau_i}{\sum_{r=1}^R \mathcal{L}(g_r|\mathbf{x})\tau_r}.$$

There may be occasions to use unequal prior probabilities, as discussed above. However, in general, by Akaike weights we mean the weights w_i without the τ_i (this assumes $\tau_i = 1/R$).

The inclusion of prior probabilities (τ_i) in the w_i is a Bayesian approach. The full Bayesian approach to model selection requires specifying τ_i on the model and a prior probability distribution on the parameters θ_i for each model. Then the derivation of posterior results (usually achievable only by Markov chain Monte Carlo

$$\mathcal{L}(g_i|x)/\mathcal{L}(g_j|x)$$

or, equivalently, the ratio of Akaike weights w_i/w_j . Such ratios are commonly used, and we will term them **evidence ratios**. Such ratios provide the evidence about fitted models as to which is better in the K-L sense.

In particular, there is often interest in the ratio w_1/w_j , where i indexes the estimated best model and j indexes the rest of the models. Such ratios are not affected by any other model, hence do not depend on the full set of R models—just on models i and j . These evidence ratios provide evidence to all other models besides i and j .

$0.431/0.024 = e^{(5.8/2)} = 18$, and we must conclude that model 6 is the K-L best model; the evidence here is reasonable for model 6.

There is a striking nonlinearity in the evidence ratios as a function of Δ_j values. Consider the ratio $w_1/w_j (\equiv w_{min}/w_j)$,

$$\frac{w_1}{w_j} \equiv \frac{1}{e^{-1/2\Delta_j}} \equiv e^{1/2\Delta_j}$$

in the comparison of the evidence for the best model versus model j . Then, we have the following table:

Δ_j	Evidence ratio
2	2.7
4	7.4
8	54.6
10	148.4
15	1,808.0
20	22,026.5

This information helps to justify the rough rules of thumb for the evidence for models being the best K-L model in the set provided some likelihood-based rules similar to these over time (Edwards (1992) and Royall (1997) for additional perspectives on evidence in a likelihood framework).

People may, at first, be frustrated that they do not have a simple point that provides a simple dichotomy to indicate what is “significant” under the Neyman–Pearson null hypothesis (where a decision is to be reached). Even knowing that statistical significance is not particularly related to biological significance, and that it is arbitrary, some investigators seem to feel comfortable b

ilarly, between underfitting and overfitting. The estimation information is a natural and simple way to view model selection. In a set of candidate models, select that fitted model where information is maximized. Proper model selection is reflected in good achieved coverage for the parameters in the model (or for predictions). Perhaps too much bias has been accepted in the tradeoff to get a false sense of high precision. This represents the worst of a highly precise, but quite biased estimate. These ideas have been in statistical thinking.

An information criterion (i.e., AIC, AIC_c , QAIC, and others) can rank the candidate models from best to worst and scale the model weights and evidence ratios. Often data do not support a single best for data analysis. Instead, suppose three models are for best, while another, larger, set of models is clearly not (underfit or overfit). Such virtual “ties” for the best approximation be carefully considered and admitted. Poskitt and Tremayne’s “portfolio of models” that deserve final consideration. Chapter 4 that there may be more than one model that is to be regarded as

Ambivalence

The inability to ferret out a single best model is not a deficiency of other selection criterion. Rather, it is an indication that the data are inadequate to reach such a strong inference. That is, the data are concerning some effect or parametrization or structure.

In such cases, all the models in the set can be used for inferences: multimodel inference.

It is perfectly reasonable that several models would serve in approximating the information in a set of data. Inference there are sometimes competing models and the data do not support only one. The issue of competing models is especially relevant to model selection uncertainty into estimators of precision. When a model has substantial support, some form of multimodel inference (averaging) should be considered (Chapter 4). The following provide some important details that must be considered in analyzing research data.

2.11.1 AIC Cannot Be Used to Compare Models on Different Data Sets

Models can be compared using the various information criteria relative, expected K-L information, only when they have been fitted to the same set of data. For example, if nonlinear regression models

AIC cannot be used to compare models where the data are in the ungrouped case (Model U) and grouped (e.g., grouped into histogram bins) case (Model G).

Data Must Be Fixed

An important issue, in general, is that the data and their distribution must be fixed and alternative models fitted to this fixed data. Information criteria should not be compared across different data sets because the inference is conditional on the data in hand.

2.11.2 Order Not Important in Computing AIC

The order in which the information criterion is computed or estimated is not relevant. Often, one may want to compute AIC_c , start with a complex model and proceed to simpler models with fewer parameters; to start with the simple models and work up to the more complex models with many parameters; this strategy might be best if numerical problems are countered in fitting some high-dimensional models. There is no issue here to proper interpretation, as opposed to the various hypotheses testing approaches where the order may be both arbitrary and the results sensitive to the choice of order (e.g., stepup (forward) vs. stepdown (backward)). Section 3.4.6 provides an example).

2.11.3 Transformations of the Response Variable

Model selection methods assume that some response variable is the subject of interest. Assuming that the scientific hypotheses concern a particular response variable, then all the models must represent exactly this variable. If the response variable is transformed, then the models must represent the transformed variable. If models in the set should all have the same response variable, then the information of mistake is illustrated by the following example. An investigator is interested in modeling a response variable y and has built 4 linear models for y , but during the model building, he decides to include a model for $\log(y)$. At that point he includes a model for $\log(y)$ as the fifth model. The information in such cases cannot be validly compared. The investigator's point, and often overlooked. In this example, one would find the best model followed by the other 4 models, each having large AIC_c values. On this result, one would erroneously conclude the importance of the nonlinearity. **Investigators should be sure that all hypotheses concern the same response variable (e.g., if the whole set of models is for $\log(y)$, no problem would be created; it is the mixing of models for y and $\log(y)$ that is incorrect).**

$$g_1(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{[y - \mu]^2}{\sigma^2}\right]$$

and another model,

$$g_2(y|\mu, \sigma) = \frac{1}{y\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \frac{[\log(y) - \mu]^2}{\sigma^2}\right]$$

Another critical matter here is that all the components of each model must be retained in comparing different probability distributions. The requirement of comparing pdfs in this spirit in Section 6.7.1. This requirement is not needed in cases like multiple regression analysis because all the comparisons are about the model structure (to select) with an assumption of normal errors for every model. There is a global model and its associated likelihood, and μ is used to represent μ as a regression function.

In other cases, it is tempting to drop constants in the log-likelihood if they do not involve the model parameters. However, alternative models may not have the same constants; this condition makes valid comparisons impossible. The simple solution here is to retain all the constants in the likelihood for all the models in the set.

2.11.4 Regression Models with Differing Error Structures

This issue is related to that in Section 2.11.3. A link between the residual sum of squares (RSS) and σ^2 from regression models with normal errors to the maximized log-likelihood value was provided in Section 2.11.3. This link is a special case, allowing one to work in an ordinary least squares regression framework for modeling and parameter estimation. It is possible to map to a likelihood framework to compute $\log(\mathcal{L}(\theta|data, model))$ for quantities under an information-theoretic paradigm.

The mapping from $\hat{\sigma}^2$ to $\log(\mathcal{L}(\theta|data, model))$ is valid only if all models in the set assume independent, normally distributed errors with a constant variance. If some subset of the models assume non-normal errors, then valid comparisons across all the models are not possible. In this case, all the models, including those with non-normal error structures, should be put into a likelihood framework since the maximum likelihood estimates of $\log(\mathcal{L}(\theta|data, model))$ and criteria such as AIC and BIC are

Tests of null hypotheses and information-theoretic approaches used together; they are very different analysis paradigms. A mistake seen in the applied literature is to use AIC to rank the best model and then “test” to see whether the best model (the alternative) is “significantly better” than the second-best model (the null). This procedure is flawed, and we strongly recommend against it (Burnham 2001c). Despite warnings about the misuse of hypothesis testing (Burnham et al. 2000, Cox and Reid 2000), researchers are still reporting on trivial null hypotheses, while failing to report effect size and confidence intervals.

Some authors state that the best model (say g_3) is *significant* relative to another model (say g_6) based on a Δ value of 4–7. Alternatively, one sees that model g_6 is *rejected* relative to the best model. These statements are poor and misleading. It seems best not to associate terms like *significant* or *rejected* with results under an information-theoretic paradigm. Instead, we concern the strength of evidence for the models in the set, using the evidence ratio (Section 2.10), as well as an adjusted R^2 , and other model diagnostics or descriptive statistics.

2.11.6 Null Hypothesis Testing Is Still Important in the Context of Strict Experiments

A priori hypothesis testing plays an important role when a randomized experiment (i.e., treatment and control groups being formally contrasted) design with random assignment) has been done and specific hypotheses have been identified. In these cases, there is a well-developed statistical theory on testing of treatment effects in such situations. We certainly acknowledge the value of traditional testing in the analysis of these *experimental* data. Still, the primary emphasis should be on the size of the treatment effects and their precision; too often the focus is on “significance,” while the treatment and control means are presented (Anderson et al. 2000 Cox and Reid 2000). No longer are we calling for estimates of effect size and associated precision, but rather statistics, P -values, and “significance.”

Akaike (1981) suggests that the “multiple comparison” problem in means should be viewed as a model selection problem, rather than as one of the many testing methods that have been developed (Burnham 1988). Here, a priori considerations would be brought to bear on a set of candidate models derived, letting information theory aid in sorting out differences in treatment means—a reformulation of the problem, instead of on testing. An alternative approach is to use effects modeling (Kreft and deLeeuw 1998).

model selection and associated parameter estimation. In the case where nuisance parameters are encountered in the recapture or resighting probabilities in capture–recapture studies. Here, it is not always clear what either the null hypothesis should be in a hypothesis testing framework. If hypotheses that are tested are naive or trivial, as Johnson et al. (2009) point out with such clarity. Should we expend resources to find out if all birds are white? Is there any reason to test formally hypotheses such as “the probability of robins is the same in cities A and B”? Of course not! We should assume that the number is different and proceed to estimate the difference and its precision: an estimation problem, not a hypothesis testing problem.

2.11.7 *Information-Theoretic Criteria Are Not a Test*

The theories underlying the information-theoretic approach to model selection and hypothesis testing are fundamentally quite different.

Criteria Are Not a Test

Information-theoretic criteria such as AIC, AIC_c , and BIC are not “tests” in any sense, and there are no associated concepts such as P -values or α -levels. Statistical hypothesis testing represents a more powerful and generally superior paradigm for the analysis of data than information-theoretic criteria.

It seems best to avoid use of the word “significant” in reporting research results under an information-theoretic paradigm.

The results of model selection under the two approaches are often very similar with simple problems; however, in more complex situations with many candidate models, the results of the two approaches can be quite different (Section 3.5). **It is critical to bear in mind that there is no inherent advantage to information-theoretic approaches to model selection over the use of null hypothesis testing for model selection.** The use of *ad hoc* (albeit a very refined set of *ad hoc* procedures in some cases) is not a test.

2.11.8 *Exploratory Data Analysis*

Hypothesis testing is commonly used in the early phases of exploratory data analysis to iteratively seek model structure and understand the data. One starts with 3–8 models, compute various test statistics for each, and find that several of the better models each have a gender effect. Thus, more models are generated to include a gender effect, and more null hypotheses are tested. Then the analyst notes that several of these models have a significant time for some set of estimable parameters; thus more models are generated, and so on. While this iterative or sequential approach is common, it is not a test.

we suggest that if the results are treated only as alternative hypotheses, a more confirmatory study to be conducted later, this might be a reasonable practice, particularly if other information is incorporated at a later stage. Still, the sequential and arbitrary nature of such testing makes us wonder whether this is really a good exploratory technique. It is readily key in on unique features of the sample data at hand. In any event, the key here is to conduct further investigation to weed out the “hunches” from the tentative exploratory work. Conclusions from an investigation has too often been ignored and the tentative ‘hunches’ published as if they were a priori results. Often, the authors do not mention the post hoc activities that led to the supposed results.

We suggest that information-theoretic approaches might be a useful exploratory tool; at least key assumptions upon which they are based are not terribly violated, and there is no arbitrary α level. The analysis using an information-theoretic criterion, instead of a significance test, eliminates inferential problems in interpreting the results, but one must still worry about overfitting and spurious correlations (see e.g. Akaike al. 2001b). The ranking of alternative models (the Δ_i and AIC_c) might be useful in the preliminary examination of data resulting from a pilot study. On these insights, one could design a more confirmatory study to address the issue of interest. The results of the pilot exploration should be used to guide the issue. While we do not condone the use of information theoretic approaches for data dredging, we suggest that it might be a more useful tool for testing in exploratory data analysis where little a priori knowledge is available. Data dredging has enough problems and risks without using an information theoretic approach that carries its own set of substantial problems.

2.12 Some History and Further Insights

Akaike (1973) considered AIC and its information theoretic interpretation as a natural extension of the classical maximum likelihood principle. Fisher (1936) anticipated such an advance over 60 years ago.

...an even wider type of inductive argument may be developed, which shall discuss methods of assigning functional form of the population.

This comment was quite insightful; of course, we might have said that A. Fisher! Akaike was perhaps kind to consider AIC an extension of the ML theory; he might just as well have said that classical likelihood is a special application of the more general information theoretic approach. Fisher believed in the importance of information theory as a tool for statistics.

that the negative of K-L information is Boltzmann's entropy; this information has been called negative entropy or "negative entropy," while max entropy is maximum disorder or maximum entropy. Conceptually,

$$\text{Boltzmann's entropy} = -\log\left(\frac{f(x)}{g(x)}\right)$$

Then,

$$-\text{Boltzmann's entropy} = \log\left(\frac{f(x)}{g(x)}\right)$$

and

$$\begin{aligned} \text{K-L} &= E_f(-\text{Boltzmann's entropy}) \\ &= E_f\left(\log\left(\frac{f(x)}{g(x)}\right)\right), \\ &= \int f(x) \log\left(\frac{f(x)}{g(x)}\right) dx. \end{aligned}$$

Thus, minimizing the K-L distance is equivalent to maximizing entropy; hence the name *maximum entropy principle* (see Jaynes 1983a, 1985 and Bozdogan 1987, Jessop 1995 for further discussion). However, maximizing entropy is subject to a constraint: the model must contain only information in the data. A good model contains the information in the data, leaving only "noise." It is the noise (entropy or uncertainty) that is maximized under the concept of the entropy maximization principle (Jaynes 1983a). Minimizing K-L information then results in an approximation to the minimum amount of information in the data. Entropy maximization is a model that maximizes the uncertainty, leaving only information in the data "maximally" justified by the data. The concepts are equivalent: maximizing K-L distance (or information loss) certainly seems the more appropriate.

The K-L information is *averaged* negative entropy, hence the name *averaged negative entropy* with respect to f . While the theory of entropy is a large subject, here readers can think of entropy as nearly synonymous with information, randomness or disorder in physical systems.

Boltzmann derived the fundamental theorem that

entropy is proportional to $-\log(\text{probability})$

Entropy, information, and probability are thus linked, although they are not necessarily multiplicative while information and entropies are additive (see Jaynes 1983a, also derived by Shannon 1948). Fritz Hasenöhl, Boltzmann's successor at Vienna University, and

mation concerning Borkzmann appears in Brush (1985, 1986). Further insights into Akaike's career are found in Findley and Pa

2.12.2 *A Heuristic Interpretation*

After Akaike's innovative derivation of AIC, people noticed an interpretation that was both interesting and sometimes misleading. The first term in AIC,

$$\text{AIC} = -2 \log(\mathcal{L}(\hat{\theta}|x)) + 2K,$$

is a measure of lack of model fit, while the second term is interpreted as a "penalty" for increasing the size of the model (or parsimony in the number of parameters). This heuristic does not do justice to the much deeper theoretical basis for AIC (in terms of distance and information theory). The heuristic interpretation has led statisticians to consider "alternative" penalty terms, and this has been unproductive (see Chapter 6). The so-called penalty term in AIC is, rather, it is the asymptotic bias-correction term. It is the asymptotic estimator of relative, expected K-L information. [If we had defined Akaike's AIC as $-\log(\mathcal{L}(\hat{\theta}|x)) + K$, the misinterpretation would have remained unchanged; some authors use this expression, but we will not. We have defined it.]

The heuristic view of the components of AIC clearly suggests a balance tradeoff and insight into how the principle of parsimony is implemented in AIC (see Gooijer et al. 1985:316). Still, we recommend using the estimate of the relative expected K-L information or distance between pairs (i.e., each g_i vs. f). Minimizing this relative, expected K-L information provides an estimated best approximating model for that pair (the *closest* approximating model to f). The relative K-L information is the element in AIC model selection.

2.12.3 *More on Interpreting Information-Theory*

Estimates of relative K-L information, the AIC difference between models, and the weights (w_i) provide a ranking of the models; thus the answer to the question which fitted model is best, which are essentially tied for best, which are clearly in an inferior class (and perhaps some that are in a superior class). These ranks are, of course, estimates based on the data. They are quite useful (cf. Section 2.7 and Sakamoto et al. 1986). The primary inference be developed using the model for which the AIC is lowest or the small number of models where there is an essential

model is > 9 – 10 AIC units from the minimum) there is 1 uncertainty and the theoretical standard errors can be used in Sections 1.2.3 and 2.14). When the results of model selection methods described in Chapter 4 can be considered. A of models and the identification of models that are nearly equal to those that are clearly poor explanations for the data at hand. Hypothesis testing provides no general way to rank models that are nested.

One must keep in mind that there is often considerable selection of a particular model as the “best” approximation served data are conceptualized as random variables; the different if another, independent set were available. It is “ability” that results in uncertain statistical inference from a set being analyzed. While we would like to make inferences robust to other (hypothetical) data sets, our ability to do so is limited, even with procedures such as AIC, with its cross-validation, and with independent and identically distributed samples. Computer-intensive resampling methods may well further assess the uncertainty of our inferences, but it remains to understand that proper model selection is accompanied by an uncertainty of uncertainty. The bootstrap technique can allow insight into this uncertainty; this and other similar issues are the subject of some of the chapters.

2.12.4 *Nonnested Models*

A substantial advantage in using information-theoretic criteria is that they are valid for nonnested models (e.g., Table 2.2). Of course, the ratio tests are defined only for nested models, and this is a substantial limitation in the use of hypothesis testing in model ranking of models using AIC helps clarify the importance of this (Burnham 1973:173); for example, some models for a particular data set and should not be used for inference.

A well-thought-out global model (where applicable) is required during the entire investigation, including the clear statement of the question to be answered, the collection of the data. This prior knowledge is then carefully used in the development of the set of candidate models (Section 1.2.4). In the ground science, the entire investigation should probably be very preliminary.

models, with some work being done in log-linear and c (e.g., factor analysis) models. Bozdogan (1987) provides a and some extensions. However, the number of published examine the performance of AIC-selected models is quite problem with the statistical literature as regards the evaluation the use of Monte Carlo methods using only very simple generating a few large effects and no smaller, tapering effects. Further Monte Carlo studies usually have a poor objective, namely, to evaluate criterion selects the simple generating model. We believe this point entirely with respect to real data analysis. Such evaluations even without regard for sample size (and often use AIC which has been used).

In Monte Carlo studies it would be useful to generate a more realistic model with several big effects and a series of small effects (Speed and Yu 1993). Then interest is refocused on finding a good approximating model and its statistical properties. The goal is to select the simple, artificial model used to generate the data. The goal to select a best approximating model for the data at hand is not the goal. The “true model” is at all complex, its use, with estimates that are less than true ones, would be poor for inference, even if it existed. If the form (but not parameter values) were known (e.g., Sakamoto 1993), a counterintuitive result occurs because the (limited) data would be used to estimate all the unknown parameters in the “true model,” which would result in a substantial loss of precision (see Figure 1.3B).

AIC reformulates the problem explicitly as a problem of finding the best of the true structure (probably infinite-dimensional, at least in the sciences) by a *model*. Model selection then becomes a simple minimization, where AIC (or more properly K-L information) is to be minimized. AIC selection is objective and represents a paradigm to that of null hypothesis testing and is free of the problems of levels, the multiple-testing problem, and the fact that some models might not be nested. The problem of what model to use is a hypothesis testing problem (Akaike 1974). However, the problem of a simple comparison of models does not justify the comparison of models (Akaike 1985 and Section 1.3.3). If one had 10 parameters, there would be 1,024 possible models, even if interactions are excluded. If sample size is $n \leq 1,000$, overfitting is a certainty. It is simply not sensible to consider such a large number of models because a model that overfits the data will almost surely result in the loss of the problem has been lost. *Even in a very exploratory analysis, it is a bad practice to consider all possible models; surely, some scientists would bear on such an unthinking approach* (otherwise, the science would be done and the work could be done by a technician).

The bootstrap is a type of Monte Carlo method used frequently in statistics. This computer-intensive approach is based on resampling from the data (Efron and Tibshirani 1993, Mooney and Duval 1993). It was first described by Bradley Efron (1979); thousands of papers have been written on the bootstrap, with various extensions and applications. Over the decades, and it has found very wide use in applied probability theory. It can be used for several purposes, particularly in the robust estimation of parameters, plugging variances or standard errors and (asymmetrical) confidence intervals. It has been used in the estimation of model selection free of bias, and estimates of precision that include model selection uncertainty.

The bootstrap has enormous potential for the biologist. It is a simple skill; however, its computer intensive nature will continue to be a challenge for large problems. We believe that at least 1,000 bootstrap samples are needed in many applications, and often 10,000 samples are needed for complex model selection. In extreme cases, reliable results could take a long time to apply the bootstrap to complex data analysis cases with large sample size and several dozen models, where the MLEs may not be found numerically.

The fundamental idea of the model-based sampling theory in statistical inference is that the data arise as a sample from some unknown probability distribution f . Uncertainties of our inference arise from the fact that if we can estimate f . The bootstrap method allows the construction of measures of our inference uncertainty by having a simple estimate \hat{f} of f and sampling from this estimated distribution. In practice, the empirical bootstrap means using some form of resampling from the actual data x to generate B (e.g., $B = 1,000$ or $10,000$) bootstrap samples; a bootstrap sample is denoted as x_b , where $b = 1, \dots, B$. Each sample data consist of n independent units, and it then suffices to draw a random sample of size n , *with replacement*, from the n units of the original bootstrap sample. However, the nature of the correct bootstrap inference can be more complex for more complex data structures.

The set of B bootstrap samples is a proxy for a set of B samples from f (in reality we have only one actual sample). The uncertainties expected from replicate real samples are inferred from the bootstrap by analyzing each bootstrap sample exactly as we first analyze the original sample. From the set of results of sample size B we measure the bootstrap uncertainties from sample to (conceptual) population (Freedman et al. 1991). In applications it has been theoretically shown (e.g., Efron and Tibshirani 1993) that the bootstrap can work well for large n (say $n > 10$), but it is not generally reliable for small n (say $n = 5$), regardless of how many bootstrap samples B are used. The bootstrap is not always successful in model selection (see Freedman et al. 1991).

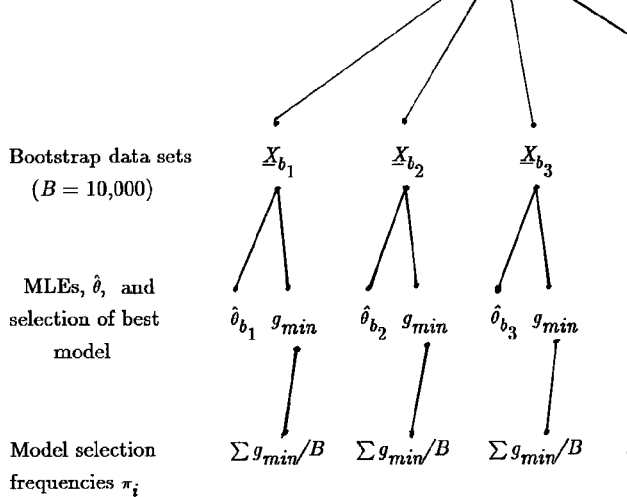


FIGURE 2.6. Diagram of the nonparametric bootstrap method as in Efron and Tibshirani (1993). The actual data set \underline{X} is sampled, using the same sample size (n); this is done B times, to obtain \underline{X}_b . Maximum likelihood theory provides estimates of the parameters for models i ($i = 1, 2, \dots, R$) and the AIC-best model (denoted by m). Its index is stored for each of the bootstrap data sets. Finally, the model selection frequencies (π_i) are computed as the sums of the frequencies where model i was best, divided by B . Of course, $\sum \pi_i = 1$.

2.13.1 Introduction

In many cases one can derive the sampling variance of an estimator from general likelihood theory. In other cases, an estimator may be complicated or may not exist in closed form. For example, the finite rate of increase (λ) can be derived from a Leslie population projection matrix. The variance of age-specific fecundity and age-specific, conditional survival probabilities. Generally, λ cannot be expressed in closed form. The bootstrap provides a variance estimation in such nonstandard cases.

Consider a sample of weights of 27 young rats ($n = 27$). The data are from Manly (1992),

57 60 52 49 56 46 51 63 49 57 59 54 56 59 57 52 52 61 59 53

The sample mean of these data is 54.7, and the standard deviation is $cv = 0.0824$. For illustration, we will estimate of the standard deviation of the sample mean. Clearly, this would be nonstandard; however, it represents a case where the bootstrap is useful.

perhaps it would not appear in the subsample, or it could appear once. Similarly, there are 3 occurrences of the weight 57 in the original data, so perhaps the bootstrap sample would have, by chance, none of them. The important point here is that a random sample of size 27 is taken without replacement from the original 27 data values. This is the first bootstrap resample. For this bootstrap sample, one computes $\hat{\mu} = \bar{x}$, the $\widehat{\text{se}}(\hat{\mu})$, and the $\text{cv} = \widehat{\text{se}}(\hat{\mu})/\hat{\mu}$, and stores that value of cv in memory.

Second, the whole process is repeated B times (where $B = 10,000$ samples for this example). Thus, we generate 10,000 bootstrap samples ($b = 1, 2, 3, \dots, 10,000$) and from *each* of these we compute the $\hat{\mu}$ and the cv and store the value of the cv.

Third, we obtain the estimated standard error of the cv (the standard deviation of the 10,000 bootstrap cv's, responding to the 10,000 bootstrap samples). The procedure for this is the same as in the case of the mean. In this case, the standard error of the cv is 0.00922, or less than 1% of the cv.

Confidence intervals can be computed in the usual way. The 95% confidence interval gives a 95% interval of (0.0640, 0.1009) for the rat data. However, the bootstrap distribution may be nonnormal and a more robust interval would be preferred. Again, the bootstrap provides a simple approach. In this case, the 10,000 $B = 10,000$ estimates of the cv in ascending order and sorted. The 2.5th percentile cut off the lower and upper 2.5 percentiles. Thus, the resulting confidence interval is asymmetric.

In the rat cv, the percentile bootstrap 95% confidence interval is (0.0640, 0.1009). This interval is about the same width as in the traditional normal theory shifted a bit toward 0. Incidentally, the mean of the 10,000 bootstrap cv's was 0.0806 (compared to the actual sample cv of 0.0824). The bootstrap provides usually adequate for the estimation of the sampling variance; however, good estimates of percentile confidence intervals require $B = 10,000$ in complicated applications.

Just as the analysis of a single data set can have many uses, the bootstrap can be used to provide insight into a host of questions. For each bootstrap sample one could compute and store the cv, the estimated covariance matrix, goodness-of-fit values, the estimate of the standard deviation factor, the model selected, confidence interval width, and so on. Inference can be made concerning these quantities, based on the B bootstrap samples.

The illustration of the bootstrap on the rat data is called nonparametric bootstrap, since no parametric distribution is assumed for the process that generated the data. We assume only that the data in the original sample were "representative" and that sample size was not small. The bootstrap is frequently used and allows assessment of bias and standard error. The use of the parametric bootstrap will be illustrated by the estimation of the inflation factor \hat{c} .

$\chi^2_{\text{gof}}/\text{df} = 3.2$. The investigators suspected some extrabinomial variation. They are surprised by the large estimate of the variance inflation factor (VIF) and suspect that the estimate is high and decide to use a parametric bootstrap to investigate their suspicion. They realize that the program RELEASE (Fleiss et al. 1987) can be used to do Monte Carlo simulations and to compute the goodness-of-fit statistics.

They input the MLEs from the real data into RELEASE and specify the parameters (ϕ_j and p_j) and use the numbers of new releases (n_j) as input. Then the amount of extrabinomial variation (i.e., the variance inflation factor, called EBV in RELEASE) is specified. In this illustration, the investigators specify meaning no overdispersion. They then run 1,000 Monte Carlo simulations to obtain the information on the estimated variance inflation factor. The average of these 1,000 values gives $\hat{E}(\hat{c})$, and this can be used as a value used to generate the data. This result provides insight into what to do about possible overdispersion in their data. The investigators could conduct several such studies for a range of values of \hat{c} to see whether $E(\hat{c}|\text{EBV}) = \text{EBV}$ and assess any systematic bias in the estimate of EBV.

This bootstrap is parametric in that parameters were specified (taken from the MLEs from real data that were available) and used to fit a model to produce Monte Carlo data. The nonparametric bootstrap does not require parameters nor a model and relies on resampling the data.

The bootstrap has been used in population biology to set confidence intervals on the median and mean life span. It is conceptually simple and has widespread use in applied statistics. Biologists planning a study of life span teaching should be familiar with the bootstrap. There is a good book on the bootstrap; see Efron and Tibshirani (1993) for an excellent introduction to the subject and a large list of references. Some valid applications of the bootstrap are tricky (even multiple linear regression), so some care is needed in complex settings!

2.13.2 The Bootstrap in Model Selection: The Bootstrap Frequency

Consider the case where data (x) with sample size n are available and several models are under consideration, each representing some hypothesis of interest. Let $B = 10,000$ bootstrap data sets, each of size n , be generated by resampling the data with replacement. MLEs of the parameters could be computed for each bootstrap sample. Then AIC_c could be computed for each of the 6 ($i = 1, 2, \dots, 6$) models and the number of times the i th model is the best (denote this by r^* , where r^* is the number of the best of the 6 models) associated AIC_c value stored for each of the 10,000 bootstrap samples. In 10,000 such analyses, one has the bootstrap frequency of the i th model being the best.

given by $\pi_i = \text{frequency}/10,000$ in this example. Of course, the AIC or TIC could have been used to estimate the π_i .

Relative frequencies for model i being selected as the best model are not the Akaike weights, but are not identical. There is no reason to expect data-based weights of evidence (as the set of w_i) to be the same as the relative frequencies at which the models are selected by an empirical procedure as being best. In general, likelihood provides a better measure of weight of evidence about parameter values, given a model (see Royall 1997), and we think that this concept (i.e., evidence about parameter values) is best represented by the likelihood of a model) rightly about a best model given an a priori set of models.

In our work we have not seen any particular advantage in using selection frequencies over the Akaike weights. Considering the large computer times required for the computation of the model selection frequencies, we prefer the Akaike weights in general. We present results in Chapters 4 and 5.

We further elaborate on the interpretation of the Akaike weights, which are conceptually different from the sampling-theory-based relative frequencies of model selection. It has been noted in the literature (see, e.g., 1994, Bozdogan 1987) that there is a Bayesian basis for interpreting the weight w_i as being the probability that model g_i is the expected best model given the data (for convenience we usually drop this “expected” from the phrase, just think of the K-L best model). Once we have accepted a model g_i given the data $\mathcal{L}(g_i|x)$, then we can compute the a posteriori probability that model g_i is the K-L best model if we assume uniform prior probabilities on the models (note that some Bayesian approaches to model selection take this approach ad hoc since it is not the full Bayesian approach). We must specify an a priori probability distribution τ_1, \dots, τ_R on the models, reflecting our belief that fitted model g_i will be the K-L best model for the data. These probabilities τ_i must be specified independent of fitting any models to the data.

2.14 Return to Flather’s Models

We now extend the example in Chapter 1 where 9 models were fitted to the accumulation curve for data from Indiana and Ohio were fitted (1992, 1996). The simple computation of AIC was done using the nonlinear regression output from program NLIN in SAS (SAS Institute Inc. 1996). In this case, apart from a constant that is the same over all models,

$$\text{AIC} = n \cdot \log(\hat{\sigma}^2) + 2K,$$

number of parameters (K), AIC values, $\Delta_i = \text{AIC}_i - \text{AIC}_{\min}$ value, adjusted R^2 values for the Indian–Ohio Major Land Resource Area for each model; the order is not relevant. Here the models are shown in order of the number of parameters (K). However, this is only a convenience example in Table 1.1.

Model	Number of parameters ^a	AIC value	Δ_i
ax^b	3	227.64	813.12
$a + b \log(x)$	3	91.56	677.04
$a\left(x/(b+x)\right)$	3	350.40	935.88
$a(1 - e^{-bx})$	3	529.17	1114.65
$a - bc^x$	4	223.53	809.01
$(a + bx)/(1 + cx)$	4	57.53	643.01
$a(1 - e^{-bx})^c$	4	-42.85	542.63
$a\left(1 - [1 + (x/c)^d]^{-b}\right)$	5	-422.08	163.40
$a[1 - e^{-(b(x-c)^d)}]$	5	-585.48	0

^a K is the number of parameters in the regression model plus 1 for σ^2 .

where $\hat{\sigma}^2 = \text{RSS} / n$ and K is the number of regression parameters (including σ^2). AIC values for the 9 models are given in Table 2.2. Model 9 is clearly the best approximating model for these data. Values of $\text{AIC}_{\min} = \text{AIC}_i + 585.48$ are also given and allow the results to be interpreted. Here, the second- and third-best models are model 6 (corresponding to Δ_i values of 163.40 and 542.63, respectively). Since these Δ values are very large, and the inference here is that model 9 is clearly the best of the candidate models considered for these data, this conclusion seems to be born out by Flather (1992), who selected this model based on a careful analysis of residuals for each model and Mallows' C_p . The remaining question is whether a still better model has been postulated with 6 or 7 parameters and increased structure. **criteria attempt only to select the best model from the set of models available; if a better model exists, but is not offered as a candidate, the information-theoretic approach cannot be expected to select the new model.**

Adjusted R^2 values are shown in Table 2.2, and while R^2 is not a measure of the proportion of the variation “explained,” it is useful in model selection (McQuarrie and Tsai 1998). In the case of the best 4 models all have an adjusted $R^2 \approx 0.99$, prompting the inference (erroneously) that all 4 models are an excellent fit to the data. The Δ_i values shows that models 6, 7 and 8 are incredibly close to model 9. The evidence ratio for the best model versus the

There are additional reasons why adjusted R^2 is poor in usefulness should be restricted to description.

2.15 Summary

Ideally, the investigator has a set of “multiple working hypotheses” and has thought hard about the background science of the issue. Background science of the matter, experience, and expertise are used to select a set of candidate models, representing each of these hypotheses. **Important philosophical issues that must receive increased attention** research problem should be carefully stated, followed by careful planning concerning the sampling or experimental design. Sample size and other issues should be considered fully before the data-gathering begins.

The basis for the information-theoretic approach to model inference is **Kullback–Leibler information**,

$$I(f, g) = \int f(x) \log \left(\frac{f(x)}{g(x|\theta)} \right) dx.$$

$I(f, g)$ is the “information” lost when the model g is used to approximate full reality or truth f . An equivalent interpretation of $I(f, g)$ is that from the approximating model g to full truth or reality. In this interpretation, we seek to find a candidate model that minimizes $I(f, g)$ among the candidate models. This is a conceptually simple, yet non-trivial task. However, $I(f, g)$ cannot be used directly, because it requires knowledge of full truth or reality and the parameters in the approximating model.

Akaike (1973), in a landmark paper, provided a way to estimate the expected $I(f, g)$, based on the empirical log-likelihood function. He expected the maximized log-likelihood value was a biased estimate of the true Kullback–Leibler information and that under certain conditions the bias was approximately equal to K , the number of estimable parameters in the approximating model g . His method, *Akaike’s information criterion*, provided a model selection to be firmly based on a fundamental theoretical door to further theoretical work. He considered AIC to be a natural extension of likelihood theory, the very backbone of statistical theory. Takeuchi (1976) derived an asymptotically unbiased estimator of the expected Kullback–Leibler information that applies in general. Under special conditions underlying Akaike’s derivation of AIC (the derivation for Takeuchi’s information criterion) requires large sample sizes. The calculation of elements of two $K \times K$ matrices in the bias-adjustment formula represents an important conceptual advance and further justifies AIC.

approaches to adjusting for this bias (the bias adjustment from the maximized log-likelihood) are summarized below.

Criterion	Bias adjustment term
AIC	K
AIC_c	$K + \frac{K(K+1)}{n-K-1}$
TIC	$\text{tr}(J(\theta)I(\theta)^{-1}) \approx K$.

These information criteria are estimates of relative, expected log-likelihood and are an extension of Fisher’s likelihood theory. AIC is easy to compute, quite effective in many applications, and widely used. When count data are found to be overdispersed, appropriate information criteria have been derived, based on quasi-likelihood (QAIC_c). If overdispersion is found in the analysis of count data, the log-likelihood function must be divided by an estimate of the overdispersion parameter (\hat{c}) to obtain the correct log-likelihood. Thus, investigators using count data analysis have several powerful methods for selecting the best model. Making inferences from empirical data to the population of interest is difficult. In practice, one need not assume that the “true model” is in the set of models (although this is sometimes mistakenly stated in the technical literature).

The AIC differences (Δ_i) and Akaike weights (w_i) are used to rank and scaling the hypotheses, represented by models. The Akaike weights (e.g., w_i/w_j) help sharpen the evidence for or against the hypotheses. All of these values are easy to compute and simple to interpret.

The principle of parsimony provides a philosophical basis for model selection. In addition, K-L information provides an objective target based on the log-likelihood. AIC, AIC_c, QAIC_c, and TIC provide estimators of relative model quality. Objective model selection is rigorously based on the log-likelihood. These methods are applicable across a very wide range of models, including biological and statistical models. We recommend presentation of the appropriate information criterion (AIC, AIC_c, QAIC_c or TIC) for each of the various models in research papers to provide full information on the evidence for each of the models.

Basic Use of the Information-Approach

3.1 Introduction

Model building and data analysis in the biological sciences pose that the investigator has some advanced education in the sciences, and statistics in particular. This requirement implies that the researcher has substantial knowledge of statistical null hypothesis approaches. Such investigators, including ourselves over the years, often find it difficult to understand the information-theoretic approach because it is conceptually so very different from the testing approach that is familiar. Relatively speaking, the concepts and practical use of the information-theoretic approach are simpler than those of statistical hypothesis testing, but much simpler than some of the Bayesian approaches to model building (Laud and Ibrahim 1995 and Carlin and Chib 1995).

The prevailing philosophy has been to use some test or criterion to select a model, from a set of models, that is somehow "best" in a particular sense. Inference is then entirely conditional on this selection. We believe that approach should be merely the beginning, a humble beginning at that. There is much more to the model building process than this initial solution. Substantive information is contained in the differences (Δ_i), since they are free from arbitrary (and unknown) constants and are directly interpretable in many cases. Both the Δ_i and the $\mathcal{L}(g_i|x)$ allow scientific hypotheses, carefully represented by model i , the discrete likelihood of model i , given the data ($\mathcal{L}(g_i|x)$), and a way to assess the relative support for the alternative models.

dence for the second, third, and *with* best model, respectively in the set. These methods go well beyond just the selection and are very useful in assessing the empirical evidence for applied scientific problems.

It will be made clear in the next two chapters that information-theoretic analysis and inference philosophies are only a midway point in the information-theoretic paradigm. As we have struggled to understand the paradigm, it has become clear to us that inference based on only a single best model is relatively poor for a wide variety of substantive reasons. Instead, we favor multimodel inference: procedures to allow formal inferences from all the models in the set. These procedures are simple to interpret and are the subjects of Chapters 4 and 5. Such multimodel inference includes model averaging, incorporating model selection uncertainty, estimates of precision, confidence sets on models, and simple measures of the relative importance of variables.

The examples below focus on the selection of a single best model, which will appear in the following chapters. However, many more examples go beyond this initial approach in terms of the evidence for the best model in the set. Methods to assess model selection uncertainty (e.g., model selection uncertainty and Akaike weights w_i) are illustrated and discussed. These methods provide the relative likelihood of model i , given the data, provide a quantitative basis for inferences concerning inferences about the actual K-L best model. Such methods should be viewed as a halfway point in understanding the information-theoretic approach where formal inferences are drawn from the data.

While the derivation of AIC (Chapter 7) lies deep in the foundations of mathematical statistics, its application is quite simple. Our initial example is a multiple linear regression model of cement hardening and drying time in the model selection literature. The remaining examples focus on more complex data sets and models. These examples illustrate the ease and power of AIC in model selection and inference. Several of these examples are used in later chapters as additional concepts and methods and several examples deal with survival models, since that has been a major research interest.

Given a model, likelihood inference provides a quantitative measure of the strength of evidence in the data regarding the plausible values of the parameters in the model (Royall 1997). Given a well-developed set of candidate models, information-theoretic methods provide a quantitative measure of the strength of evidence in the data regarding the plausible values of the parameters in the model. Information criteria can be computed and interpreted in the aid of subjective judgment (e.g., α -levels or Bayesian probabilities). The candidate models has been derived.

approach was the clarification of the importance of model assumptions. For substantial, prior information on the system being studied,

At some early point in the analysis of count data, the global model should be assessed using standard methods. This assessment should accompany continuous data (see Carroll and Rupert 1988). Generally, no concept of overdispersion in continuous data. Residual variation should receive careful attention. One should look for outliers, highly leveraged points, symmetry, trends, and autocorrelation in residuals (McCullagh and Nelder 1989). There are many diagnostic procedures that should be used to aid in the modeling of count data. If, after proper attention to the a priori considerations, the model fits poorly, then information-theoretic methods will select a small set of poor-fitting models. This undesirable situation promotes poor science that went into the modeling and definition of count data models. Lack of fit of the global model should be a flag that more consideration must be given to the modeling, based on the nature of the questions being asked and the design of the data collection. The effort must be classed as exploratory and very tentative. Some data dredging, leading perhaps to some tentative model conclusions. Treated as the results of a pilot study, then more data collected and the analysis could proceed in a more confirmatory manner using the techniques we outline in this book.

Computer programs for likelihood methods nearly always calculate the log-likelihood at its maximum, and the appropriate information can be easily computed by hand, if necessary. Similarly, the MLE of σ^2 from standard output of LS programs allows calculation of information criteria from this estimate in most cases. While some packages currently print AIC, relatively few print the value of the log-likelihood and this is unfortunate (see Example 1 below, where AIC is used because the ratio n/K is small).

3.2 Example 1: Cement Hardening Data

The first example is a small set of data on variables thought to influence the heat evolved during the hardening of Portland cement (Wolfe 1964:649). These data represent a simple use of multiple linear regression (see Section 1.2.2). This data set (the “Hald data”) has been used by many authors (e.g., Hald 1952:635–649, Seber 1977, Daniel and Wood 1981:294–342 and 629–673, Stone and Brook 1986, McCulloch 1993, Hjorth 1994:31–33, Ronchetti and Staud 1996, Ibrahim 1996, and Sommer and Huggins 1996) and will be used here.

($3\text{CaO} \cdot \text{SiO}_2$), x_3 = tetracalcium aluminato ferrite ($4\text{CaO} \cdot \text{Al}_2\text{O}_3 \cdot \text{Fe}_2\text{O}_3$), silicate ($2\text{CaO} \cdot \text{SiO}_2$)] are used to predict the dependent variable evolved per gram of cement after 180 days of hardening.

x_1	x_2	x_3	x_4	y
7	26	6	60	78.5
1	29	15	52	74.3
11	56	8	20	104.3
11	31	8	47	87.6
7	52	6	33	95.9
11	55	9	22	109.2
3	71	17	6	102.7
1	31	22	44	72.5
2	54	18	22	93.1
21	47	4	26	115.9
1	40	23	34	83.8
11	66	9	12	113.3
10	68	8	12	109.4

important points. The data include 4 predictor variables and a response variable of 13 (Table 3.1). The predictor variables (as a percentage of cement) are x_1 = calcium aluminate ($3\text{CaO} \cdot \text{Al}_2\text{O}_3$), x_2 = tricalcium silicate ($3\text{CaO} \cdot \text{SiO}_2$), x_3 = tetracalcium aluminato ferrite ($4\text{CaO} \cdot \text{Al}_2\text{O}_3 \cdot \text{Fe}_2\text{O}_3$), x_4 = dicalcium silicate ($2\text{CaO} \cdot \text{SiO}_2$), while the response variable is y = evolved per gram of cement after 180 days of hardening. (1971) provide further details on these data for the interest of the reader. The primary focus of this chapter is “approximating model to use?”

The small size of the sample necessitates the use of cross-validation; however, we will present comparable values for AIC in the next chapter. We use an obvious notation for denoting what variables are in each model. That is, if variables x_1 and x_3 are in a particular model, we denote it as model {13}; each model has an intercept (β_0).

3.2.1 Set of Candidate Models

Because only 4 variables are available, the temptation is to consider all possible models ($2^4 - 1 = 15$) involving at least one of the predictor variables. However, due to the small sample size, we will consider this example as large. Lacking any personal knowledge concerning the physics of cement hardening, we will consider the full set of models, including the null model {0} with $K = 5$ parameters. While we generally advocate consideration of all possible models of the x_i (but no interaction terms), of the predictor variables), this approach will allow some

values for both AIC and AIC_c followed by the Akaike weights (w_i), are ordered in terms of Δ_i for AIC_c .

Model	K	$\hat{\sigma}^2$	$\log(\mathcal{L})$	Δ_i AIC	Δ_i AIC _c
{12} ¹	4	4.45	−9.704	0.4346	0.0000
{124}	5	3.69	−8.478	0.0000	3.13
{123}	5	3.70	−8.504	0.0352	3.17
{14}	4	5.75	−11.370	3.7665	3.33
{134}	5	3.91	−8.863	0.7528	3.88
{234}	5	5.68	−11.290	5.6072	8.74
{1234}	6	3.68	−8.469	1.9647	10.53
{34}	4	13.52	−16.927	14.8811	14.44
{23}	4	31.96	−22.519	26.0652	25.63
{4}	3	67.99	−27.426	33.8785	31.11
{2}	3	69.72	−27.586	34.2052	31.43
{24}	4	66.84	−27.315	35.6568	35.22
{1}	3	97.37	−29.760	38.5471	35.77
{13}	4	94.39	−29.558	40.1435	39.70
{3}	3	149.18	−32.533	44.0939	41.32

¹ Here, $\log(\mathcal{L}) = -n/2 \cdot \log(\hat{\sigma}^2) = -9.7039$, $AIC_{min} = -2 \log(\mathcal{L}) + 2K = AIC + \frac{2K(K+1)}{n-K-1} = 32.4078$.

others in the published literature (e.g., Draper and Smith 1966 and Hoeting and Ibrahim 1996). We note, however, that not only a single variable might have been excluded on a priori basis, but the cement involves a mixture of at least two compounds that must be present. We will extend this example in Chapter 4 to examine the issue of model uncertainty and other issues.

3.2.2 Some Results and Comparisons

The use of AIC_c suggests model {12} as the best approximation to these data (Table 3.2). The estimated regression coefficients for this model are

$$\hat{E}(y) = 52.6 + 1.468(x_1) + 0.662(x_2)$$

where the estimated standard errors of the 3 estimated parameters (the *model*) are 2.286, 0.121, and 0.046, respectively (this result is from Hald 1952). The adjusted $R^2 = 0.974$ and the ML estimate of σ^2 for the AIC_c -selected model. The second-best model is {124}, which is 3.13 units from the best model (Table 3.2). Other candidate models are also clearly many of the models represent poor approximations to the data (at least the models in Table 3.2 with Δ_i values > 10). No

global model (model {1234}).

The Akaike weight for the best model is not large, relative to the other models. The ratio of the weights for the best 4 next-best models ranges from only 4.8 to 7; this is not surprising. Model {12} is likely best if other replicate samples were available.

Using a type of cross-validation criterion (Q_{cv}), Hjorth (1981) selected model {124} with $K = 5$ for these data. Here, his result is

$$\hat{E}(y) = 71.6 + 1.452(x_1) + 0.416(x_2) - 0.2$$

where the estimated standard errors are 14.142, 0.117, 0.117, respectively. Model {124} has an adjusted $R^2 = 0.976$ and Draper and Smith (1981:325–327) used cross-validation and the Mallows selection criterion, which is quite similar to Q_{cv} , and also selected model {124}. Note, had AIC been used, ignoring the ratio $n/K \approx 2$, model {124} would have been selected (Table 3.2); AIC_c should be used if this ratio is small.

Is there any basis to say that AIC_c selected a better approximation to the truth than Hjorth's cross-validation procedure or AIC or the Mallows procedure? It is difficult to answer conclusively because truth is not known. The regression coefficient on x_4 is not "significant" under a hypothesis testing scenario ($t = 1.36$, 9 df), and the estimated standard error on the regression coefficient for x_2 increased by a factor of 1.186 compared to model {12}. The adjusted R^2 statistics for model {124} is 0.976 (vs. 0.974), but it has one additional parameter. The regression coefficient between x_1 and x_3 was -0.824 , while the coefficient on x_4 was -0.973 . Just on the basis of this latter correlation, it would be wise to allow both x_2 and x_4 in the same model (if n were only 13, perhaps there would be more support for including x_4). While not completely compelling, it would seem that AIC_c selected a better parsimonious model in this case. An additional, negative argument is the computer-intensive nature of Hjorth's cross-validation procedure compared to the information-theoretic approach. With more variables or sizes or more variables, or with more models to consider, these approaches may often become computationally too "costly."

Draper and Smith (1981) used Mallows's C_p statistic and selected model {12}, in agreement with AIC_c (this might be fortuitous, because a simple version of C_p or Q_{cv} is available). They further point out that the constant (approximately 98%) for any i ; thus the $X'X$ matrix is theoretically singular. Small rounding errors were evident, since the percentage data were expressed as integers, leading to a barely nonsingular matrix. At best, model {1234} would be a poor approximation to the truth. They also warn against the unthinking use of multiple regressions and present a detailed analysis of forward, backward, and stepwise regressions.

and retaining only x_1 and x_2 . This represents an improvement over model {12}, but it does not merely add new variables, without looking to see whether x_3 has become redundant. Draper and Smith (1981) provide a review of the various older model selection alternatives and offer recommendations (but do not discuss any of the information-theoretic approaches). They provide an intensive analysis of the cement data of Box and Jenkins (1976) and include detailed computer output in two large appendices.

Another analysis approach involves computation of the eigenvalues and eigenvectors of the elements on the (centered) $X'X$ matrix and examination of the eigenvalues for the 4 explanatory variables (see Draper and Smith 1981 and Brooks 1990). The principal component eigenvalues are 1.57607, 0.18661, and 0.00162. Approximately 95.3% of the variance is contained in the first 2 eigenvectors, while 99.96% is contained in the first 3 eigenvectors. These results certainly suggest that the global model is overfitted to the data (i.e., 4 predictor variables are redundant). In addition, the results suggest that 2 predictor variables will suffice (given $n = 13$). Of course, the use of the percentage eigenvalues requires some judgment and, furthermore, relatively few biologists are knowledgeable about principal component eigenvalues and eigenvectors. *We believe that the investigators do not understand the methods leading to the results of their work; this is a common problem with some advanced methods.* Such understanding seems to be a prerequisite for the information-theoretic approaches.

One could ask whether there is a need for model selection when there are only 4 predictor variables (i.e., why not merely take the global model, estimate the parameters and use it for inference?). This simple strategy is illustrated here. First, note that this global model is overfitted to the data relative to model {12}, and is therefore a poor approximation to the true model. The estimates of parameters for the global model are

$$\hat{E}(\hat{y}) = 62.4 + 1.551(x_1) + 0.510(x_2) + 0.102(x_3)$$

where the estimated standard errors, given this model, are 0.755, 0.709, and 0.709, respectively. These standard errors are large because the matrix is nearly singular (the percentage coefficients of variation for the parameters were 4.3, 8.2, and 6.9 under model {12}, compared to 11.2, 11.2, and 11.2, respectively, under model {1234} (see Wood and Thomas 1986)). The regression coefficient for x_1 might be judged as “significant” in a statistical testing sense, and the model is clearly overfit (see Figure 1). The global model has an adjusted $R^2 = 0.974$ and $\hat{\sigma} = 1.918$. Surely a parsimonious model such as {12}, would better serve the analyst in this case.

Loss of precision is expected in using an overfit global model. In addition, there is also a nonnegligible probability that even the simple model parameter may be incorrect in such cases. It seems some

If all the predictor variables are mutually orthogonal (which is not be so bad. Orthogonality arises in controlled experiments and levels are *designed* to be orthogonal. In observational data, there is a high probability that some of the predictor variables will be dependent. Rigorous experimental methods were just being developed at the time these data were taken (about 1930). Had such data been widely available and the importance of replication understood, it might have been possible to break the unwanted correlations and establish cause and effect.

With only a single data set, one could use AIC_c and select the best model for inference. However, if several other independent data sets were available, would the same model be selected? The answer is *perhaps*. In general, there would be variation in the selected model from one data set to the next, as there would be variation in parameter estimates over different data sets. If the same model is used for analysis. The fact that other data sets are available, the use of other models leads us to the issue of model selection.

Based on simulation studies, we are usually surprised by the results. There is a lot of variation in selecting a parsimonious model for a given problem. It is probably the case that in many real-world problems there is a lot of model selection uncertainty. We generated 10,000 bootstrap samples to estimate model selection uncertainty. The parameters of the models shown in Table 3.2, were estimated and AIC_c was calculated for each bootstrap sample. The following summary shows the relative frequencies (π_i) from applying AIC_c (models not shown) to each of the 10,000 bootstrap samples. Here, $\hat{\pi}$ are the estimated selection probabilities. Also shown are the Akaike weights (w_i) for each data set:

Model	K	Bootstrap Sel. Freq. $\hat{\pi}_i$	Akaike weights w_i
{12}	4	0.5338	0.567
{124}	5	0.0124	0.118
{123}	5	0.1120	0.116
{14}	4	0.2140	0.107
{134}	5	0.0136	0.081
{234}	5	0.0766	0.007
{1234}	6	0.0337	0.003
{34}	4	0.0039	0.000.

As might be expected with such a small sample size, the results varied substantially, and model {12} was selected as the best model in 53.38% of the bootstrapped samples. Model {14} was selected 21.40% of the time. It is a quite reasonable result that models {12} and {14} are the most frequently selected.

3.2.3 A Summary

In summary, the simple approach of using AIC_c appears to be a parsimonious model as the basis for inference from the data. AIC_c sharpens the inference about which parsimonious model is best to AIC. A priori information could have resulted in fewer models and generally strengthened the process (note, that Hald (1956) did only an analysis of x_1 and x_2 and presented the analysis of the other variables several pages later). It seems likely that model {14} as a variable might have been excluded from serious consideration. It must have been known about cement in the late 1920s. Since it is that Woods et al. (1932) had knowledge of the negative relationship between x_2 and x_4 ; after all, model {14} was their second-best model. Both of both x_2 and x_4 in the same model (where the correlation was the over parametrized global model. An important feature of the theoretic approach is that it provides a ranking of alternative models. Some inferences to be made about other models that might be made. In addition, the rankings suggest that some models that remain. Models {24}, {1}, {13}, and {3} for the cement data). The AIC_c is to focus the evidence for or against the various models.

carefully defining a small set of candidate models, based on what is known about the problem, cannot be overemphasized.

An investigator with, say, 10 explanatory variables can learn much from the data and a multiple linear regression analysis. Some substantial supporting science that can be used to help select a set of models to consider. In this case, there would be $2^{10} = 1024$ more if transformations or interaction terms were allowed. It would surely be a risk. The analysis, by whatever method chosen, is considered exploratory and the results used to design further experiments leading to a more confirmatory analysis, based on some a priori information.

3.3 Example 2: Time Distribution of an Insecticide Added to a Simulated Ecosystem

This example concerns the addition of the insecticide DDT to a laboratory system that simulates a pond of water. The original study by Smith (1966) and his colleagues; our main reference for this example is Smith (1966).

Diaz and Neely note (1975) that the determination of the distribution of this chemical introduced into an ecosystem is an environmental issue. They go on to mention that "... a model describing each step of the process would be extremely important, however, to try to find a suitable model to predict the long-term environmental consequences." This is consistent with Akaike's and the one recommended here is used because it rests on systems of first-order differential equations. Parameters, given a model, are estimated by least squares. These can easily be used to compute AIC values to aid in selecting the best approximating model.

The active ingredient of DURSBAN® is 0,0-Diethyl (4-methylpyridyl) phosphorothioate, which was labeled with radioactivity on the pyridyl ring and added at a level of 1 mg/6 gal in a 10-gal aquarium (Figure 3.1). This aquarium contained 2 inches of soil (100 g/liter), plants (salvinia, anacharis, milfoil, and water cucumber). Samples of the various components were analyzed for radioactivity at different time periods following the addition of DURSBAN. Each time period yielded a sample size (n) of 36. The results are in percentages from the crude radioactivity measurements.

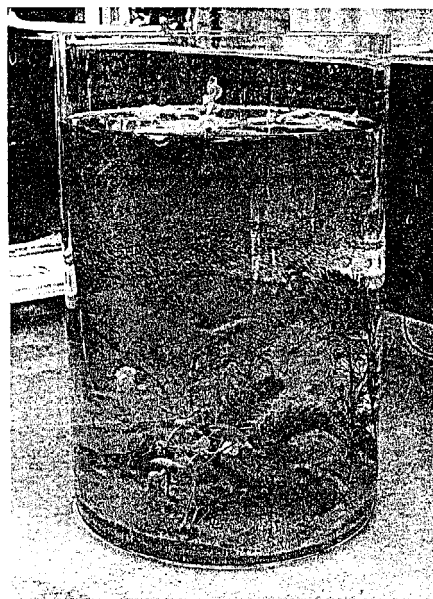


FIGURE 3.1. Glass aquarium used in the studies of DURSABAN.

Time after DURSBAN® addition (hours)	Percent radioactivity	
	Fish	Soil & Plants
0	0	0
1.5	15.2	35.2
3.0	19.0	46.0
4.0	19.3	56.0
6.0	20.7	61.0
8.0	23.0	60.5
10.0	24.2	59.3
24.0	21.2	51.5
48.0	23.0	38.3
72.0	22.7	38.3
96.0	20.5	36.3
120.0	17.3	38.3

1975:150). The authors of the study assumed that the measurements were normally distributed, with zero means and a constant standard deviation (we take this to mean the actual measurement error of the measurements).

3.3.1 Set of Candidate Models

Blau and Neely (1975) had a great deal of knowledge about the system they exploited this in a priori model building. They began by assuming that an equilibrium exists between DURSBAN® in the water and the soil components (B), and a direct uptake of the chemical by the fish. This led to their Model 1 (Figure 3.2), which was represented by a set of three differential equations, where the rate parameters to be estimated are k_1 , k_2 , and k_3 :

$$\begin{aligned} dx_A(t)/dt &= -k_1x_A(t) + k_2x_B(t) - k_3x_C(t), \\ dx_B(t)/dt &= k_1x_A(t) - k_2x_B(t), \\ dx_C(t)/dt &= k_3x_A(t), \end{aligned}$$

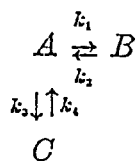
with initial conditions $x_A(0) = 100$, $x_B(0) = 0$, and $x_C(0) = 0$. This is a three-compartment model (Brown and Rothery 1993) and is similar to the model used in the field studies. Blau and Neely (1975) used $x_A(t)$, $x_B(t)$, and $x_C(t)$ to denote the amount at time (t) of A , B , and C , respectively, with the restriction that

$$x_A(t) + x_B(t) + x_C(t) = 100.$$

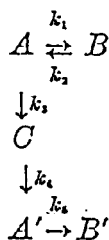
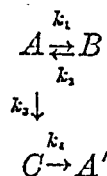
They used nonlinear least squares to estimate model parameters and their analytic methods were quite sophisticated. The parameter estimates for this model were $\hat{k}_1 = 0.510$, $\hat{k}_2 = 0.800$, $\hat{k}_3 = 0.00930$.

$k_2 \downarrow$

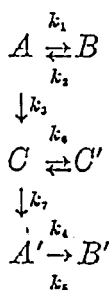
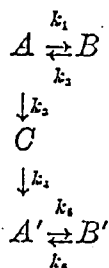
C



Model 2a



Model 3a



Model 4a

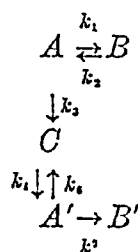


FIGURE 3.2. Summary of models used by Blau and Neely (1975) in the DURSCHAN[®] in a simulated pond ecosystem.

(this is their residual sum of squares (RSS) divided by $n - K$ (σ^2); thus, $K = 4$ for this model.

Blau and Neely (1975) built six other models, each based on the system, but also based on examination of the residuals (there are some inconsistencies here that we were unable to resolve; we will use the material from their paper). While some data did support their main derivation of additional models seemed to support their hypotheses about the processes. They were well aware

of fit tests and separated “pure error” from the remaining tests. Model selection was accomplished by statistical hypothesis tests (see Section 1.2.2) and examining the RSS. They found Model 4a (see Table 3.4) to be the best and also found some support for Model 4b.

3.3.2 Some Results

Analysis of these data under an information-theoretic paradigm was done by Blau and Neely’s (1975) Table II, since they provide values for each of their seven models. Due to the relationships between the models and ML theory (see Section 1.2.2),

$$\log(\mathcal{L}(\hat{k}, \hat{\sigma}^2 \mid data)) = -n/2 \cdot \log(\hat{\sigma}^2)$$

where $\hat{\sigma}^2 = \text{RSS} / n$. Then,

$$\text{AIC} = -2 \cdot \log(\mathcal{L}(\hat{k}, \hat{\sigma}^2 \mid data)) + 2K$$

and

$$\text{AIC}_c = \text{AIC} + \frac{2K(K+1)}{n-K-1}.$$

These computations were done by hand on a simple calculator in approximately 20 minutes. The results of this extended analysis are shown in Table 3.4 and suggest that Model 4a is the best to use for inference (in line with Blau and Neely (1975)). Only Model 4b is a competitive model with a value of 7.611 and seems relatively implausible for these data. The evidence ratio for model 4a vs. 4b is 0.978/0.022, indicating strong support for 2-way transfer between the viscera and the essential difference between models 4 and 4b is $c \nrightarrow d$.

Carpenter (1990) used these data and seven models to compare Bayesian analysis with equal Bayesian prior probabilities for each model.

TABLE 3.4. Summary of model selection statistics (the first three columns are from Blau and Neely 1975). Statistics for the AIC_c -selected model are shown in bold.

Model	RSS	$\log(\mathcal{L}(\hat{k}, \hat{\sigma}^2 \mid data))$	K	AIC	AIC_c
1	5374	−90.105	4	188.209	189.4
2a	1964	−71.986	5	153.972	155.9
2b	848	−56.869	5	123.737	125.7
3a	208.3	−31.598	6	75.196	78.0
3b	207.9	−31.563	7	77.127	81.1
4a	58.6	−8.770	8	33.540	38.8
4b	79.4	−14.238	7	42.475	46.4

Model 1 is a poor second. In this example, a ranked list of the various methods might be expected to be in somewhat of a tie. This example illustrates that it is often easy to perform model selection on complex systems, based on information provided in prior analysis clearly shows that five of the seven models have similar support, and inferences from these models would likely be poor. The third-best model (3a) has an evidence ratio of 3×10^9 , while the best model (1) has an evidence ratio of 2×10^{33} . Clearly, these models are given the data available.

Blau and Neely's (1975) results are interesting, and would likely be the best model. The evidence ratio for the second-best model is also reasonable to base inference on just the best model in the set. One is often comfortable with the concept that inferences can be made from a model; in a sense, the inference here is the model.

After a final model is chosen it is often wise to examine the model using standard methods. Such examination may reveal issues with the model; in this sense, science never "stops."

Formal *statistical* inferences include the following: (1) the liberation between DURSBN and the soil and plant system; (2) the uptake by a shorter uptake of DURSBN by the fish; (3) fish tend to excrete DURSBN; (4) the liberated material (metabolite) is again taken up by the soil and plants; (5) fish have two compartments: the flesh and the flesh; (6) the final sink for DURSBN is the soil and (7) plants readily dissipate the metabolite as degraded H_2O . Estimates of the various transfer rates are given by the model. Estimates of precision are available as standard errors or confidence intervals.

3.4 Example 3: Nestling Starlings

We generated a set of Monte Carlo data to illustrate model selection. We discussed with a much more complicated example of an example of a model. Thus, in a sense, the generating model is "truth"; we will use a model that is realistic for the moment, but mitigate it by including many parameters and a wide variety of tapering treatment effects. In addition, we will use a global model that has four fewer parameters than the generating model. The generating model is not in the set of candidate models. The example contains many so-called nuisance parameters (parameters that are not of interest). This is the only example in Chapter 3 where "truth" is not known. Interesting insights can be gained from this knowledge. The model is what parsimonious approximating model can be used for inference. It will lead to valid inference about the structure of the system.

3.4.1 Experimental Scenario

We generated data to mimic the experiment conducted (1988) (also see Burnham et al. 1987:343–348). The research to the survival effects of an organophosphate pesticide administered to European starlings (*Sturnus vulgaris*). We assume for illustration that the field experiment is designed using artificial nest boxes placed on a small island. Fledgling birds are assumed not to leave the island and are banded in the early fall months when the experiment is conducted (1988). Nest boxes are monitored during the nesting season to determine when hatching. All nestlings are leg-banded with uniquely numbered bands following hatching, and half of those nestlings are randomly assigned to the treatment group and the remaining birds assigned to a control group. We will assume that 600 nestling starlings are banded and released from each box (i.e., the number of starlings originally released in each box). All nest boxes contain 4 young birds (thus 2 treatment and 2 control). We assume these to be of nearly uniform size and age and that they move about and behave independently. Starlings randomly assigned to be in the treatment group receive an oral dose of pesticide. Starlings in the control groups are given pure corn oil under otherwise identical conditions. Colored leg bands provide a unique identification for each bird and therefore its group membership, on each weekly resighting. Band collection will be assumed to begin after a 4-day period following hatching; for simplicity, we assume that no birds die due to handling or band marking but before resighting efforts begin a week later. Starlings are potentially resighted during the following 9 weeks; sampling occurs on the island and is done on each Friday for 9 weeks. Thus, there are 10 occasions; occasion 1 is the initial marking and release of the birds, and 9 resighting occasions.

The pesticide is hypothesized to affect conditional survival probabilities (parameters of interest) and resighting probabilities (the nuisance parameters). However, the pesticide industry's position is that only mortality effects are likely, while environmental groups suspect that there may be both (short-term) and chronic (long-term) effects on survival probabilities. We assume that the resighting probabilities might also be affected by the treatment. The set of candidate models might span the range of the conditional survival probabilities; of course, one might design the experiment to include several control groups released at different, independent locations (islands), and use these as the basis for empirical estimates of treatment effect and prediction. See Burnham et al. 1987 for a discussion of experiments of this general type. We focus on an example of the model selection issue and not

tional survival probability (ϕ) and resighting probability (p) for treatment (t) and control (c) groups at week i :

$$\begin{aligned}\phi_{ti} &= \phi_{ci} - (0.1)(0.9)^{i-1} \text{ for } i = 1, \dots, \\ p_{ti} &= p_{ci} - (0.1)(0.8)^{i-2} \text{ for } i = 2, \dots,\end{aligned}$$

using program RELEASE (Burnham et al. 1987). These models assume a smooth temporal tapering of effect size due to the treatment on survival and resighting probabilities. That is, each week the effect of pesticide is diminished. We used the initial per-week survival and resighting probabilities for the control group as 0.9 and 0.8, respectively. The initial survival and resighting probabilities for the control group were $\phi_{ci} \equiv \phi_c \equiv 0.9$ and $p_{ci} \equiv p_c \equiv 0.8$. The data were collected from week 3.5 for each treatment and control group.

3.4.3 Set of Candidate Models

Define ϕ_{vi} as the conditional probability of survival for group v ($v = t$ for treatment and c for control) from week i to $i + 1$, and p_{vi} as the conditional probability of resighting for treatment group t ($i = 2$ to 10). The set of models that seem reasonable might include: no treatment effects (g_0), a model for an acute effect on survival probability ($g_{1\phi}$), and a model for an acute effect on both survival probability and the first resighting probability (denote that the first occurs at week 2) (model g_{2p}). This initial line of a priori models led to three models:

Model	Parametrization
g_0	All $\phi_{ti} = \phi_{ci}$ and all $p_{ti} = p_{ci}$ (no treatment effects)
$g_{1\phi}$	g_0 , except $\phi_{t1} \neq \phi_{c1}$ (an acute effect on ϕ_{t1})
g_{2p}	$g_{1\phi}$, except $p_{t2} \neq p_{c2}$ (acute effects on ϕ_{t1} and p_{t2})

Chronic effects might arise from starlings that are in the control group. Effects of the pesticide; these starlings might be more susceptible to the pesticide (this would be revealed in lessened survival during the summer period). Chronic effects, if they exist, might be revealed in the control group. That is, one might expect chronic effects to diminish the survival of the starlings in the control group. Agreement is reached, that chronic effects, if they exist, should not last more than one week.

10). The data given for each group (c) are the number of starlings released after last being released at time i . R_i = the number of birds released all of those released in weeks 2, . . . , 9 were merely rereleased. Each $(R(i) - \sum_j m_{ij})$ is modeled as a multinomial distribution with sa

Week	$R(i)$	Observed Recaptures for Treatment					
		$m(i, j)$					
		$j = 2$	3	4	5	6	7
1	300	158	43	15	5	0	
2	158		82	23	7	1	
3	125			69	17	6	
4	107				76	8	
5	105					67	2
6	82						5
7	81						
8	70						

Week	$R(i)$	Observed Recaptures for Control					
		$m(i, j)$					
		$j = 2$	3	4	5	6	7
1	300	210	38	5	1	0	
2	210		157	20	8	2	
3	195			138	24	2	
4	163				112	24	
5	145					111	1
6	139						10
7	124						
8	115						

Define $S_i = \phi_{ii}/\phi_{ci}$ for $i = 1$ to 7 as the measure of conditional survival probability, compared to the control. (The control group will experience some mortality as the birds age, but here the interest is in any *additional* mortality incurred by the pesticide treatment.) The parameters S_i ($i = 1, 2, \dots, 7$) are 0.929, 0.943, 0.954, 0.964, and 0.971, respectively. With no chronic effects, one expects $S_2 < S_3 < S_4 < \dots < S_7 < 1$, as compared to the parameters above (of course, the unconstrained *estimates* based on some approximating model, might not follow this pattern). Here, it seems reasonable to consider the presence of chronic effects, in addition to the hypothesized acute effects. Thus, the model includes chronic effects on both conditional survival and resighting

$g_{2\phi}$	g_{2p} , except $\phi_{t2} \neq \phi_{c2}$ (chronic effect on
g_{3p}	$g_{2\phi}$, except $p_{t3} \neq p_{c3}$ (chronic effect on
$g_{3\phi}$	g_{3p} , except $\phi_{t3} \neq \phi_{c3}$ (more chronic ef
g_{4p}	$g_{3\phi}$, except $p_{t4} \neq p_{c4}$ (more chronic ef
\vdots	
$g_{7\phi}$	All ϕ_{vi} and p_{vi} differ by treatment group 7 weeks

This last candidate model ($g_{7\phi}$) allows chronic treatment effects on conditional survival and resighting probabilities up through the seventh week, in addition to the acute treatment effects on ϕ_{t1} and p_{t2} . This model is as complex as our global model, and it has 30 parameters. The treatment effects taper through the ninth week; thus, the generating model is not in this set of models and has more parameters than the global model (g_0).

Model g_0 has 17 parameters, while model $g_{7\phi}$ has 30 parameters. The simplest model would have a constant survival and resighting probabilities for each group ($g_{\phi,p}$) and thus no treatment or week effects on either survival or resighting probabilities. This model would have 17 parameters (ϕ and p). Alternatively, a four-parameter model could allow the four parameters to differ by treatment group (ϕ_t , ϕ_c , p_t , and p_c). Because of the relatively large sample size in this example, these models might be considered simple and unlikely to be useful based on initial biological information. We might well exclude these from the set of candidate models. **biological support should not be included in the set of candidate models.** However, as an example, we will include these simple models in the set of candidate models and note that these models might well be viewed as simple models if the initial sample size released were 60 instead of 600.

The effective sample size in these product multinomial models is the number of starlings released (or rereleased) at each week. [The effective sample size in these product multinomial models is a complicated issue that we will not divert attention to this matter here, except to say that here we are using the term in the context of AIC_c . Technical notes on this subject may be found in KPBS.] In this example, $n = 2,583$ releases (a resighting counts as a release, including “recaptured and rereleased”). Because 600 starlings were released at week 1 (the nest boxes), the remaining 1,983 starlings were resighted at least once. Because of the large effective sample size, AIC_c is unnecessary; however, if one chose always to use AIC , no problems would be encountered because AIC_c approaches AIC as n/K gets large.

A statistician on the research team suggests adding models that allow for the possible tapering treatment effects on conditional survival and resighting probabilities.

insights concerning possible long-term enroute treatment. Employing a type of sine transformation on the parameters α and β is used here. In this transformation, the parameter (θ) , representing the probability of survival, is assumed to be between 0 and 1) to be modeled as a function of the covariate (e.g., X) is replaced by the expression $\sin(\alpha + \beta X)$. The new parameters α and β are the intercept and slope parameters in the covariate model. The transformation utilizes one-to-one increasing or decreasing sigmoid functions and the logistic link function in generalized linear models. In particular, the submodel $g_{\sin p_t}$ were defined for the dynamics of starling group:

$$\begin{aligned} g_{\sin \phi_t} & \quad \sin(\phi) = \alpha + \beta(\text{week}), \\ g_{\sin p_t} & \quad \sin(p) = \alpha' + \beta'(\text{week}). \end{aligned}$$

These submodels each have only 2 parameters (intercepts α and β and α' and β') and assume that $\sin(\phi_t)$ or $\sin(p_t)$ is a linear function of the conditional survival of starlings in the treatment group with respect to the summer period progresses, eventually approximating the control group).

These above two submodels for the treatment group and the four submodels below for the control group:

$$\begin{aligned} g_{\phi_{ci}} & \quad \phi \text{ is allowed to differ for each week; hence } \phi_i \text{ for } (i = 1, \dots, 8). \\ g_{\phi_c} & \quad \phi \text{ is assumed constant across weeks.} \\ g_{p_{ci}} & \quad p \text{ is allowed to differ for each week; hence } p_i \text{ for } (i = 2, \dots, 9). \\ g_{p_c} & \quad p \text{ is assumed constant across weeks.} \end{aligned}$$

For example, a model can be developed using $g_{\sin \phi_t}$ for the treatment group and model g_{ϕ_c} for the conditional survival of the control group. This part of the model has 3 parameters; α , β , ϕ_c , for the resighting probabilities. Thus, one could consider the treatment group and model $g_{p_{ci}}$ for the control group as a function of the resighting probabilities. This would add the parameters $p_{c1}, p_{c2}, \dots, p_{c10}$, for a total of $K = 14$ parameters. As an illustration, a rich mixture of candidate models in Table 3.6 (a set of 24 models). If this were a real situation, still other a priori models might be carefully supported with biological reason. If this experiment only 60 nestlings, then several simple models should be included. The high-dimensional models would be deleted. This set of 24 models will serve as a first example where there is some substantial

a pesticide. (All values are scaled by the additive constant $-17,437.1$ in the best model.) Akaike weights (w_i) are also shown.

Model	AIC	No. Parameters	Δ_i
$g_{7\phi}$ (global)	4,495.409	30	27.1
g_{7p}	4,493.619	29	25.1
$g_{6\phi}$	4,491.649	28	23.1
g_{6p}	4,489.889	27	22.1
$g_{5\phi}$	4,491.679	26	23.1
g_{5p}	4,491.929	25	24.1
$g_{4\phi}$	4,490.199	24	22.1
g_{4p}	4,489.029	23	21.1
$g_{3\phi}$	4,489.629	22	21.1
g_{3p}	4,492.619	21	24.1
$g_{2\phi}$	4,501.809	20	34.1
g_{2p}	4,517.019	19	49.1
$g_{1\phi}$	4,523.489	18	55.1
g_0	4,532.599	17	64.1
$g_{\sin \phi_t, \phi_{ci}, \sin p_t, p_{ci}}$	4,485.669	21	17.1
$g_{\sin \phi_t, \phi_{ci}, \sin p_t, p_c}$	4,475.249	14	7.4
$g_{\sin \phi_t, \phi_c, \sin p_t, p_{ci}}$	4,479.359	14	11.1
$g_{\sin \phi_t, \phi_c, \sin p_t, p_c}$	4,467.779	6	0.0
$g_{\sin \phi_t, \phi_{ci}, p_{ti}, p_{ci}}$	4,488.629	28	20.1
$g_{\sin \phi_t, \phi_{ci}, p_{ti}, p_c}$	4,478.209	21	10.1
$g_{\sin \phi_t, \phi_c, p_t, p_{ci}}$	4,484.699	13	16.1
$g_{\sin \phi_t, \phi_c, p_t, p_c}$	4,473.119	5	05.1
$g_{\phi_t, \phi_c, p_t, p_c}$	4,770.479	4	302.1
$g_{\phi, p}$	5,126.609	2	356.1

3.4.4 Data Analysis Results

As one would expect with simulated data, they fit the generating model well ($\chi^2 = 35.5$, 36 df, $P = 0.49$). [A large literature of goodness-of-fit testing in this class of models exists (e.g., Burnham and Anderson 1998; Pollock et al. 1990); we will not pursue the details of such tests here.] Since the data were simulated such that no overdispersion was present, the value of the overdispersion factor c could be computed under the null hypothesis from the results of the goodness-of-fit test, $\hat{c} = \chi^2/\text{df} = 0.99$. The global model $g_{7\phi}$ has fewer parameters than the generating model, so these data well ($\chi^2 = 35.4$, 30 df, $P = 0.23$). The value of \hat{c} for this model was 1.18, reflecting no overdispersion in this case. (The generating model, which is known to be true in this instance); after all, it, by definition, is of “truth.” In practice, one cannot usually distinguish between

particular, one might consider using the modifications to AIC_c (i.e., $QAIC = -2 \log(\mathcal{L})/1.18 + 2K$). We will mention this at a later point. The critical information needed for selection, model ranking and scaling the other models is shown in Table 3.6.

The interpretation of the 24 models for the experimental data (Table 3.6) can be sharpened by examining the Akaike weights. The AIC-selected model ($g_{\sin \phi_t, \phi_c, \sin p_t, p_c}$) is 0.906, while the second-best model ($g_{\sin \phi_t, \phi_c, p_t, p_c}$) has a weight of 0.063 and the third-best model ($g_{\sin \phi_t, \phi_c, p_t}$) has a weight of 0.022. The sum of the weights for the top three models is less than 0.01. In this case, one is left with strong support for the first model, with fairly limited support for the second-best model (evidence ratio of best vs. second-best ≥ 14). The evidence ratio for the third-best model is about 41.2. Thus, the data support one model as convincingly as the others. It seems to be little need to attempt model averaging or bootstrapping (see Table 4) to gain further robustness in inferences from these models. In addition, the use of conditional standard errors for the AIC-selected model, will likely suffice. Note that bootstrapping in this case is very, very difficult. Software development would be a very, very difficult task and computer time on a Pentium 1PC would likely take a very long time. The Akaike weights provide a distinct advantage in completing this simulated starling experiment.

The model with the minimum AIC value was $g_{\sin \phi_t, \phi_c}$ with parameters $(\alpha, \beta, \alpha', \beta', \phi_c, \text{ and } p_c)$. Using estimates of the MLEs of the survival and resighting parameters, one can derive MLEs of the survival and resighting probabilities. The MLEs for the treatment survival probabilities were as follows:

i	ϕ_{ti}	$\hat{\phi}_{ti}$	$\widehat{se}(\hat{\phi}_{ti})$
1	0.800	0.796	0.021
2	0.810	0.810	0.016
3	0.819	0.824	0.014
4	0.827	0.838	0.014
5	0.834	0.851	0.160
6	0.841	0.864	0.019
7	0.847	0.876	0.022
8	0.852	0.887	0.026
9	0.857	0.898	0.029

The survival parameter for the control group was 0.90, and the AIC-selected model was 0.893 ($\widehat{se} = 0.008$). These estimates are very close to the parameter values, and one can correctly infer the direction of the effect of the treatment on weekly survival probabilities. The 9 estimates of survival probability for the treatment group are very close to the MLEs of α and β in the submodel $\sin(\phi_t)$.

the conditional survival and resighting probabilities for both groups. The conditional survival and resighting probabilities for the control group were constant over weeks in this model, but not in the treatment group. The AIC-selected model captures the generating process. Figure 3.3 illustrates the similarity between the estimates from the global model ($g_{7\phi}$), and the AIC-selected model in terms of the treatment effect, $1 - S_i$.

Part of the reason that this analysis was successful was the reasoning that led to modeling the treatment effects, rather than trying to estimate the week-specific treatment effects (i.e., the S_i). Modeling allowed substantial insight into the tapering, which was not the case. Note: The two simplest models ($g_{\phi_i, \phi_c, p_i, p_c}$ with $K = 2$) were not at all plausible ($\Delta_i = 302.70$ and 356.13), so that these models would not normally have been considered for this experiment, since they lacked any reasonable biological support for the sample size involved. Of course, had sample size been very large, these models might have been more reasonable to include in the model selection process.

If sample size is small, one must realize that relatively little information is probably contained in the data (unless the effect size is large), and the data may provide few insights of much interest. It is easy to routinely err by building models that are far too complex for the data.

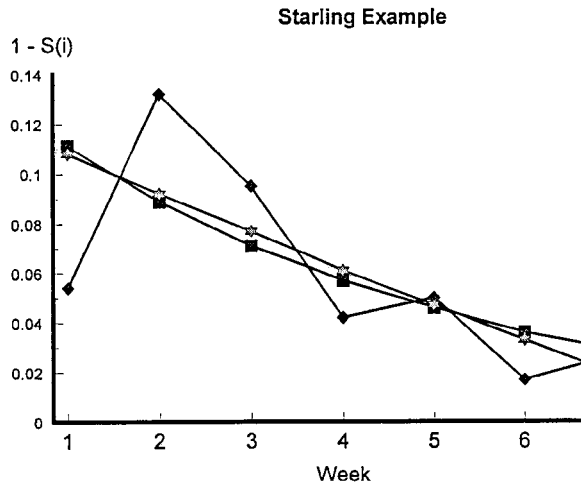


FIGURE 3.3. Treatment effect ($1 - S_i$, for week $i = 1, \dots, 7$) for the generating models ($g_{9\phi}$) with 34 parameters, compared with estimates from the global model ($g_{7\phi}$) with 30 parameters and the AIC-selected model with 30 parameters.

before analysts get a feeling for modeling based on some known about the science of the problem of interest.

3.4.5 Further Insights into the First Fourteen Ne

If only the first 14 models (Table 3.6) had been defined a concerning which model to use would have been far less cl these 14 models is over 25 units from the AIC-selected m not have been known. Second, 7 models have AIC values best of the 14. Thus, some additional steps would be nec model selection uncertainty into inference for these exp analysis was based on just the first 14 models.

We now examine further the results that would have be set of candidate models included just the first 14 models in tial theory (e.g., the estimators exist in closed form) and RELEASE, Burnham et al. 1987) exist for this sequenc allowing the illustration of a number of deeper points. F that these 14 models are clearly inferior to the models hy treatment effects (a diminishing linear treatment effect link function) for birds in the treatment group (e.g., the b model g_{4p} , is 21.25 AIC units above the selected model an compared to only 6 parameters for the AIC-selected mode to the importance of a good set of candidate models. Se chronic effects could not be identified by model g_{4p} (i.e., t effects on survival in the later time periods, S_4 , S_5 , and values provide clues that at the very least, models $g_{4\phi}$ (th (therefore, S_5) are also somewhat supported by the data models have AIC values within 1.17 and 0.86, respective fact, models g_{3p} through model $g_{6\phi}$ have fairly similar A and Figure 3.4). Unless the data uniquely support a particu not take the resulting model as *the* answer for the issue a that the particular data set can provide. Perhaps more tha be considered for inference from the 14 models (Chapter

The program RELEASE (Burnham et al. 1987) allows ap values of estimators and theoretical standard errors to be models in this class (i.e., the 14 appearing at the top of Tab allow insight into why the more minor chronic effects w model g_{4p} (the model estimated to be the best among the

i	$1 - E(\hat{S}_i)$	$\widehat{se}(1 - \hat{S}_i)$	$(1 - E(\hat{S}_i)) / \widehat{se}(1$
4	0.057	0.053	1.08
5	0.046	0.055	0.84
6	0.036	0.057	0.63

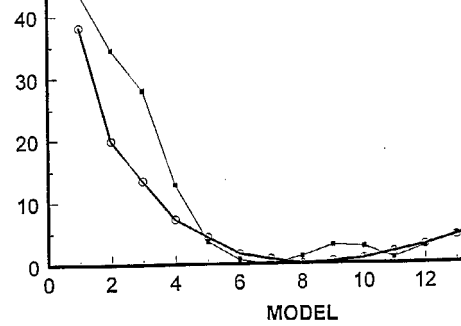


FIGURE 3.4. Estimated theoretical (heavy line) and sample Δ_i values for the 14 models used for the starling experiment. The estimated ($n = 50$) theoretical, expected AIC values (shown as open circles) are minimum for g_{5p} , while the realized AIC value from the sample data is minimum for g_{4p} (Table 3.6). Generally, there is good agreement between the theoretical and sample values (here plotted as Δ_i values).

The expected treatment effect size (i.e., $1 - E(\hat{S})$) was 0.5 for the standard errors were of a similar magnitude or larger than the standard errors of the mean for the first two columns above. The larger effects (i.e., S_1 and S_2) were not identified; however, at some point, the effect size is too small to identify with confidence from the information contained in the data. Still, if one had only the first 14 models and had used g_{4p} , inference from the data in this example would have been reasonable, but hardly optimal. The acute and larger chronic effects were convincingly identified. Comparison of AIC values for models g_{4p} and g_{6p} ($K = 28$) would have provided reasonable evidence for extended chronic treatment effects. Still, having to estimate the standard errors would lead to imprecise estimators, compared to those using the full model ($g_{\sin \phi_t, \phi_c, \sin p_t, p_c}$). AIC, AIC_c , and $QAIC_c$ are fundamental tools for a basis for a unified approach to the statistical analysis of data in the biological sciences. Further details concerning this class of models are provided by Anderson et al. (1994), Burnham et al. (1994), Burnham et al. (1998), and Anderson et al. (1998).

3.4.6 Hypothesis Testing and Information-Theoretic Approaches Have Different Selection Frequencies

At this point it is illustrative to examine briefly how information-theoretic selection compares to traditional approaches based on hypothesis testing. Thus, Monte Carlo methods were employed to generate

under model $g_{9\phi}$ with 34 parameters, which was parametrized to reflect the effect on both conditional survival and resighting probabilities for

Model	Hypothesis Testing			Information
	Stepup	Stepdown	Stepwise	
1 g_0	0.6	0.0	0.0	0.0
2 $g_{1\phi}$	17.4	0.0	13.9	0.0
3 g_{2p}	14.8	0.3	13.5	0.2
4 $g_{2\phi}$	26.8	1.5	26.3	1.2
5 g_{3p}	16.3	2.8	16.9	2.7
6 $g_{3\phi}$	14.6	6.9	16.1	6.8
7 g_{4p}	5.9	7.5	7.1	8.5
8 $g_{4\phi}$	2.7	11.9	3.8	13.5
9 g_{5p}	0.8	10.3	1.3	12.0
10 $g_{5\phi}$	0.2	13.3	0.7	14.3
11 g_{6p}	0.0	10.9	0.2	11.3
12 $g_{6\phi}$	0.0	12.9	0.2	11.3
13 g_{7p}	0.0	10.2	0.1	8.8
14 $g_{7\phi}$	0.0	11.5	0.1	9.5

pendent samples (data sets) using the same methods as with the original set of simulated data on nestling starlings. The numbers released and all parameter values were identical to generate the first set of data. Six methods were used to infer the best model: The first 3 methods involve well-known selection hypothesis testing (stepup or forward selection, stepdown selection, and stepwise selection), each using $\alpha = 0.05$. Three information methods were also used on each of the data sets: AIC, AICc, and \hat{c} as a variance inflation factor, estimated for each simulation. The results (Table 3.7 and Figure 3.5) show substantial differences in selection frequencies for the various methods.

Stepup selection, on average, selects model $g_{2\phi}$ with 26.8 parameters (on average was 20.3 parameters selected). These results are similar to the stepwise approach, which selects, on average, model g_{3p} with 16.9 parameters (the average number of parameters was 20.6). Given that the true model is under model $g_{9\phi}$ with 34 parameters, these methods seem to select simple models that miss most of the chronic treatment effects. Stepdown selection resulted, on average, in model g_{6p} (mean 26.9 parameters). These are quite different model selection frequencies than the other methods. Of course, the selection frequencies would change if a different (arbitrary) α level (say, 0.15 or 0.01) had been used. If the treatment impacts differed (i.e., a different model used to

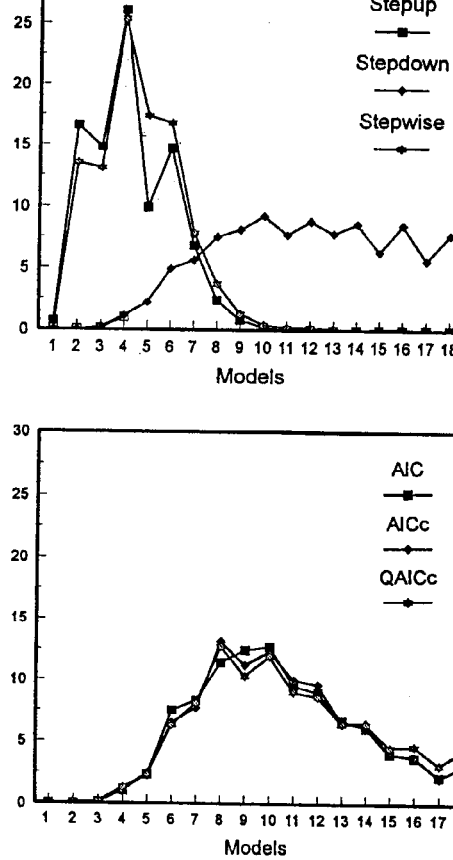


FIGURE 3.5. Model selection probabilities for three hypothesis testing methods and three information-theoretic approaches based on 50,000 Monte Carlo generated data, generated under model $g_{9\phi}$.

The practical utility of hypothesis testing procedures is limited in model identification (Akaike 1981b:722).

In this example, AIC selection averaged 25.4 parameters (model g_{5p}). Both AIC and the stepdown testing did reasonably well in identifying the larger chronic effects. The differences between AIC and stepdown testing are trivial, as one would expect from the large sample size (Table 3.7 and Figure 3.5). Even the use of QAICc results in a small difference in this example because, the estimated variance is near 1 ($\hat{c} = 1.18$).

in this simulated example. In addition, one must note that uncertainty exists in model choice for all six approaches (Figure 3.4.6). The material in Chapter 4 is particularly important, since the uncertainty should be incorporated into estimates of precise estimators. In this example the statistical hypothesis testing is an alternative to selection based on estimating the relative model quality. In general, we recommend strongly against the use of null model selection.

3.4.7 Further Insights Following Final Model Selection

Selection of the best model and the relative ranking of all models is objective, given a set of candidate models, and can be in the aid of subjective judgment. The formal data-based selection is a key part of the analysis. In the example, model $g_{9\phi}$ was selected. Thus “truth” is known and serves here as a basis for comparison. The model selection process does not try to select the model that generated these data. The goal is to select the model which model is the “best approximating model” for analysis. In the sense of having the smallest K-L distance from approximating the truth. Further information concerning the statistical properties of the models in the $g_0, \dots, g_{9\phi}$ class are given in Anderson et al. (2001).

The startling example illustrates an ideal a priori strategy to explore some potential realities *after* the analysis has been completed. We select model $g_{\sin \phi_t, \phi_c, \sin p_t, p_c}$ as the best (with 6 parameters). We must also perhaps consider models

$$\begin{aligned} g_{\sin \phi_t, \phi_c, p_t, p_c} & (\Delta_i = 5.34 \text{ with only 5 parameters}) \\ g_{\sin \phi_t, \phi_{ci}, \sin p_t, p_c} & (\Delta_i = 7.47, \text{ with 14 parameters}) \end{aligned}$$

before making some final inferences (at this time, the model selection process does not address the variance component due to model uncertainty). The model $g_{\sin \phi_t, \phi_{ci}, \sin p_t, p_c}$ is somewhat inconsistent. Why would the weekly variation in conditional survival for the birds in the treatment group be different from the only smooth time trends in the treatment group? Perhaps this model should lead to a thorough review of field methods in an effort to determine if there is a treatment effect. Perhaps this model should not have been in the set of candidate models. It may be picking up random variation in the data. The model $g_{\sin \phi_t, \phi_c, p_t, p_c}$ has Δ_i is 7.47, and this model has 8 additional parameters over the best model. It would seem that this model is a relatively poor approximation of the truth, though it might play a role in estimating the variance component due to model selection uncertainty.

After the analysis of the data to this point, suppose that the researcher asks about models where there is no treatment effect.

be added to the set of candidate models and more AIC values question brings up several points. First, if this suggestion wining the estimates of $p_{t2}, p_{t3}, p_{t4}, \dots, p_{p7}$ vs. p_{c2}, p_{c3}, p_{c4} that there seemed to be little difference between success pairs, then this is a form of data dredging, and any subsequent clearly detail the process by which the additional models v**encourage full investigation of the data to gain all possible information only want investigators to reveal the extent of any data to be placed.** Second, if that suggestion was made on conceptual grounds by studying the intermediate results, then the new class of models to the list, AIC computed, the Δ_i , w_i , and evidence ratio computed and inferences made. However, in this second case, the treatment what faulted for not considering the set of models more fully in place.

3.4.8 Why Not Always Use the Global Model for Inference?

Some might argue that the global model (or another model with parameters representing likely effects) should always be used for inference. After all, this model has been carefully defined by existing considerations and has all the effects thought to be reasonable. The principle of parsimony and endure the various issues concerning model selection, selection uncertainty, etc.? We can illustrate the problem using the starling data and the global model ($g_{7\phi}$) with 30 parameters. The results are given here in detail for estimates of S_i for $i = 1, \dots, 7$.

i	\hat{S}_i	$\widehat{\text{se}}(\hat{S}_i)$	CI_L	CI_U
1	0.946	0.044	0.858	1.033
2	0.868	0.052	0.766	0.969
3	0.905	0.058	0.792	1.017
4	0.958	0.057	0.847	1.069
5	0.950	0.047	0.859	1.042
6	0.983	0.051	0.883	1.084
7	0.973	0.059	0.858	1.088

The poor precision is illustrated by the upper confidence interval for 6 of the 7 include the value of 1 (i.e., no treatment effect). The limits would be truncated at 1.0). The average coefficient of variation for $\hat{\phi}_{ti}$ under model $g_{7\phi}$ is 4.75% vs. 2.08% under the AIC-selected model.

Attempts to select a properly parsimonious model for inference, based on rewards, primarily an approximating model that has a reasonable tradeoff between bias and variance. The tradeoff between bias and variance of model selection where expected K-L information loss is minimized.

of model parameters that are unnecessarily imprecise (Section 3.2.2), and this weakens the inferences made. In fact, to show the real patterns that can be validly inferred from the smooth decrease in ϕ_{ii} and thus in S_i . Sometimes there have 50, 100, or even 200 parameters, and this makes interpretation difficult. One cannot see patterns and structure, since there are so many parameters, most estimated with poor precision. Thus, some analysts have used analyses of these estimated parameters in order to “see the forest for the trees.” This has rarely been done correctly, since the estimators usually have high partial sampling correlations, making simple analysis results difficult to interpret. It is better to embed the reduced model in the log-likelihood function and use information-theoretic criteria to select a simple, *interpretable* model that fits the information in the data.

3.5 Example 4: Sage Grouse Survival

3.5.1 Introduction

Data from sage grouse (*Centrocercus urophasianus*) banded in Colorado, provide insights into hypothesis testing and model selection criteria in data analysis. The example is taken from Zabludowski (1985). Additional details are found there. Here we will use data on subadult (birds less than 1 year old) and adult (birds more than 1 year old) males banded during the breeding season (first week of March through the first week of May), from 1973 through 1987. Sage grouse are hunted by hunters, and nearly all of the band recoveries were from hunters who reported a banded bird and reported it to the Colorado Division of Wildlife. In 1973, 1,777 subadult and 1,847 adult males were banded. In 1987, the numbers of band recoveries were 312 and 270, respectively. This example is basic theory for modeling and estimation for these types of data, as found in Brownie et al. (1985).

Two types of parameters are relevant here: S_i is the survival probability relating to the annual period between banding i and $i+1$. r_i is the conditional probability of a band from a bird being recovered given that the bird died in year i . In the model building it is common to use i as a subscript to denote age (subadult vs. adult) and t to denote time (e.g., S_{a*t} denotes survival probabilities that vary by both age and time). In some models assume that r_i is a constant, thus resulting in a single parameter S_{15}).

Year Banded	Number banded	Recoveries by hunting season											
		73	74	75	76	77	78	79	80	81	82	83	84
1973	80	6	4	6	1	0	1	0	0	0	0	0	0
1974	54		6	5	2	1	0	0	0	0	0	0	0
1975	138			18	6	6	2	0	1	0	0	0	0
1976	120				17	5	6	2	1	1	0	0	0
1977	183					20	9	6	2	1	1	0	0
1978	106						14	4	3	1	0	0	0
1979	111							13	4	0	1	0	0
1980	127								13	5	3	0	0
1981	110									13	5	0	0
1982	110										7	0	0
1983	152											1	0
1984	102												0
1985	163												0
1986	104												0
1987	117												0

		73	74	75	76	77	78	79	80	81	82	83	84
		7	4	1	0	1	0	0	0	0	0	0	0
1973	99												
1974	38		8	5	1	0	0	0	0	0	0	0	0
1975	153			10	4	2	0	1	1	0	0	0	0
1976	114				16	3	2	0	0	0	0	0	0
1977	123					12	3	2	3	0	0	0	0
1978	98						10	9	3	0	0	0	0
1979	146							14	9	3	3	0	0
1980	173								9	5	2	0	0
1991	190									16	5	0	0
1982	190										19	0	0
1983	157											1	0
1984	92												0
1985	88												0
1986	51												0
1987	85												0

3.5.2 Set of Candidate Models

The biological objective of this study was to increase the survival process of sage grouse. Zablan (1993) used a global model, with 58 parameters. For the purpose of this study, the set of candidate models includes the following:

1	S, r	2	Constant S
2	S_t, r	16	Year-dependent S
3	S_a, r	3	Age-dependent S
4	S_{a+t}, r	17	Age- and year-dependent S , no
5	S_{a*t}, r	31	Age- and year-dependent S , in

Models with r year-dependent (t):

6	S, r_t	16	Constant S
7	S_t, r_t	29	Year-dependent S
8	S_a, r_t	17	Age-dependent S
9	S_{a+t}, r_t	30	Age- and year-dependent S , no
10	S_{a*t}, r_t	44	Age- and year-dependent S , in

Models with r age-dependent (a):

11	S, r_a	3	Constant S
12	S_t, r_a	17	Year-dependent S
13	S_a, r_a	4	Age-dependent S
14	S_{a+t}, r_a	18	Age- and year-dependent S , no
15	S_{a*t}, r_a	32	Age- and year-dependent S , in

Models with interaction terms (denoted by the $a * t$) all have its own set of time-dependent parameters. The addition of “+” exclude interaction terms; e.g., for model S_{a+t} the difference between subadult and adult survival parameters, year estimates of survival probabilities for subadults and adults, a logit scale and separated by β_0 (see below). A logit transformation made on S , and age (a) and year (as a dummy variable), function,

$$S_{a+t} \text{ denotes } \text{logit}(S) = \beta_0 + \beta_1(a) + \beta_2(t_1) + \beta_3(t_2)$$

This approach is similar to logistic regression. However, it is embedded in the log-likelihood function. Such models are biologically realistic and can be biologically realistic. Models without interaction terms; for example, model S_{a+t}, r_a has 18 parameters; for model S_{a*t}, r_a .

These 15 models plus models $\{S_a, r_{a*t}\}$ and $\{S_a, r_{a+t}\}$ and $\{S_{a*t}, r_{a*t}\}$ seem like sound initial choices; however, further iterations might lead one to exclude models with many parameters due to the relatively sparse data available (Table 3.8). We realized that a set of models would ideally be fine-tuned in a real-world application if a long-term increase or decrease in survival was hypothesized. We introduce submodels for survival such as

$$\text{logit}(S) = \beta_0 + \beta_1(a) + \beta_2(T) + \beta_3(a * T)$$

survival probability in these years (s_{aj} , s_{sj}) could have been used to differentiate them from the survival probability in more recent years. We assume that a great deal of thought has been put into the construction of candidate models.

A primary interest in banding studies is often to estimate survival probabilities and assess what external covariates might influence survival. (1993) modeled sage grouse survival using 4 year-dependent covariates (cov_t): winter precipitation (wp), winter temperature (wt), winter precipitation (sp), and spring temperature (st) (she provided the definitions of these variables; we will not need to note the units). Submodels with survival probabilities of the form

$$\text{logit}(S_t) = \beta_0 + \beta_1(cov_t) \quad \text{or} \quad \text{logit}(S_{t+a}) = \beta_0 + \beta_1(cov_t)$$

could be constructed. Such submodels for survival have one parameter (an intercept and one slope coefficient for the first submodel) and one age effect, and 1 slope coefficient for the second submodel. These models provide some insights into biological correlates, which themselves

3.5.3 Model Selection

Zablan's analysis was done using the programs ESTIMATOR (Brownie et al. 1985) and SURVIV (White 1983). Zablan fit the global model $\{S_{a*}, r_{a*}\}$ fit the data well ($\chi^2 = 34.34$ with 30 df, $p = 0.25$) and she computed a variance inflation factor from the global model $\{S_{a*}, r_{a*}\} = 1.14$. Her calculations are in agreement with those of the global model $\{S_{a*}, r_{a*}\} = 87.85$, 80 df, $\hat{c} = 1.10$). There was no evidence of overdispersion; thus there was no compelling reason to use a larger sample size for parameter estimation in these surveys is than the number of birds banded, which equaled 3,624 in this case. Zablan used 58 parameters, giving the ratio $n/K = 3,624/58 = 62$; this ratio has been safely used instead of AIC_c . We used the program AIC (Burnham 1999, White et al. 2001) to compute MLEs of the parameters, their conditional covariance matrix, the maximized value of the log-likelihood function, AIC_c , Δ_i , and w_i for each of the 17 candidate models. The weather covariate and 4 models with one of the weather covariates.

AIC_c selected model $\{S_a, r_a\}$ with 4 parameters (Table 3.5.3). This model assumes that conditional survival and reporting probabilities are constant over years (subadult vs. adult), but constant over years. Here, the male survival probability was 0.407 ($\hat{se} = 0.021$), while the estimated survival probability for subadults was higher, at 0.547 ($\hat{se} = 0.055$). The respective coefficients of variation were 5.2 and 10.0. An inference here is that female grouse survive at a higher rate than male adults; perhaps

Number	Model	$\log(\mathcal{L})$	K	AIC_c
Without environmental covariates:				
1	S, r	-2, 215.564	2	4435.13
2	S_t, r	-2, 205.074	16	4442.30
3	S_a, r	-2, 215.096	3	4436.20
4	S_{a+t}, r	-2, 203.797	17	4441.76
5	$S_{a* t}, r$	-2, 199.277	31	4461.11
6	S, r_t	-2, 204.893	16	4441.94
7	S_t, r_t	-2, 194.611	29	4447.71
8	S_a, r_t	-2, 204.526	17	4443.22
9	S_{a+t}, r_t	-2, 193.633	30	4447.84
10	$S_{a* t}, r_t$	-2, 188.531	44	4466.17
11	S, r_a	-2, 214.717	3	4435.44
12	S_t, r_a	-2, 204.544	17	4443.26
13	S_a, r_a	-2, 211.357	4	4430.72
14	S_{a+t}, r_a	-2, 204.544	18	4439.96
15	$S_{a* t}, r_a$	-2, 196.065	32	4456.72
16	$S_a, r_{a* t}$	-2, 197.572	32	4459.73
17	$S_{a* t}, r_{a* t}$	-2, 174.557	58	4467.03

of breeding and increased predation on breeding males. probabilities (the \hat{r}) for first-year subadult birds were also for adult birds (0.227 ($\hat{se} = 0.031$) and 0.151 ($\hat{se} = 0$). The use of the AIC_c -selected model does not indicate that dependent variation in the parameters, only that this variation is small in the sense of a bias–variance tradeoff and K-L information. The estimated next-best model ($\Delta_i = 4.41$) without a covariate with only two parameters, while the third-best model ($\Delta_i = 4.41$) with three parameters. There is relatively little structure revealed by this is not surprising, since the data are somewhat sparse.

The AIC -selected model assumed that the survival and reproduction varied by age class. Thus we considered four models: a linear function of age (subadult vs. adult) and one of the age classes (wp , wt , sp , or st), while retaining the age-specific report. The results were interesting and, at first, suggest that each of the four models is nearly tied with the AIC_c -selected model $\{S_a, r_a\}$, an important point, illustrated below:

Number	Model	$\log(\mathcal{L})$	K	AIC_c
With environmental covariates:				
18	S_{a+wp}, r_a	-2, 210.828	5	4431.6
19	S_{a+wt}, r_a	-2, 211.334	5	4432.6
20	S_{a+sp}, r_a	-2, 210.819	5	4431.6
21	S_{a+st}, r_a	-2, 210.802	5	4431.6

probabilities are related to temperature or precipitation. In the examination, it can be seen that the value of the maximized log-likelihood was very similar to the best model in Table 3.9, without any covariate. The inclusion of models with a covariate has not improved the fit to the data. The best model without a covariate has 4 parameters, while the covariate models have 5 parameters. This difference in the number of parameters explains most of the difference in the AIC_c values between the best model and the 4 models with a covariate. Hence, upon comparison, there is virtually no support for any of the covariates from the data. This leads to a point that is important in general.

Models Within Two Units of the Best Model

Models having Δ_i within about 0–2 units of the best model were examined to see whether they differ from the best model. The results show that they have essentially the same values of the maximized log-likelihood as the best model.

In this case, the larger model is not really supported; it is rather “close” only because it adds 1 parameter and therefore has a Δ_i of 2 units, even though the fit, as measured by the log-likelihood, was not improved.

Further insights into the sage grouse data can be obtained by examining the weights for the first 17 models in Table 3.9 (ignoring the last 10 models). In this case, the weight for the AIC_{min} model $\{S_a, r_a\}$ is 0.77, the second-best model $\{S, r\}$ has a weight of 0.085 (evidence ratio of 9.1), and the third- and fourth-best models had weights of 0.073 and 0.050, while the other models were nearly zero (the sum of the Akaike weights for the remaining models was < 0.02). The annual variation in the probabilities was small (temporal process variation $\hat{\sigma}_s = 0.0279$ for subadults); thus model $\{S_a, r_a\}$ seems reasonable. Models 1–4 all had fewer parameters than the AIC_{min} model. Thus, conditional sampling variances from those models were smaller than from the AIC_{min} model. In addition, these three models had small Akaike weights. This leads to some trust in the conditional sampling variances as a reasonable reflection of the precision of the parameter estimates.

3.5.4 Hypothesis Tests for Year-Dependent Survival

Zablan (1993) computed a likelihood ratio test between the global model $\{S_{a*}, r_{a*}\}$ (the global model) using the program BROWNE (Brown, 1993) and the evidence of year-dependent survival ($\chi^2 = 46.78$, 26 df, $P < 0.001$). MARK provides similar results ($\chi^2 = 46.03$, 26 df, $P < 0.001$). The test allowed a fairly general structure on the reporting probabilities. The results seemed convincing and provided evidence that survival

for year-dependent survival probabilities (i.e., S vs. S_T) and the test result been even more “significant.” However, of reasoning, the test of $\{S_a, r_a\}$ vs. $\{S_{a*}, r_a\}$ gives $\chi^2 : P = 0.336$. Still other testing strategies are possible, and might be deemed the best.

Given a believed year-dependence in annual survival (1993) asked whether this variability was partially explained by covariates, with or without an age effect. However, she was unable to find a relationship between annual survival probabilities and any covariate using likelihood ratio tests (the smallest P -value for the test was 0.194). She used model $\{S_{a*}, r_{a*}\}$ (the global model) as a baseline.

3.5.5 Hypothesis Testing Versus AIC in Model Selection

An apparent paradox can be seen in the results for the male survival analysis. This allows us to further compare alternative paradigms of hypothesis testing and AIC for model selection. The test between the test $\{S_a, r_a\}$ and $\{S_{a*}, r_{a*}\}$ attempts to answer the questions, “Given the data on the recovery probabilities, is there evidence that survival is (in addition to age) year-dependent?” The answer provided by the test is 0.009. But this answer is seemingly in contrast to the AIC-selected model, where there is no hint of time-dependence (Table 3.9). The Δ_i values for models $\{S_a, r_{a*}\}$ and $\{S_{a*}, r_{a*}\}$ are 30.400, respectively. AIC lends little support for a best model that includes year-dependent survival or reporting probabilities.

The answer to this paradox is interesting and important. The null hypothesis that $S_1 = S_2 = S_3 = \dots = S_{14}$ for all years is obviously false, so why test it? This is not properly a hypothesis test (see Johnson 1995 and Yoccoz 1991 for related issues). The

Model Interpretation

Sometimes, the selected model contains a parameter that varies with time, or areas, or age classes (i.e., $\theta = \theta_1 = \theta_2 = \dots$). This should not imply that there is *no* variation in this parameter. The parsimony and its bias/variance tradeoff finds the actual variation to be relatively small in relation to the information contained in the data. It “costs” too much in lost precision to add estimates of the individual θ_i . As the sample size increases, then at some point the estimates of the individual parameters would likely be false. Just because a parsimonious model contains a parameter that varies across strata does not mean that there is no variation in the parameter across the strata.

not relate directly to the issue of what approximating model is best for inference. One model selection strategy that has often been used is to do likelihood ratio tests of each structural factor (e.g., for each of the parameters S and r) and then use a model where the factors that were “significant” at, say, $\alpha = 0.05$. However, this strategy would suggest that this strategy would lead to a model with good properties (i.e., small bias, good precision, and achieved coverage at the nominal level).

Clearly, one must also worry about the multiple testing problem. Even if the fact that many such tests would not be independent, in the present case, then likelihood ratio tests are not chi-square distributed. Even if the test statistics may not be chi-square distributed for small sample sizes such as these, where sample sizes are far from asymptotic, then where many models contain nuisance parameters (the r_i), the α -level is arbitrary as well. Many of the models in Table 2.7.1 are thus such that likelihood ratio tests are not possible between these models. There is a certain lack of symmetry (this is again related to the difference between the null and alternative hypotheses and how this might relate to the “best approximating model” (see Section 2.7.2)). A very good problem here is how the test results are to be incorporated into the best model for statistical inference. This problem becomes acute when there are many (say, $> 8-10$) candidate models. Using just the set of models for sage grouse data, one would have 136 potential likelihood ratio tests. If some of these models were not nested, prohibiting a test between them. With 136 (or even 36) test results there is no theory or understanding of what the best model should be and no rigorous, general way of deciding (e.g., which model is second-best? Is the second-best model as good as the best, or substantially inferior?). Finally, what is to be done if the test results are inconsistent, such as those found in Section 3.

The biological question regarding annual survival probability can be stated as, “How much did annual survival vary during the study? Has survival varied little over the 14 or 15 years, or has there been a large variation in these annual parameters? Such questions are not ones of hypothesis testing (see Franklin et al. 2002). The proper inquiry should be on the amount of variation among the parameters (S_1, S_2, \dots, S_{14}) for each of the two age classes; we will assume that the deviation among these parameters by σ_s . Of course, if we have the S_i , then $\hat{\sigma}_s = \left(\sum_{i=1}^{14} (S_i - \bar{S})^2 / 13 \right)^{1/2}$.

We next ask why the AIC procedure did not pick up the model where survival varied by year? The reason is simple; AIC attempts to select the best approximating model for the observed data. In the sense of minimizing loss or a tradeoff between bias and variance, it was poor performance.

els $\{S_a, r_{a*}\}$ and $\{S_{a*}, r_{a*}\}$ to model the variation in S_t (Figure 1.3A) that the difference in AIC values for model $\{S_a, r_a\}$ versus $\{S_{a*}, r_{a*}\}$ was 29.01, suggesting that model $\{S_a, r_{a*}\}$ is highly overfit. The differences among survival probabilities are large enough to be included in the *model selection* problem, not one of hypothesis testing.

Estimates of the 30 survival probabilities under model $\{S_a, r_a\}$ are shown in Table 3.10. The average of the 15 estimates of adult survival is nearly the same as that from the 4-parameter model selected by the AIC. However, the average percent coefficient of variation for \hat{S}_t for model $\{S_{a*}, r_a\}$ compared to only 5.2 for \hat{S}_t in the AIC-selected model indicates that the best estimate of survival for a particular year is merely \hat{S}_t (from model $\{S_a, r_a\}$).

The situation was similar for subadult survival; the average percent coefficient of variation for \hat{S}_t for model $\{S_{a*}, r_a\}$ was 0.548, compared to 0.547 for the AIC-selected model. The respective average percent coefficients of variation were 35.6% for models $\{S_{a*}, r_a\}$ and $\{S_a, r_a\}$. In summary, 54 (or even 58) parameters “cost too much” in terms of increased variability (Figure 1.3B and Table 3.10) and reflect substantial overfitting. The coefficient of variation for \hat{S}_t for model $\{S_{a*}, r_a\}$ was worse than for model $\{S_{a*}, r_{a*}\}$ was used; coefficients of variation were 35.6% for adult survival and 30.5% for subadult survival. The model suggested by the AIC approach had 58 parameters, while the AIC-selected model had 14 parameters. This illustrates the cost of using overparametrized models. The results of hypothesis tests clearly show “significance” for adult survival (and reporting) probabilities. Models $\{S_a, r_{a*}\}$ and $\{S_{a*}, r_a\}$ are general for these data and lie far to the right of the bias–variance trade-off in Figure 1.2. Zablan recognized the problems in using models with many parameters and commented, “While significant differences were found between the recovery rates of males and of both age classes, and between the estimates had unacceptably wide confidence intervals.”

3.5.6 A Class of Intermediate Models

The researcher could use the AIC-selected model $\{S_a, r_a\}$ with 4 parameters and then proceed to obtain an estimate σ_s of the variance of parameters and then proceed to obtain an estimate σ_s of the variance of parameters for each age class, using model $\{S_{a*}, r_{a*}\}$ or $\{S_{a*}, r_{a+t}\}$ if desired. This parameter is embedded into the likelihood framework, allowing for the estimation of σ_s ; this becomes a “random effects” model. However, the model is a “variance components,” and consistent estimates of σ_s can be obtained. For example, say, model $\{S_{a*}, r_a\}$, following, for example, Anderson et al. (1987:62–66); and Burnham et al. (1987:260–269). This approach is based on the simple partitioning of the total variance (

to first year survival of subadults. The model assumes that the subadults that survive to the second year after banding and thus have the same year-dependent mortality as the birds banded as adults.

Year(<i>i</i>)	\hat{S}_i	Standard Error	95% Confidence Interval	
			Lower	Upper
1	0.462	0.128	0.238	0.686
2	0.500	0.092	0.327	0.673
3	0.357	0.073	0.230	0.484
4	0.412	0.074	0.277	0.547
5	0.464	0.073	0.328	0.599
6	0.507	0.069	0.375	0.640
7	0.465	0.066	0.340	0.590
8	0.357	0.062	0.246	0.468
9	0.397	0.063	0.282	0.512
10	0.340	0.061	0.233	0.447
11	0.321	0.063	0.212	0.430
12	0.358	0.073	0.231	0.485
13	0.171	0.071	0.071	0.271
14	0.339	0.133	0.138	0.540
15	0.549	0.129	0.305	0.793
1	0.725	0.114	0.462	0.988
2	0.629	0.152	0.321	0.937
3	0.524	0.106	0.323	0.725
4	0.528	0.112	0.316	0.740
5	0.566	0.093	0.383	0.749
6	0.446	0.120	0.237	0.655
7	0.386	0.117	0.193	0.579
8	0.513	0.110	0.307	0.719
9	0.497	0.118	0.282	0.712
10	0.615	0.111	0.389	0.841
11	0.547	0.101	0.351	0.743
12	0.368	0.121	0.173	0.563
13	0.440	0.107	0.251	0.629
14	0.744	0.104	0.498	0.990
15	0.695	0.111	0.450	0.940

^aBased on a back transformation of the interval endpoints on a logit scale (see Appendix 1).

additive components: the variance in the population parameter estimates ($\text{var}(\hat{S}_i)$), the conditional sampling variance ($\text{var}(\hat{S}_i \mid \text{model})$). The approach assumes that the true S_i are independently and identically distributed random variables. In this case, these assumptions are weak and the effect somewhat mitigated by the fact that one has ML estimates of the sampling covariance matrix $\text{var}(\hat{S}_i)$ directly from the estimates \hat{S}_i ; thus, by subtracting the conditional sampling variance from the total variance, one can obtain an estimate of the variance in the population parameter estimates.

The estimate of ϕ_s would provide some insight into survival parameters; this would be done in the context of the S_i vary. Exact details of the optimal methodology would be left for the future; however, some unpublished results seem exciting: $\hat{\phi}_s = 0.0426$, 95% likelihood interval $[0, 0.106]$, and $cv = 4.9\%$ on S . Thus, one can infer that the true annual survival probabilities were fairly stable (the variation in the true annual survival probabilities was fairly small, about 10%). Thus, the large variation in the estimates of annual survival probabilities under model $\{S_{a*}, r_a\}$ (Table 3.9) is due primarily to sampling error. The large estimated standard errors suggest. Additional details, including the estimates of annual survival probabilities, appear in Chapter 3 (see White and White 2002).

Ideally, the number of parameters in the various candidate models should be small (see Section 2.7.2). In the grouse model, the survival probabilities without year-dependent survival would have two (if age is included) parameters, while a model with year-dependent survival would have as many as 30 parameters (15 for each of the two age groups). The differences in the number of parameters between certain candidate models are not ideal, and one should consider intermediate models with a small number of parameters. Zablan's (1993) various covariate models are an example; here an intercept and slope parameter on one of the covariates would introduce 2 parameters (3 with age) instead of 15 (30 with age). With the hypothesis testing approach, AIC-selection shows that the covariate models to be essentially tied with the AIC-selected model (Table 3.9).

Hypothesis testing and AIC are fundamentally very different approaches to model selection and in drawing inferences from a set of data. For example, AIC_c tries to select a model that well approximates the data in the data. That selected model then provides estimates of S_1, S_2, \dots, S_{15} for each age group in the sense of K-L in the sense of bias versus variance tradeoff). That is, an estimate of average survival probability model $\{S_a, r_a\}$ (i.e., 0.407 for adults) would be used to estimate S_5 for adult grouse (hence, $\hat{S}_5 = 0.407$, 95% confidence interval $[0.328, 0.606]$), and this estimate would have better inferential properties than the estimate of $\{S_{a*}, r_a\}$, whereby S_5 would be estimated using the year-specific survival probabilities (see Table 3.10, where this estimate is given as 0.464 with 95% confidence interval of 0.328 to 0.606). If inference about the condition of the fifth year is made from the general model $\{S_{a*}, r_{a*}\}$, then the precision is worse yet (95% confidence interval of $[0.328, 0.606]$). In the last two cases, the precision is relatively poor (e.g., compare the results of A and C for further insights).

This example illustrates the use of information-theoretic analysis of count data displayed as multidimensional contingency tables (Sakamoto 1982). Schoener (1970) studied resource utilization by two species of the genus *Anolis* on several islands in the Lesser Antilles, Caribbean. Here we use his data collected on *Anolis gramineus* near Whitehouse, on Jamaica, as provided by Bishop et al. (1970). The data have been analyzed by Fienberg (1970), Bishop et al. (1970), Bishop and Nelder (1989) and Qian et al. (1996), and the reader is encouraged to compare the approaches given in these papers to that presented here.

In his general studies of species overlap, Schoener (1970) studied the distribution of species of lizards in an area of trees and shrubs that had been recently cleared. The height (< 5 or ≥ 5 ft) and diameter (< 2 or ≥ 2 in) of the trees (sunny or shaded), and time of day (roughly early morning or afternoon) were recorded for each of the two species of lizards. The data were taken on each individual only once per “census”; data were not collected if a lizard was disturbed, and the census route was varied each day. Table 3.11 (from McCullagh and Nelder 1989:128–135) summarizes the observations of lizards appear in the table with 48 cells. The data here differ from those in the first edition due to errors (both in the data and in computing AIC_c .)

TABLE 3.11. Contingency table of site preference for two species of lizards, *A. opalinus* (denoted by g and o, respectively) on the island of Jamaica (Schoener 1970).

Insolation (I)	Diameter (D) in	Height (H) ft	Time of Day (T);			
			Early morning		Midday	
			g	o	g	o
Sunny	≤ 2	< 5	20	2	8	0
		≥ 5	13	0	8	0
	> 2	< 5	8	3	4	0
		≥ 5	6	0	0	0
Shaded	≤ 2	< 5	34	11	69	20
		≥ 5	31	5	55	4
	> 2	< 5	17	15	60	32
		≥ 5	12	1	21	3

pertise in lizard ecology and behavior, decisions concerning the selection of candidate models will necessarily be somewhat controversial. In addition, attention on modeling and model selection issues as an extension of the work on several inference issues. In reading Schoener's (1970) paper, the literature he cited, it would seem that a model with all the main effects (D, I, T, and S) might serve as a starting point for model selection. Several second-order interactions might be suspected, e.g., D * I, D * T, and I * T. If the two species are partitioning their resources, then the terms H * S, D * S, I * S, and T * S included should be reasonable. If designed and data were collected, it was probably evident that the response was affected by several variables as well as some interactions. This suggests that a model with all main effects and second-order interactions should be considered. Then issues remain concerning possible higher-order interactions. On biological grounds, it might seem reasonable to consider interactions such as H * D * I, H * D * T, and H * D * S; or further interactions such as H * T * S, and I * T * S. Finally, the second-order term D * I might be important; thus some models without this term were considered. We use the short set of models in Table 3.12 for illustrative purposes. In the late 1960s, T. W. Schoener and his colleagues, including me, could have developed a better set of a priori candidate models.

3.6.2 Comments on Analytic Method

We used a loglinear model with Poisson errors following Agresti (1984, pages 453–456), and the analysis was made conditional on the total frequencies, $\sum n_j$. Specifically, we used the SAS program GENMOD (SAS, 1985). The likelihood is, $\mathcal{L}(\mu | n_j, \text{model}) = \prod_{j=1}^{48} \frac{\mu_j^{n_j} e^{-\mu_j}}{n_j!}$. The log-likelihood for the global model can be expressed as

$$\log(\mathcal{L}) = \sum_{j=1}^{48} (n_j \cdot \log(\mu_j) - \mu_j) - \sum_{j=1}^{48} \log(n_j!)$$

where n_j is the number of observations in cell j with Poisson mean μ_j , $j = 1, 2, \dots, 48$. The purpose of the modeling is to put some structure on the 48 means. Then one has the log-linear model for the analysis of variance. Thus, β is the vector of effects and the intercept. The final term in the log-likelihood is a constant; thus SAS computes the deviance term, and the resulting log-likelihood is positive. AIC is computed in the same manner, even though the AIC values are scaled by $\sum_{j=1}^{48} \log(n_j!)$, which is such arbitrary, additive constants are not present in the log-likelihood. The software packages allow ML estimates from discrete data. We will provide a number of relevant analysis options (see summary of results).

if overdispersion is thought to exist in count data. We might comment, “In practice, we learn more from *estimating* densities than from *testing hypotheses* about their values.”

3.6.3 Some Tentative Results

A model with only the five main effects might be a starting point for the information in these data. Alternative models with various two- and three-way interaction terms are shown in Table 3.12. None of the first-, second-, or third-order interactions are supported by the data; in fact, the model with any support is the model with merely the main effects. This is in agreement with McCullagh and Nelder (1989), but different from Fienburg (1970) and Bishop et al. (1975).

Note that if another model were to be added, the Δ_i values would have to be recomputed (but not the AIC values). Such values are relative to the minimum AIC model, given a set of candidate models. The models derived from the data via a set of models are effectively a subset of the set of models considered. Strict experimentation might be required for additional insights into the issue of resource utilization in the field.

If our original interest had been only on differences in resource utilization between two species, we could examine a model with all five main effects (H, D, I, T, S) vs. a similar model without a species effect (H, D, I, T, S). Only two candidate models would have been considered. The difference in AIC values between the two models is 163.70, indicating strongly that the two species are utilizing resources differently. Alternatively, we could make a similar comparison between two models, but also include the relevant second-order interactions.

H, D, I, T, S, H*D, H*I, H*T, H*S, D*I, D*T, D*S
vs.

H, D, I, T, H*D, H*I, H*T, D*I, D*T, I*S

Here, the difference Δ_i is 84.63 and again clearly indicates that the two species are utilizing their habitat differently. Other alternative models might be pursued if more were known about the study design and the species involved.

Bishop et al. (1975) performed a number of hypothesis tests for the data and this resulted in a model with all the main effects plus the two-way interaction terms H*D, H*S, D*S, T*S and the third-order interaction H*D*S. If the study is considered to be merely exploratory, then one might consider a wider class of candidate models, and many models with various combinations of interaction terms might be found. However, the number of possible models would be 2,000, depending on what rules might be applied concerning the inclusion of lower-order effects if higher-order effects are included. With powerful computing equipment, an exhaustive study of the data might be possible.

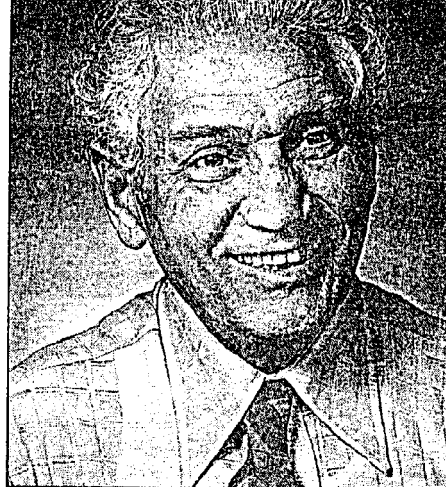
The model with the minimum AIC_c is shown in bold.

Model	$\log(\mathcal{L})$	K	AIC_c
1 All main effects, H D I T S	1,181.08	7	-2,347.1
2 All main effects and second-order interactions	1,181.86	21	-2,311.1
3 Base ⁽¹⁾ but drop DT	1,180.52	19	-2,320.1
4 Base ⁽¹⁾ plus HDI, HDT and HDS terms	1,182.97	25	-2,311.1
5 Base ⁽¹⁾ plus HDI, HDS, HIT, HIS, HTS, and ITS	1185.75	30	-2,300.1
6 Base ⁽¹⁾ plus HIT, HIS, HTS, and ITS	1185.48	28	-2,311.1
7 Base ⁽¹⁾ plus HIS, HTS, and ITS	1184.01	26	-2,311.1
8 Base ⁽¹⁾ plus HIT, HIS, HTS, and ITS, but drop DT	1184.15	26	-2,311.1
9 Base ⁽¹⁾ plus HIT, HIS, and ITS, but drop DT	1183.29	24	-2,311.1
10 Base ⁽¹⁾ plus HIT and HIS, but drop DT	1182.18	22	-2,311.1

⁽¹⁾ “Base” is a model with the five main effects plus all second-order interactions.

this $2 \times 2 \times 2 \times 3 \times 2$ table is nearly prohibitive (see Agresti 1996). This illustrates again the importance of a set of good a priori hypotheses. This study is somewhat exploratory. Of course, data dredging can be useful in finding a model that “fits” *these* data. However, the goal of such studies usually stress an inference about the population or process generating the data description. Thus, the results from intensive data dredging are often as tenuous.

The rigorous analysis of multidimensional contingency tables is often problematic because of the large number of possible models. Good (1978:19) suggest that conclusions drawn from contingency tables should be only exploratory. While a good a priori set of models seems to be difficult to forecast higher-order interactions in many situations, several problematic issues are associated with contingency table analysis. Analytic approaches exist (Manly et al. 1993), and alternative link functions besides a log link and Poisson errors can be considered (Sokal and Rohlf 1989 for additional information).



Solomon Kullback was born in 1907 in Brooklyn, New York. He graduated from the City College of New York in 1927, received an M.A. degree in mathematics from the University of Chicago in 1930, and completed a Ph.D. in mathematics at the George Washington University in 1933. He was known to all who knew him, had two major careers: one in the Department of Mathematics at the University of Chicago (1930–1962) and the other in the Department of Statistics at George Washington University (1962–1972). He was chairman of the Statistics Department from 1962 to 1972. The latter part of his professional life was spent in the National Security Agency, and much of the work done during this time is still classified. Clearly, most of his studies on information theory were done during this time. Many of his results up to 1958 were published in his book *The Theory of Information*. Additional details on Kullback may be found in the book by Kullback and Anonymus (1997).

3.7 Example 6: Sakamoto et al.'s (1986) example

Here we return briefly to the example used in Section 1.4. Sakamoto et al.'s (1986) book. Ten data sets (each with $n = 21$) were generated from a simple model

$$y = e^{(x-0.3)^2} - 1 + \epsilon.$$

Sakamoto et al. used the simple polynomials from order 0 to 7 as candidate models to visually illustrate the concepts of uniformity (see Figure 1.4). Because the sample size ($n = 21$) is small relative to the dimension of the largest model in the set ($K = 7$), AIC_c should be used in the analysis of data in this example. AIC_c was computed for each of the ten data sets and then averaged for each of the 7 models. Then Δ_i was computed and derived from these averages:

Quadratic	4	-29.63	2.88	0
Cubic	5	-26.80	5.71	0
4th-order	6	-23.79	8.72	0
5th-order	7	-19.26	13.25	0
$\min f(x) = 0.3$	2	-32.51	0.0	0

Clearly, the model based on the additional (hopefully a) that $f(x)$ is 0 at $x = 0.3$ is the best of the set (compare results with the plots in Figure 1.4). The Akaike weights sharpen the inference and suggest that only the quadratic is a competitor to the special model for these simulated data. The ratio between the best and second-best models is $0.77/0.001$; this ratio between the best and third-best model is only $0.001/0.000001$. These comparisons are in line with the visual images in Figure 1.4. They reinforce understanding of the information-theoretic quantities in a simple example. More complex data and models defy simple plots and require such K-L information and various information criteria before

3.8 Example 7: Models of Fish Growth

Shono (2000) presented a reanalysis of data on the growth of rainbow salmon (*Oncorhynchus masou*) from Kiso et al. (1992). Shono presented a comparison of model selection under AIC, AIC_c , BIC (Bayesian criterion, Schwarz 1978); we will focus on model selection based on the inferences made by using Δ_i , w_i , and evidence ratios.

Kiso et al. (1992) estimated the parameters of three standard models using length and age data over a period of 2 to 19 months. We give a brief review of the information; further details can be found in Shono and, particularly, in Shono (2000). The sample size was $n = 100$ but was large relative to the number of parameters in the best model (where $K = 5$). The models for length $L(t)$ as a function of months, were

von Bertalanffy	$L(t) = L_\infty[1 - \exp\{-\kappa(t - t_0)\}]$
Gompertz	$L(t) = L_\infty \exp[-\exp\{-\kappa(t - t_0)\}]$
Logistic	$L(t) = \frac{L_\infty}{1 + \exp\{-\kappa(t - t_0)\}}$

The basic model parameters are L_∞ and κ . Each model also has t_0 in the exponent, and this was itself modeled in four ways to capture concerning the seasonal pattern in growth. These submodels

$$\begin{array}{ll} \text{Type 2} & F(t) = t + \frac{\theta_1}{2\pi} \sin 2\pi(t - t_1) + \frac{\theta_2}{4\pi} \sin 4\pi(t - t_1) \\ \text{Type 3} & F(t) = t + \frac{\theta_1}{2\pi} \sin 2\pi(t - t_1) + \frac{\theta_3}{6\pi} \sin 6\pi(t - t_1) \end{array}$$

The unknown parameters in these submodels are θ_1 , θ_2 , and θ_3 , in the case, $\theta_i \geq 0$. The θ_i are the amplitudes of sine curves with one, two, and three cycles/year, respectively, and t_1 is the starting time of the curves. They define Type 1 as a modified type where the growth rate changes once per year. Type 2 allows growth rate changes in a constant number of cycles per year, whereas Type 3 allows growth rate changes in a constant number of 1 and 3 cycles per year (Kiso et al. 1992:1780).

Each of the three growth models included four submodels, giving a total of 12 models. The results, taken partially from Table 3.13, are given in Table 3.13. As shown by Shono (2000), the von Bertalanffy with the Type 1 seasonal effect. However, looking at the Δ_i values, we note that the Gompertz model with Type 2 is essentially tied ($\Delta = 0.72$); the evidence ratio is 0.27/0.13. The Gompertz Type 3 ($\Delta = 1.13$) and the Gompertz Type 2 are also close competitors. Three of the best four models in the top form. Even the worst of the top four models is still quite good. The ratio with the best model is 0.27/0.13 = 2.1. Clearly, the Type 2 model, relative to the other types. In fact, all nine models with a seasonal pattern have Δ_i values < 6, suggesting that some seasonal pattern is important in salmon growth. Beyond this, there is considerable uncertainty. This is a clear case where inference based only on the best model is risky. Inference, including prediction, should probably include twelve models or, at least, the nine models allowing seasonal pattern. In addition, estimates of precision should allow for the high uncertainty in selection. These are subjects treated in Chapter 4.

3.9 Summary

The purpose of the analysis of empirical data is not to find the best model, not at all. Instead, we wish to find a best approximating model to the data, and then develop statistical inferences from this model. The model is *the* inference from the available data. We search for a “best model,” but rather for a *parsimonious model* giving an account of the interpretable information in the data at hand. Data analysis is a question, “What level of model complexity will the data support?” The question under- and overfitting are to be avoided. Larger data sets tend to

Formulae	Type	K	AIC_c	Δ_i
von Bertalanffy	Basic	2	135.82	117.26
Gompertz	Basic	2	50.56	32.00
Logistic	Basic	2	54.05	35.49
von Bertalanffy	1	4	18.56	0.00
Gompertz	1	4	19.28	0.72
Logistic	1	4	21.02	2.46
von Bertalanffy	2	5	21.03	2.47
Gompertz	2	5	20.08	1.52
Logistic	2	5	24.52	5.96
von Bertalanffy	3	5	21.30	2.74
Gompertz	3	5	19.69	1.13
Logistic	3	5	24.25	5.69

complex models, and the selection of the size of the model between bias and variance.

The analysis of data under the information-theoretic approach is simple. That is, the computational aspects are simple, and to understand and interpret if one has the value of the maximum likelihood ($\log(\mathcal{L})$) or the residual sum of squares (RSS) for each model, quantities are routinely printed by nearly any commercial software. Computation of AIC , AIC_c , or $QAIC_c$ from either is so simple to the point that it can easily be done by hand. Similarly, the differences (Δ_i), Akaike weights (w_i), and evidence ratios are simple to calculate. The general approach is flexible enough to be used in a wide range of practical situations in the life sciences. These are all approaches that go back to Kullback–Leibler information and have a deep theoretical foundation. The easy part of the information-theoretic approaches includes the computational aspects and the clear understanding of these results (e.g., evidence).

The hard part, and the one where training has been so common, is the thinking about the science of the matter before data analysis. It has been too easy to collect data on a large number of things with the hope that a fast computer and sophisticated software will find the important things—the “significant” ones (the “just the numbers” approach). A major effort should be mounted to understand the nature of the problem, a critical examination of the literature, talking with others about the problem, and thinking deeply about alternative hypotheses. “test” dozens of trivial matters (is the correlation zero? is the treatment zero? are ravens pink?, Anderson et al. 2000), the concerted effort to provide evidence on *meaningful* questions.

The compartment models of DOKSÖDERN are a simple of first-order differential equations, and they represent a high complexity. The real data on cement hardening and fish growth are examples but provide insights into the interpretation of the solution of the analysis of the simulated data from Sakamoto. Some comparisons with the graphical material in Section 3.2 of resource partitioning in *anolis* lizards must be considered. Burnham et al. (1996) give a comprehensive example of the survival probabilities of the Northern Spotted Owl (*Strix nebulosa*) caused by widespread clear-cutting in national forests. They use a large data set over eleven geographic areas in a political region of national importance.

Researchers often attempt to perform some further, separate (multiple linear or logistic regression) of the estimates from the model in an effort to understand the structure of the process leading into its behavior. However, such external analyses are not enough. The estimates for the starling and grouse data have a multivariate variance structure; the variances are unequal (not constant), the estimates are dependent (not independent). When the analysis is such further "external" analysis of the parameter estimates is clear that a properly parsimonious, easily interpretable model has not been achieved, and hence the analysis has partially failed.

We recommend carefully developing a set of candidate models for the science of the issue (e.g., embed additive submodels for the data into the log-likelihood using an appropriate link function), and then find the MLEs under these models. Then one can focus on model selection to find a properly parsimonious model(s) (the models in the set of candidates that are "close" to truth in the K-L information sense) that will lead to good inference. In contrast, if only the high-dimensional global model is used, and estimates of parameters obtained by ML or LS, then the analysis is virtually defeated, because a parsimonious interpretation may be impossible, patterns often cannot be found, and estimates are not precise. Zablan (1993) used a global model with 58 parameters to analyze the sage grouse data and correctly observed, "... such a model has unacceptably wide confidence intervals."

Statistical analysis of empirical data should not be just a description of the data given only a set of data (the numbers). The cement hardening data (Section 3.2) have too frequently been analyzed without examining the treatment of the data: the important a priori considerations are not taken into account here. What was known about the chemistry of cement hardening in the 1930s, then at least by the 1960s)? For example, can cement harden well with only a single ingredient? If not, this might lead to single-variable models out of consideration. Given the complexity of the data, the

variables x_1 and x_3 is similar. Surely, these conditions were not considered, or could have affected the models in the set to be considered. It is not clear how the researchers have known, or bothered to find out, that $\sum_{j=1}^4 x_{ij}$ is a constant for each of the thirteen observations and therefore excluded the four-variables approach, based on these a priori grounds? Instead, the researchers have been the modus operandi, and “all possible models” were tried. “Let the computer find out” is a poor strategy for researchers. It is better to bother to think clearly about the problem of interest and then try to find out. *The sterile analysis of “just the numbers” will continue to be a major barrier to progress in the sciences.*

Researchers often resort to using a computer program to try all possible models and variables automatically. Here, the computer will discover the important variables and relationships (the “all possible models” approach void of any thinking or science). Coors et al. (1988) conducted stepwise linear regression analyses using AIC (SAS Institute 1988:786) in a study of elk (*Cervus canadensis*). They found that this approach “. . . provided results that were statistically significant, but ecologically unrealistic, unstable due to multicollinearity, and overfitted (too many variables).” The literature is full of such failed studies; just because a selection criterion does not mean that valid inferences can be drawn. The primary mistake here is a common one: the failure to consider a priori models, each representing a *plausible* research hypothesis.

The presentation of results in scientific publications should be based on the model selection process used in arriving at a set of candidate models. The model selection process should be based on the study design and the alternative research hypotheses. The presentation and discussion of the $\log(\mathcal{L})$ values, K , the appropriate selection criterion for each model, Δ_i , and w_i is recommended (Burnham and Anderson 2001d). Evidence ratios should be presented with other model selection criteria to allow a comprehensive assessment of the alternative hypotheses. This information allows the merits of each model to be contrasted. If no model selection criteria were used, data dredging was done following the formal analysis, then the results should be clearly noted and the tentative insights from these activities should be discussed.

All Possible Models

Unthinking approaches have been the common modus operandi. “Let the computer find out” is a poor strategy and usually reflects the fact that the researchers did not bother to think clearly about the problem of interest and the study setting.

The sterile analysis of “just the numbers” will continue to be a major barrier to progress in the sciences.

Formal Inference From More Model: Multimodel Inference

4.1 Introduction to Multimodel Inference

Model selection is most often thought of as a way to select a model, then inference is conditional on that model. However, in multimodel approaches are more general than this simplistic concept. Given a set of models, specified independently of the sample, we can perform formal inferences based on the entire set of models. Here, inference is on all the models in the set and this has several advantages. It does reinforce the importance of having a good set of models to represent the scientific hypotheses of interest. Part of multimodel inference includes ranking the fitted models from best to worst, based on a criterion, and then scaling to obtain the relative plausibility of each model, by a weight of evidence (w_i) relative to the selected best model. Conditional sampling variance ($\text{var}(\hat{\theta}|x, g_i)$) from each model, weighted by w_i , unconditional inferences about precision can be obtained for the entire set of models. Model-averaged parameter estimates and their unconditional sampling variances can be easily computed. Model selection uncertainty is a substantial subject in its own right, well beyond the scope of determining the best model.

By *unconditional*, we mean not conditional on any particular model. However, inference is still conditional on the full set of models considered. For spending the time to arrive at a good set of models, based on theory or hypothesized about the science underlying the study, the set of models should be small and well justified, at least for con-

in so doing? Second, how can model selection uncertainty be incorporated into estimates of precision? Third, how can the importance of predictor variables in analyses such as linear or logistic regression be assessed? Finally, how can a confidence set of models be chosen? K-L best model? More research is required to develop and evaluate methods for these issues, but we provide several approaches. Four examples are provided to illustrate the use of multiple models in formal inference.

4.2 Model Averaging

4.2.1 Prediction

Consider model-based inference for prediction, where R models are considered, each having the parameter θ as the predicted value of Y . Model i allows an estimate of the parameter, $\hat{\theta}_i$. If one of the models is K-L best (e.g., if its $w \geq 0.90$), then inference could proceed conditionally, on the selected best model. However, it is often the case that a single model is clearly superior to some of the others in the set. If the predicted value ($\hat{\theta}$) differs markedly across the models (i.e., if the $\hat{\theta}_i$ for the models $i = 1, 2, \dots, R$), then it is risky to base prediction on the selected model. An obvious possibility is to compute a weighted average of the predicted value, weighting the predictions by the Akaike weights.

Model Averaging

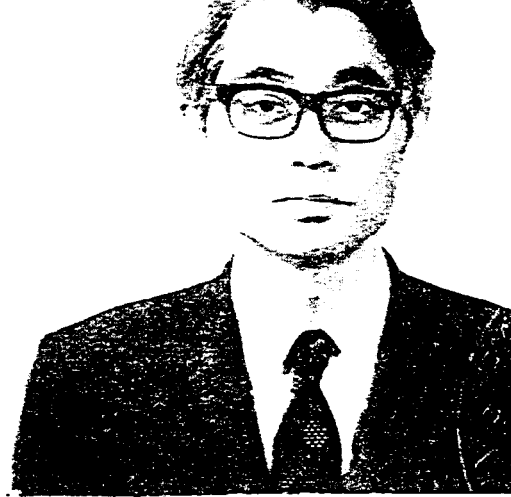
This concept leads to the model averaged estimates,

$$\hat{\theta} = \sum_{i=1}^R w_i \hat{\theta}_i,$$

where $\hat{\theta}$ denotes a model averaged estimate of θ . Alternatively, the Akaike weights are used to provide the estimated model selection frequencies. Model averaging can be done using,

$$\hat{\theta} = \sum_{i=1}^R \hat{\pi}_i \hat{\theta}_i.$$

This type of model averaging is useful for prediction problems where a particular parameter (e.g., γ an immigration parameter) is of interest for all the models in the set. Prediction is an ideal way to validate a model because each model in a set, regardless of its parametrization, can be used to make a predicted value.



Hirotugu Akaike was born in 1927 in Fujinomiya-shi, Shizuoka-ken. He received B.S. and D.S. degrees in mathematics from the University of Tokyo in 1950 and 1953, respectively. He worked at the Institute of Statistical Mathematics for 15 years and became its Director General in 1982. He has received many awards, prizes, and honors in theoretical and applied statistics (deLeeuw 1992, Parzen 1994). The *“Proceedings of the First US/Japan Conference on the Frontiers of Statistics: An Informational Approach* (Bozdogan 1994) commemorated Professor Akaike’s 65th birthday. Bozdogan (1994) records that the idea of a connection between the Kullback-Leibler discrepancy and the empirical log-likelihood function occurred to him on the morning of March 16, 1971, as he was taking a seat on a commuter train.

4.2.2 Averaging Across Model Parameters

If one has a large number of closely related models, such as those arising from stepwise or based variable selection (e.g., all subsets selection), design selection, or model selection, the model is unsatisfactory because that “best” model is often not stable. That is, the model estimated to be best would vary from data set to data set. If multiple replicate data sets would be collected under the same conditions, then in this situation, model averaging provides a relatively more stable inference.

The concept of inference being tied to all the models considered in the model selection bias effects on linear regression coefficient estimation in the all subsets selection. For the linear regression coefficient β_j of the predictor variable x_j there are two versions of model averaging: one is the estimate $\hat{\beta}_j$ where β_j is averaged over all models in which x_j is included.

$$\beta_j = \frac{\hat{\beta}_j}{w_+(j)},$$

$$w_+(j) = \sum_{i=1}^R w_i I_j(g_i),$$

and

$$I_j(g_i) = \begin{cases} 1 & \text{if predictor } x_j \text{ is in model } g_i \\ 0 & \text{otherwise.} \end{cases}$$

Here, $\hat{\beta}_{j,i}$ denotes the estimator of β_j based on model g_i . It is merely the sum of the Akaike weights over all models where predictor variable j is explicitly in the model. Note, $w_+(j)$ is an average value about whether variable x_j is in (or not in) the model.

Thus, $\hat{\beta}_j$ is a “natural” average to consider, as it only averages over models where an unknown β_j parameter appears. Note, however, it ignores evidence about models g_i wherein $\beta_{j,i} \equiv 0$.

An alternative way to average over linear regression models is that variable x_j is “in” every model, it is just that in some models the corresponding β_j is set to zero, rather than being considered unknown. When on model g_i being selected, model selection has the effect of setting β_j from zero (Section 1.6). Thus, a second model-averaged estimator, $\tilde{\beta}_j$, is suggested:

$$\tilde{\beta}_j = w_+(j) \hat{\beta}_j.$$

This $\tilde{\beta}_j$ actually derives from model averaging over all models, even where x_j is not in a particular model, it is because $\beta_{j,i} \equiv 0$ in those models. The resultant average is identical to $w_+(j) \hat{\beta}_j$. The estimate $\hat{\beta}_{j,i}$. The resultant average is identical to $w_+(j) \hat{\beta}_j$. $w_+(j)$ serves to *shrink* the conditional $\hat{\theta}_j$ back towards zero. This serves to ameliorate much of the model selection bias.

Investigation of this general idea and its extensions are an ongoing project.

One point here is that while $\hat{\beta}_j$ can be computed ignoring those models where x_j does not appear, $\tilde{\beta}_j$ does require fitting all R of the models. Improved inference requires fitting all the a priori models. This is a type of model averaging. When possible, one should use all the models, via model averaging and selection bias adjustment, rather than a “select the best model and ignore the others” strategy.

There are several advantages, both practical and philosophical, to model averaging, when it is appropriate. Where a model averaging approach is used it often has reduced bias and, sometimes has better predictive performance.

motivating ideas). Bayesian model averaging is easy to be difficult to implement in practice. Information-theoretic averaging are easy both to understand and implement, even a large number of models, each with potentially many parameters.

While there are many cases where model averaging is useful, model averaging structural parameter estimates in some cases. While it is often appropriate to average slope parameters in regression models, structural parameters in nonlinear models

$$E(y) = (a + b_x)/(1 + cx) \quad \text{or} \quad E(y) = a(1 - [1$$

should not be averaged. For example, a weighted average of models of any of the parameters a , b , c , or d would not be useful. Model averaging the predicted expected response variable value of x , across models, is advantageous in reaching a conclusion that is not conditional on only a single model.

It is important to realize that the expected value of the model parameter estimate, $E(\hat{\theta})$, is not necessarily the same as θ from absolute truth. In sampling theory the estimator $\hat{\theta}$ ($\equiv \hat{\theta}_i$ for the selected model by sample), arrived at in the two-stage process of model selection and parameter estimation given the model, is by definition an unbiased estimator of $E(\hat{\theta})$ as given by (4.1 or 4.2). Therefore, the unconditional expected value of $\hat{\theta} \equiv \hat{\theta}$ is to be computed with respect to $E(\hat{\theta})$. Any remaining bias in $\hat{\theta}$ cannot be measured or allowed for in model selection and averaging. This is part of the intent of having a good set of models and sound averaging to render this bias negligible with respect to the unconditional expected value.

Model-averaging ideas are well developed from the Bayesian perspective (see Madigan and Raftery 1994, Draper 1995, Raftery 1997, Hoeting et al. 1999; Newman 1997 provides an application of model averaging). This has not yet been commonly adapted into applied frequentist statistics. A theoretical basis for these approaches and ideas appears in the work of Buckland et al. 1997 and the Bayesian references just above.

4.3 Model Selection Uncertainty

An understanding of statistical inference requires that one understands the process that generates the sample data we observe. For a given computer simulation study, data are observed on some particular data set. A second, independent data set could be observed on the same data set. If, under nearly identical conditions, the new data set would be observed, the first. Clearly both data sets would contain information about the process that generates the data.

data observed. However, inferences must not be overly tied to the (single) data set observed. That is, we would like to be robust, with respect to the particular data set observed, we tend to avoid problems associated with both underfitting (overinterpreting) the limited data we have. Thus, we would like to make inferences about the process as if a large number of data sets were also available.

With only a single data set, one could use AIC, and so on, for inference. However, if several other independent data sets were available, would the same model be selected? The answer is that not necessarily, but generally, there would be variation in the selected model across data sets, just as there would be variation in parameter estimates given that the same model is used for analysis. The fact that the same model is used for analysis. The fact that the same model might suggest the use of other models leads us to model selection and hence another variance component that should be included in the precision of parameter estimates (Section 1.7).

If an analyst selects a model using AIC_c (or using some other criterion such as cross-validation) and makes estimates of the same parameter in that model, he invariably does so using the selected model. The estimated precision will then likely be too small because the variance component due to model selection is omitted. The standard errors computed conditional on the selected model are small, confidence intervals will be too narrow, and achieved coverage will be below the nominal level. Chatfield (1995b) reviews this issue. Also see Rencher and Pun (1980), Chow (1981), Hurvich and Tsai (1988), Pötscher (1991), Goutis and Casella (1995), and Kabaila (1998).

This section presents a variety of methods that can be used to assess the uncertainty associated with model selection, either with respect to the model, or regarding uncertainty about selected variables; and measures of unconditional precision (e.g., sampling variance, standard errors, confidence intervals) for parameter estimators, rather than measures of sampling uncertainty conditional on a selected model. This research is encouraged to better understand the properties of these approaches.

There are three general approaches to assessing model selection: (1) theoretical studies, mostly using Monte Carlo simulations; (2) bootstrap applied to a given set of data; and (3) utilizing test statistics (i.e., Δ_i) and model weights w_i from the set of models. Insights can be obtained about model selection and associated uncertainty by extensive Monte Carlo simulations of model selection (e.g., Hurvich and Tsai 1998). Use of the bootstrap and $\hat{\pi}_i$ values applied to a single data set; hence they represent our focus here. The bootstrap requires 10,000 samples for reliable results, and it could take many

development and understanding.

Monte Carlo investigations generate 1,000 to 10,000 in (sometimes, 100,000 or even a million samples are needed, ating model. These data sets are then analyzed, the log-likelihood to obtain the MLEs, and model selection is done to identify for each sample. Finally, one can summarize resultant results, models selected and other information of interest, such as as well as conditional and unconditional variances of parameters over models. For results to apply fully to the K-L model envisioned for real data, the generating model (which is the relation study) should be complex and not contained in the set of models, g_1, \dots, g_R (i.e., g_i or $g_i(\underline{x} | \underline{\theta})$). We present some of Monte Carlo simulation results in Chapters 5 and 6.

Many, if not most, simulation studies on model selection are under these conditions (e.g., Wang et al. 1996). Rather, they are because (1) a simple generating model is used (so no model, only small K), (2) the set of models considered contains the true (i.e., contains “truth”), and (3) the model selection goal is to find the generating model (hence to select “truth”). None of these are of real data-analysis problems; hence we discount the results of such studies as appropriate guides to real-world model selection (see 6 for more details on this common error).

The fundamental idea of the model-based sampling approach to statistical inference is that the data arise as a sample from a probability distribution, f , and hence the uncertainties of interest can be measured if we can estimate f . There are ways to construct an estimator of (in essence) f from the sample data. The fundamental bootstrap method (Section 2.13) is that we compute measures of uncertainty from that estimated sampling distribution of \hat{f} .

4.3.1 Concepts of Parameter Estimation and Model Selection Uncertainty

Statistical science should emphasize estimation of parameters and measures of estimator uncertainty. Given a correct model (assumes $g = f$), an MLE is reliable, and we can compute measures of its sampling variance and a reliable confidence interval (e.g., likelihood interval; see Royall 1997). If the model is selected independently of the data at hand, and is a good approximation of the true, large, then the estimated sampling variance is essentially correct, and an appropriate confidence interval will essentially achieve its nominal coverage. This would be the case if we used only one model, deciding

here, the selection process is expected to introduce additional sampling uncertainty into any estimated parameter; hence sampling variances are too small: They are conditional on the model selected, and do not reflect model selection uncertainty. One result is that confidence intervals can be expected to have less than nominal coverage.

Consider a scalar parameter θ , which may be used in all models considered, but is in the selected model, and therefore θ_i given model g_i . Here, the subscript i denotes the model selected, with the understanding that this parameter means the same θ in which it appears. There is a conceptual true value of θ . However, the value of θ that we would infer, in the sense of (for large sample size) from model g_i applied to the data, may differ from the true value. Given model g_i , the MLE, $\hat{\theta}_i$, has a conditional sampling variance $\text{var}(\hat{\theta}_i | g_i)$ and hence a conditional sampling variance $\text{var}(\hat{\theta}_i | g_i)$. We use the notation $\text{var}(\hat{\theta}_i | g_i)$ to be functionally and numerically identical to $\text{var}(\hat{\theta}_i | g_i)$; notation is more traditional; we use the former notation to emphasize the importance of assuming the model in its variance calculation; other models are also being considered.

There is a concept of the true value of θ : It is the value that we would compute based on knowing truth, f , even though θ need not be in f . To the extent a model, g_i , is wrong (i.e., $g_i \neq f$), the value of θ when this K-L best value is determined for θ under assumption f is, even when data are generated by f , if those data are analyzed under model g_i , we will infer (for large n) that the value of θ inferred from model g_i this possible “bias” (i.e., $\theta_i - \theta$) is not of great concern. This bias can be dominated by the conditional sampling standard error of $\hat{\theta}_i$ (this domination is one feature of a “good” model). The value of θ is one source of model selection uncertainty into $\hat{\theta}_i$; that is, the value of θ in models in an unknown manner. In many situations the model selected is something to us, and we will then take $\hat{\theta}_i$ derived only from model g_i , as the most meaningful estimator of θ . This is what is done, and seems sensible, so much so that the alternative seems at first strange, but is an alternative to getting an estimate of θ from multiple models. Model averaging arises in a natural way from the unconditional sampling variance of $\hat{\theta}_i$.

Another problem arising from model selection uncertainty is where the estimate of sampling variance is derived from the fitted model (e.g., multiple linear regression). In this case, the estimate is a negatively biased estimator of conditional sampling variance. Related concepts are required to understand this issue. For example, the structure is Poisson, binomial, multinomial (includes continuous distributions).

(e.g., for the binomial, the model parameter is p and the variance is $\text{var}(\hat{p}) = \frac{p(1-p)}{n}$).

Underestimation of sampling variance due to structural seem to be a serious problem in cases where such a theoretical variance is known. In particular, this is true if we also use a variance estimator \hat{c} , applied for all models in the set of R models (so selection is based on \hat{c}) to adjust for any modest structural lack of fit or misspecification.

The second situation occurs often in regression models. Sampling variance, σ^2 , is functionally unrelated to the true parameter value, so there is no true replication. Then σ^2 must be estimated from the fitted model. In this case there is neither true replication nor a basis to infer σ^2 , such as there is in models for count data. If we have a structural component of the model to the data, we will get a biased estimate of σ^2 , hence estimated sampling standard errors of any $\hat{\theta}_i$ will be biased. This is likely to be a compensating increase in the factor $(\hat{\theta}_i - \theta_i)$ at the level of the regressor values can eliminate this problem, but we do not have such true replication.

The ideas of classical sampling theory can be used to correct for the sampling variance of $\hat{\theta}$ resulting from the two-stage process. In the first stage, selection, then (2) using $\hat{\theta} \equiv \hat{\theta}_i$ given that model g_i was selected. If the process carried out many times, m , each time on an independent sample j we get $\hat{\theta}_j$ as our estimator of θ . This conceptual process is the same as selecting a model in repetition j , but we do not need, hence avoid, the indexed notation (such as $\hat{\theta}_{i,j}$) to denote both sample j and model i given sample j .

The estimated unconditional sampling variance, $\widehat{\text{var}}(\hat{\theta})$, is the variance of $\hat{\theta}$. It would be $\sum(\hat{\theta}_j - \bar{\hat{\theta}})^2/(m-1)$; $\bar{\hat{\theta}}$ is the simple average of the $\hat{\theta}_j$ (hence the $\hat{\theta}_j$ have been averaged over selected models). The estimator represents the total variation in the set of m values of $\hat{\theta}$. Within and between-model variation is included. This sampling variance can be partitioned into R subsets, one for each model which is selected. Each contains all the $\hat{\theta}$'s computed under selected model i . Then from the i th subset of the $\hat{\theta}$ values an estimate of the conditional sampling variance of $\hat{\theta}$ when model g_i was selected. Formal mathematics involves partitioning the above $\widehat{\text{var}}(\hat{\theta})$ into R components and then combining to get a theoretical unconditional sampling variance gives us σ^2 as a weighted combination of conditional variances, variances of $\hat{\theta}_i$ for variation among $\theta_1, \dots, \theta_R$. The weights involved are the selection probabilities. Relevant formulas are given in the next section. I mention one more issue.

ter can be (but need not be) made based on just those models in which the parameter appears. An example is variable selection in linear regression: y versus p regressors, x_1, \dots, x_p (plus an intercept). There are 2^p models, but each regressor appears in only half of these models. Thus if regressor variable x_j , hence parameter β_j , appears in the AIC best model, we could restrict ourselves to just that subset of models that contain β_j in order to directly estimate the unconditional variance of $\hat{\beta}_j$. All the above (and below) considerations about conditional variances with regard to a particular parameter selection apply to just the subset of models that include the parameter.

We have emphasized models as approximations to truth. The “wrong” model is technically called model misspecification (see Chapter 7). The usual theoretical sampling variances of MLEs, $\text{var}(\hat{\theta}_i)$, are conditional on the model, but only trivially so if the model is a good approximation to truth. The theory that gives the correct conditional (on the model) sampling variance of $\hat{\theta}_i$ in the event of model misspecification (Chapter 7 gives the details). However, the correct estimator of $\text{var}(\hat{\theta}_i | g_i)$ is then so messy and variable (a type of instability) that it generally seems preferable to use the theoretical estimator supplied by the usual model-specific theory (which assumes that the model is correct). This simplification is especially defensible when done in conjunction with some model selection procedures intended to minimize both serious overfitting and model selection bias. We believe AIC is suitable for this selection purpose and that the correct consideration is thus to get reliable unconditional sampling variances (and confidence intervals) for MLEs after model selection.

4.3.2 Including Model Selection Uncertainty in Sampling Variance

We continue to assume that the scalar parameter θ is of primary interest and is considered. This will often be the case for our full set of a priori models, and is always the case if our objective is prediction with the model, as interpolation or extrapolation with a generalized linear model. If our focus is on a model structural parameter that appears in many models in our full set of models, then we can restrict ourselves to that subset of models that make the sort of inferences considered here about the parameter. In the latter case we simply consider the relevant subset as the full set of models under consideration.

In repeated (conceptual) samples there is a probability π that the AIC best model is the true model. Presentation of a defensible way to augment $\widehat{\text{var}}(\hat{\theta})$ to include model selection uncertainty involves the idea of model averaging.

$$\theta = \sum_{i=1}^R \pi_i \theta_i, \text{ and its estimator } \hat{\theta} = \sum_{i=1}^R \pi_i \hat{\theta}_i$$

In some theory development we use π_i rather than $\hat{\pi}_i$, but such as $\bar{\theta}$ and $\hat{\hat{\theta}}$.

The theoretical, unconditional sampling variance of $\hat{\theta}$ is given by

$$\text{var}(\hat{\theta}) = \sum_{i=1}^R \pi_i \left[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2 \right]$$

This result follows directly from frequentist sampling theory. See Section 4.3.1 that if we had m independent samples and selected one sample to each sample to get $\hat{\theta}_j$, $j = 1, \dots, m$, then the estimator would be

$$\widehat{\text{var}}(\hat{\theta}) = \sum_{j=1}^m (\hat{\theta}_j - \bar{\hat{\theta}})^2 / (m - 1).$$

Here, j indexes the sample that $\hat{\theta}_j$ came from (whatever the model), whereas i indexes that $\hat{\theta}_i$ arose from model i (whatever the model). This notation allows us to focus on different aspects of the model selection process without a notation so complex that it hinders understanding. As m become infinite, the above estimator of $\text{var}(\hat{\theta})$ converges to the unconditional sampling variance of $\hat{\theta}$. By first grouping the $\hat{\theta}$ values by model and then taking the needed limit as $m \rightarrow \infty$.

Readers less interested in the derivation of an estimator of the model-averaged estimate may want to skip to the following. The term $\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2$ is just the mean square error of $\hat{\theta}_i$; in one sense the unconditional variance of $\hat{\theta}$ is just an average of these errors. Specifically,

$$E[(\hat{\theta}_i - \hat{\hat{\theta}})^2 | g_i] = \text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2$$

and we recommend thinking of the above quantity as the conditional variance of $\hat{\theta}_i$, given model i , when $\hat{\theta}_i$ is being used as an estimator. The incorporation of model selection uncertainty into the variance of $\hat{\theta}$ is a new thinking like this. The matter arises again when we need to compute the covariance, $E[(\hat{\theta}_i - \bar{\hat{\theta}})(\hat{\theta}_j - \bar{\hat{\theta}}) | g_i]$, that also allows for model selection uncertainty.

One might think to estimate the augmented sampling variance as $\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\hat{\theta}})^2$. Such an estimator is not supported by theory and is likely to be both biased (it could be bias-corrected) and inconsistent; however, we have not investigated this possible estimator.

are estimated (they are not the major source of estimation variance estimator), we can evaluate $E(\widehat{\text{var}}(\hat{\theta}))$ to bias-correct $\widehat{\text{var}}(\hat{\theta})$. This correction involves the sampling variance, $\text{var}(\hat{\theta})$, of the model-averaged estimator.

$$E(\widehat{\text{var}}(\hat{\theta})) = \text{var}(\hat{\theta}) + \sum \pi_i \text{var}(\hat{\theta}_i | g_i) - \text{var}(\hat{\theta})$$

which leads to

$$\text{var}(\hat{\theta}) = \text{var}(\hat{\theta}) + \sum \pi_i E(\hat{\theta}_i - \hat{\theta})^2.$$

It seems that we cannot avoid estimating $\hat{\theta}$ and $\text{var}(\hat{\theta})$ even to compute $\widehat{\text{var}}(\hat{\theta})$ that we are seeking. First, note that efforts to evaluate $\widehat{\text{var}}(\hat{\theta})$ are circular, thus useless. Anyway, at worst we would only estimate $\widehat{\text{var}}(\hat{\theta})$ without (much) bias, which we can clearly do. Second, note that as one might expect, if our goal is to estimate θ , then $\hat{\theta}$ is to be preferred to $\hat{\theta}_i$ because it will have a smaller variance. However, given that our goal is to estimate θ , there is no reason to claim that $\hat{\theta}$ is the superior estimator as compared to $\hat{\theta} \equiv \hat{\theta}$.

From Buckland et al. (1997) we will take the needed variance

$$\text{var}(\hat{\theta}) = \left[\sum_{i=1}^R \pi_i \sqrt{\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2} \right]^2$$

with the estimator as

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{i=1}^R \hat{\pi}_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]^2$$

Formula (4.5) entails an assumption of perfect pairwise correlation of $\hat{\theta}_i - \hat{\theta}$ and $\hat{\theta}_h - \hat{\theta}$ for all $i \neq h$ (both i and h index model-averaged estimators). This assumption is unlikely; however, it will be highly conservative. The value of $\rho_{ih} = 1$ is conservative in that $\text{var}(\hat{\theta})$ computed from (4.5) is likely to be too large if this assumption is in error. Also, by justifying the use of (4.5) into (4.6) a further upward bias to (4.5) results. Thus from (4.5) and (4.6) $\widehat{\text{var}}(\hat{\theta}) = \widehat{\text{var}}(\hat{\theta}) + \sum \pi_i (\hat{\theta}_i - \hat{\theta})^2$ with $\widehat{\text{var}}(\hat{\theta})$ from (4.5). This introduces a much positive bias. Hence, we are now suggesting just using (4.6) from (4.6).

All simulations we have done so far, in various contexts, have shown the bias of this estimator:

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{i=1}^R \hat{\pi}_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]^2$$

(4.6) and (4.7) is acceptable. Improved estimation of the variances under model selection may be possible. However, to give practical solutions to some problems under model selection, the expectation that improvements will be further explored.

The $\hat{\pi}_i$ in (4.7) (and the equivalent (4.6)) will usually be different from the weights, w_i . In general, $w_i \neq \pi_i$; rather, w_i can be considered as an estimate of π_i , but (4.7) seems robust to slightly imprecise values of π_i . Alternatively, one can use the bootstrap estimates, $\hat{\pi}_i = b_i/B$; if one has bootstrap samples, the analytical formulas above can be used.

As a final part of this section we give some details of the derivation of (4.5) in a more restricted context than was used in Buckman (2009). Specifically, we do not assume that the R models are random; rather, we just condition on the set of R models. This is reasonable; provided; hence, inferences are conditional on just this set of models. Each $\hat{\theta}_i$ is considered as an estimator of $\bar{\theta}$, and it is this conditional expectation that is critical to getting a variance formula that includes model selection.

Ignoring that the π_i in (4.5) need to be estimated, the variance can be expressed as

$$\begin{aligned} \text{var}(\hat{\theta}) &= \sum_{i=1}^R (\pi_i)^2 \text{E}[(\hat{\theta}_i - \bar{\theta})^2 | g_i] \\ &\quad + \sum_{i=1}^R \sum_{h \neq i}^R \pi_i \pi_h \left[\text{E}(\hat{\theta}_i - \bar{\theta})(\hat{\theta}_h - \bar{\theta}) | g_i, g_h \right] \end{aligned}$$

From above we know that

$$\text{E}[(\hat{\theta}_i - \bar{\theta})^2 | g_i] = \text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2$$

In order to coherently allow for model selection uncertainty in the definition of a correlation, we must interpret this expression for $\text{var}(\hat{\theta})$ as

$$\text{E}(\hat{\theta}_i - \bar{\theta})(\hat{\theta}_h - \bar{\theta}) | g_i, g_h) = \rho_{ih} \sqrt{\text{E}[(\hat{\theta}_i - \bar{\theta})^2 | g_i] \text{E}[(\hat{\theta}_h - \bar{\theta})^2 | g_h]}$$

hence,

$$\begin{aligned} &\text{E}(\hat{\theta}_i - \bar{\theta})(\hat{\theta}_h - \bar{\theta}) | g_i, g_h) \\ &= \rho_{ih} \sqrt{[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2][\text{var}(\hat{\theta}_h | g_h) + (\theta_h - \bar{\theta})^2]} \end{aligned}$$

$$\begin{aligned} \text{var}(\hat{\theta}) = & \sum_{i=1}^R (\pi_i) \left[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2 \right] \\ & + \sum_{h \neq i}^R \sum_{i=1}^R \pi_i \pi_h \rho_{ih} \sqrt{[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2][\text{var}(\hat{\theta}_h | g_h) + (\theta_h - \bar{\theta})^2]} \end{aligned}$$

We have no basis to estimate the across-model correlation (other than the bootstrap, but then we do not need theory for that). The expression simplifies if we assume that all $\rho_{ih} = \rho$:

$$\begin{aligned} \text{var}(\hat{\theta}) = & (1 - \rho) \left[\sum_{i=1}^R (\pi_i)^2 [\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2] \right] \\ & + \rho \left[\sum_{i=1}^R \pi_i \sqrt{[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2]} \right]^2 \end{aligned}$$

From (4.8), if we further assume $\rho = 1$, then we get (4.5):

$$\text{var}(\hat{\theta}) = \left[\sum_{i=1}^R \pi_i \sqrt{[\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \bar{\theta})^2]} \right]^2$$

Unconditional Variance Estimator

Then, using the Akaike weights (w_i) instead of the model selection frequencies (π_i) and using estimates instead of parameters, we get a useful result,

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]^2$$

where $\hat{\theta}$ is model-averaged estimate (4.1).

This estimator of the unconditional variance can be computed using the MLE $\hat{\theta}$ from the selected model or for the model averaged

If only a subset of the R models is used, then the w_i must be based on just these models (thus these new weights must be recomputed). If one has the estimated model selection frequencies ($\hat{\pi}_i$), then the estimator

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{i=1}^R \hat{\pi}_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]^2$$

is useful. In either case, $\widehat{\text{se}}(\hat{\theta}) = \sqrt{\widehat{\text{var}}(\hat{\theta})}$.

parameter estimators in a model, or (what is the same) conditional covariance between two estimators as based on model. A formula for the unconditional $\text{var}(\hat{\theta})$ values is

$$\text{var}(\hat{\theta}) = \left[\sum_{i=1}^R \pi_i \sqrt{\text{var}(\hat{\theta}_i | g_i) + (\theta_i - \hat{\theta})^2} \right]$$

we need something useful for the analogous

$$\begin{aligned} & \text{cov}(\hat{\theta}_1, \hat{\theta}_2), \\ & \hat{\theta}_1 = \sum_{i=1}^R \pi_i \hat{\theta}_{1i}, \\ & \hat{\theta}_2 = \sum_{i=1}^R \pi_i \hat{\theta}_{2i}. \end{aligned}$$

For a useful estimation formula the Akaike weight, w_i , w

We propose the formula below for $\widehat{\text{cov}}(\hat{\theta}_1, \hat{\theta}_2)$:

$$\begin{aligned} &= \bar{r}_{1,2} \left[\sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_{1i} | g_i) + (\hat{\theta}_{1i} - \hat{\theta}_1)^2} \right] \left[\sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_{2i} | g_i) + (\hat{\theta}_{2i} - \hat{\theta}_2)^2} \right] \\ &= \bar{r}_{1,2} \widehat{\text{se}}(\hat{\theta}_1) \widehat{\text{se}}(\hat{\theta}_2), \end{aligned}$$

where $\bar{r}_{1,2}$ is

$$\bar{r}_{1,2} = \sum_{i=1}^R w_i \times r_{1,2|i},$$

and $r_{1,2|i}$ is the estimated sampling correlation between model i . This model-conditional sampling correlation can be estimated from the conditional-on-model- i sampling variance-covariance matrix (based on large sample means).

Now if we are considering the difference $d = \hat{\theta}_1 - \hat{\theta}_2$ (Gelman et al., 2001), then

$$\widehat{\text{var}}(d) = A^2 + B^2 - 2\bar{r}_{1,2}AB,$$

where

$$\begin{aligned} A &= \sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_{1i} | g_i) + (\hat{\theta}_{1i} - \hat{\theta}_1)^2} \\ B &= \sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_{2i} | g_i) + (\hat{\theta}_{2i} - \hat{\theta}_2)^2} \end{aligned}$$

work and experience are needed with the formula versus bootstrap formula performance.

4.3.3 Unconditional Confidence Intervals

The matter of a $(1 - \alpha)100\%$ unconditional confidence interval is considered. We have two general approaches: the bootstrap (see, e.g., Efron, 1997), or analytical formulas based on analysis results from the bootstrap set. The analytical approach requires less computing; hence, it is preferred.

The simplest such interval is given by the endpoints $\hat{\theta}_i \pm \widehat{\text{se}}(\hat{\theta}_i)$, where $\widehat{\text{se}}(\hat{\theta}_i) = \sqrt{\widehat{\text{var}}(\hat{\theta}_i)}$. One substitutes the model-averaged estimator for the estimator used. A common form used and recommended for the interval is $\hat{\theta}_i \pm t_{\text{df}, 1-\alpha/2} \widehat{\text{se}}(\hat{\theta}_i | g_i)$. When there is no model selection ignored, it is clear what the degrees of freedom (df) are for the t -distribution here. For (4.7) it is not clear what the degrees of freedom are, but it is, however, that we are focusing on situations where sample sizes are large enough that the normal approximation will be applicable. These confidence intervals are based on the assumption that $\hat{\theta}_i$ has a normal distribution.

We will hazard a suggestion here; it has not been evaluated, but a similar procedure worked in a different context. If for each i we have degrees of freedom df_i for the estimator $\widehat{\text{var}}(\hat{\theta}_i | g_i)$, then for small degrees of freedom one might try using the interval $\hat{\theta}_i \pm z_{1-\alpha/2} \widehat{\text{ase}}(\hat{\theta}_i)$ where the adjusted standard error estimator is

$$\widehat{\text{ase}}(\hat{\theta}_i) = \sum_{i=1}^R \hat{\pi}_i \sqrt{\left(\frac{t_{\text{df}_i, 1-\alpha/2}}{z_{1-\alpha/2}} \right)^2 \widehat{\text{var}}(\hat{\theta}_i | g_i) + \widehat{\text{var}}(\hat{\theta}_i)}$$

In cases where $\hat{\theta}_i \pm z_{1-\alpha/2} \widehat{\text{se}}(\hat{\theta}_i)$ is not justified by a normal distribution (as judged by the conditional distribution of $\hat{\theta}_i$), interval coverage can be based on a transformation of $\hat{\theta}_i$ if a suitable transformation is known. Log and logit transforms are commonly used in the context of general linear models. In fact, in general linear models, the parameter $\underline{\theta}$ will be linked to the likelihood by a set of functions $\underline{\theta} = \underline{W}(\underline{\beta})$. Then it is $\underline{\beta}$ that is directly estimated, and it is $\underline{\theta}$ that has the simple normal-based confidence limits on components. If $\underline{\theta}$ is used. An interval constructed from a component of $\hat{\underline{\beta}}$ and its sampling variance ((4.5) applies) can be back-transformed to the corresponding component of $\underline{\theta}$.

The above methods are justified asymptotically, or if the normal distribution applies to $\hat{\theta}$. However, “asymptotically” means

there is no model selection is the profile likelihood interval (Hsu and Hsu 1999, Sprott 2000). We suggest here an adaptation that widens the likelihood interval to account for model selection.

Let the vector parameter $\underline{\theta}$ be partitioned into the component γ_i and the rest of the parameters, denoted here by $\underline{\gamma}$. Then $\mathcal{L}(\theta_i | \underline{x}, g_i)$ as a function of θ_i (the subscript denotes the model used) can be written as

$$\mathcal{P}\mathcal{L}(\theta_i | \underline{x}, g_i) = \max_{\gamma_i | \theta_i} \left[\mathcal{L}(\theta_i, \underline{\gamma}_i | \underline{x}, g_i) \right]$$

almost always $\mathcal{P}\mathcal{L}(\theta_i | \underline{x}, g_i)$ has to be computed numerically. The profile deviance as

$$\mathcal{PD}(\theta_i) = 2 \left[\mathcal{P}\mathcal{L}(\hat{\theta}_i | \underline{x}, g_i) - \mathcal{P}\mathcal{L}(\theta_i | \underline{x}, g_i) \right]$$

The large sample profile likelihood interval ignoring model selection uncertainty is the set of θ_i that satisfy the condition $\mathcal{PD}(\theta_i) \leq \chi^2_{1,1-\alpha}$. This interval is approximately a $(1 - \alpha)100\%$ confidence interval.

We propose an interval that is a version of (4.9) adjusted for model selection uncertainty: the set of all θ_i that satisfy

$$\mathcal{PD}(\theta_i) \leq \left[\frac{\widehat{\text{var}}(\hat{\theta})}{\widehat{\text{var}}(\hat{\theta}_i | g_i)} \right] \chi^2_{1,1-\alpha}.$$

It suffices to solve (numerically) this inequality for the endpoints, $\hat{\theta}_{i,L}$ and $\hat{\theta}_{i,U}$. In the event that we are not doing model selection, it seems logical to use the resultant confidence interval from the full set of models.

All of the above was assuming that the parameter of interest θ occurs in all of the R models in the full set of models. Often this will not be the case. There will be a subset of size $Q < R$ of the models in which θ occurs. Conceptually, the parameter θ does not occur in the other $R - Q$ models, even as a value equal to zero. In this event we suggest applying the above theory to just that subset of Q models. The $R - Q$ models does not appear seem totally uninformative about the value of θ . The situation is similar to model selection as in linear all subset selection because there we can assume that every structural parameter is in each model, but some are zero. $\theta \equiv 0$ in a model.

In the case that $Q = 1$ (θ is unique to one model in the full set) none of the above results can be used. In this case it seems there is not a direct way to include model selection uncertainty about the value of θ . An approach we can envision here is to use the conditional sampling variance estimator, $\widehat{\text{var}}(\hat{\theta} | g_i)$, by some

$$\frac{\widehat{\text{var}}(\hat{\theta})}{\widehat{\text{var}}(\hat{\theta}_i | g_i)}$$

and found that it can vary greatly by parameter. Thus, estimating an inflation factor for $\hat{\theta}$ based on a different parameter in different models (from the one that model θ appears in) seems problematic. Fundamentally, it is not clear that we should inflate the conditional variance of a parameter unique to just one model in the selected model. At all we can, and should, do in this case is note the uncertainty in the model that model is likely to be the K-L best model in the full set, and use the model-specific conditional sampling variance for that model. The intervals for θ are then constructed based on just the one model that appears (e.g., profile likelihood interval, or other parameter-specific interval) $\hat{\theta} \pm z_{1-\alpha/2} \widehat{\text{se}}(\hat{\theta}|g)$.

Bootstrap construction of an unconditional confidence interval for a parameter in the selected model is not fundamentally different from bootstrap-based interval construction without model selection. This has been much discussed in the statistical literature (see, e.g., Efron and Tibshirani 1993, Mooney and Duval 1993, Hjorth 1994).

First one generates a large number, B , of bootstrap samples, and applies model selection to each bootstrap sample. A model from the original set are fit to each bootstrap sample, and one of the models is selected as best. Only the estimated parameters from that model are kept in an output set of parameter estimates for each bootstrap sample (the index of the selected model for each sample). Hence, for a parameter not in the selected model g_i for bootstrap sample b , the variance is zero for bootstrap sample b (the subscript b and $\hat{\theta}_b^*$ denote the parameter estimate for bootstrap sample; the model used to get this $\hat{\theta}^*$ varies over bootstrap samples). For any parameter in common over all models, there will be non-zero values in the output data set. In either case the variation in the output values of $\hat{\theta}_b^*$, $b = 1, \dots, m$ ($\leq B$), reflects both model selection and within-model sampling variation of $\hat{\theta}$ given a model.

As noted in Section 4.2.2, the model selection frequency can be estimated from these bootstrap results; but our focus here is on unconditional uncertainty and confidence intervals (Efron and Tibshirani 1993). The average of all m values of $\hat{\theta}_b^*$, $\bar{\hat{\theta}}^*$ is an estimator of the model-average parameter. Hence the empirical variance of the set of m values of $\hat{\theta}_b^*$

$$\text{var}(\hat{\theta}^*) = \sum (\hat{\theta}_b^* - \bar{\hat{\theta}}^*)^2 / (m - 1),$$

$\text{var}(\hat{\theta}) = \text{var}(\hat{\theta}^*)$, the simplest confidence interval is $\hat{\theta} \pm 2\sqrt{\text{var}(\hat{\theta})}$, such an interval fails to make full use of the value of the finding upper and lower interval estimates that allow for a distribution for $\hat{\theta}$ under model selection.

The simple, direct bootstrap-based confidence interval (Efron and Tibshirani 1993). Order the bootstrap values from smallest to largest and denote these ordered values by $\hat{\theta}_{(b)}^*$, $b = 1, \dots, B$. A $(1 - \alpha)100\%$ confidence interval select the $\alpha/2$ lower and $(1 - \alpha/2)$ upper percentiles of these ordered bootstrap estimates as $\hat{\theta}_L$ and $\hat{\theta}_U$. If m may not occur at integer values of b , but if m is large, it suffices to use $\hat{\theta}_L = \hat{\theta}_{(l)}^*$ and $\hat{\theta}_U = \hat{\theta}_{(u)}^*$, where $l = [m \cdot \frac{\alpha}{2}]$ and $u = [m \cdot (1 - \frac{\alpha}{2})]$. Possibly better, unconditional intervals after model selection bootstrap are considered by Shao (1996) for regression probabilities. This can be done with model averaging in a straightforward way. Let $\hat{\beta}$ or $\tilde{\beta}$ and model averaging should be done for each bootstrap sample. One obtains $\hat{\beta}^*$ or $\tilde{\beta}^*$ for each of the B bootstrap samples. Then the standard error and confidence intervals from these results.

Note that B needs to be at least several hundred for the bootstrap to begin to work well, and we recommend 10,000 (and at least 1,000 if the parameter of interest is in every model, then $m = B$, with m random). If the parameter is truly not in every model, m is random. m can be made larger to ensure that a sufficient sample size, m , of bootstrap samples is obtained.

4.4 Estimating the Relative Importance of Variables

Data analysis is sometimes focused on the variables to include in the selected model (e.g., important vs. unimportant). Variables are the focus of model selection for linear or logistic regression. If an investigator uses stepwise analysis to arrive at a final model, a conclusion is drawn that the variables in this model are important and other variables are not important. While common, this is problematic among other issues, fails to fully consider model selection. In this section we provide simple methods to quantify the evidence for the importance of a variable in the set.

Consider 10 models based on combinations of a number of variables. Assume that the selected best model includes x_1 with a weight of only 0.3. There is considerable model selection uncertainty, and hence there would seem to be only weak evidence for the importance of variable x_1 based on the selected best model. But one must

that situation would suggest that x_1 is a very important predictor. One of this importance is to sum the Akaike weights (or the best subset of models that include variable x_1 . This idea is a to model selection whenever it is equated to variable selection in nonlinear models of any type.

Consider the hypothetical example of three regressors, search for the best of the eight possible models of the simple linear type: $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \beta_3x_3 + \epsilon$. The possible combinations of x_1, x_2, x_3 that define the eight possible models are shown below, along with their Akaike weights w_i (a 1 denotes that x_i is in the model; other

x_1	x_2	x_3	w_i
0	0	0	0.00
1	0	0	0.10
0	1	0	0.01
0	0	1	0.05
1	1	0	0.04
1	0	1	0.50
0	1	1	0.15
1	1	1	0.15.

While the selected best model has weight of only 0.50, the sum of 0.5 of being the actual K-L best model here, the sum of Akaike weights for the model containing x_1 is 0.79. This is evidence of the importance of x_1 in the models considered. Variable x_2 was not included in the best model, but this should not suggest that it is of zero importance. The sum of the weight of evidence support is 0.35. Finally, the sum of the Akaike weights for the model containing x_3 is 0.85. Thus the evidence for the importance of x_3 is substantially more than just the weight of evidence support. We can order the three predictor variables in this example by their importance: x_3, x_1, x_2 with importance weights of 0.85, 0.79, and 0.15. With other methods recommended here, we see that we can make a more informed selection to go well beyond just noting the best model from the set.

Relative Variable Importance

Estimates of the relative importance of predictor variables can be made by summing the Akaike weights across all the models in which a variable j occurs. Thus, the relative importance of variable j is the sum $w_+(j)$.

The larger the $w_+(j)$ the more important variable j is relative to other variables. Using the $w_+(j)$, all the variables can be ranked by importance.

The direction and magnitude of effect size should often be estimated by averaging estimates with appropriate measures of precision.

of an models that include the pair of variables. For the pair of evidence for the importance of this pair is 0.19. For pair of evidence for importance is 0.23, while for the pair x_1 evidence is 0.65 (compared to 0.5 for the selected model procedures apply when assessing the relative importance

When assessing the relative importance of variables us it is important to achieve a balance in the number of model variable j . For example, in the numerical example above variables appeared in four models. This balancing puts each footing.

To summarize, in many contexts the AIC selected best some variables and exclude others. Yet this inclusion or does not distinguish differential evidence for the importance the model. The model weights, w_i or $\hat{\pi}_i$, summed over all a given variable provide a better weight of evidence for the variable in the context of the set of models considered.

4.5 Confidence Set for the K-L Best Model

4.5.1 Introduction

There exists a concept of a confidence set for the K-L best data, just as there is a confidence interval for a parameter b data. For a 95% confidence set on the actual K-L best model (not unique) approach is to sum the Akaike weights from until that sum is just ≥ 0.95 ; the corresponding subset of confidence set on the K-L best model. In this example (as indexed the models as 1 to 7 in order of decreasing weight set is models $\{1, 2, 3, 4, 5\}$, which has sum of weights = approach to a confidence set of models we are interpreting as a posterior probability (i.e., given the data and the set that model i is the K-L best model (see Section 6.4.5). approach but it is easy to understand.

There is another approach to developing a confidence on the idea of a Δ_i being a random variable with a sa In particular, let index value *best* correspond to the actual model in the set. There is always a K-L best model in the set that ties might occur). It is thus model g_{best} that we should analysis; we just do not happen to know a priori the value of conceptual interest is

$$\Delta_p = \text{AIC}_{best} - \text{AIC}_{min} .$$

confidence interval on θ . A pivotal quantity is one whose distribution is independent of any unknown parameters, a t -distributed variable. The “ p ” in the Δ defined by (4.12) denotes that this Δ is a p -value rather than an actual Δ_i that we can compute from the data.

It is not exact to consider $\Delta_p = \text{AIC}_{\text{best}} - \text{AIC}_{\text{min}}$ as a p -value, but it seems a useful approximation in some contexts. The approximation is useful in is one of complex truth; tapering effect sizes; being good approximations to truth, with full truth not known; and being used; and a lot of nested sequences of models (as in the example in Chapter 3). Monte Carlo studies on the above approximation we have done many of these and results support the conclusion that in this context, the sampling distribution of this Δ_p has substantial overlap with the 95th percentile of the sampling distribution of Δ_p is generally less than 10, and in fact generally less than 7 (often closer to 4 in simple situations). This means that an alternative rule of thumb for an approximate confidence set on the K-L best model is the subset of all models g_i for which Δ_i value that is roughly in the range 4 to 7. In fact, the Δ_i for the K-L best model is not competitive as a candidate for the K-L best model. Δ_i is probably somewhere between 2 and 10 in many situations. Δ_i is not large, while a $\Delta_i = 10$ is strong evidence against g_i being the K-L best model in the set of models considered, if sample size is large. These guidelines, rough as they are, are useful.

We review this interpretation of evidence from the Δ_i in Chapter 5. For independent, sample sizes are large, and models are nested,

Δ_i	Level of Empirical Support
0–2	Substantial
4–7	Considerably less
> 10	Essentially none

Models with $\Delta > 10$ represent very strong evidence that g_i is not the K-L best model. The reader should not take these guidelines too literally, as there are situations to which they do not apply well (such as small sample size or dependent observations). Likewise, in the absence of models, these guidelines may not hold.

We had these guidelines well in mind when we encountered the guidelines for the Bayes factor. The Bayes factor is a Bayesian measure of the relative data-based likelihood of one model versus another, without considering any priors on the set of models (Berger and Raftery 1996a); it is somewhat analogous to $\exp(-\frac{1}{2}\Delta_i)$. Raftery (1996b:165) presents a similar scale for interpretation of Δ_i as evidence for the simpler model being considered.

A third reasonable basis for a confidence set on model selection is likelihood-based inference (see e.g., Edwards 1992, Az

of model i versus the estimated K-L best model g_{best} as a measure of relatively more plausible models. Thus our confidence set consists of models for which the ratio

$$\frac{\mathcal{L}(g_i|x)}{\mathcal{L}(g_{min}|x)} > \text{cutoff},$$

where the cutoff value might be $\approx 1/8$. Models where the ratio is greater than $1/8$ are in the confidence set and are deemed plausible. No direct sampling theory interpretation is required and we return to the idea of the selected subset of models including g_{best} with a preset, known, long-run inclusion relative frequency. This procedure has the advantage that the cutoff remains unchanged by the addition or deletion of a model (of course, a new model will be added to the confidence set). Thus, a confidence set based on the relative likelihood has a desirable invariance property. In contrast, any change in the data can alter the confidence set when summing the Akaike weights.

We have presented three approaches to finding a confidence set of models: (1) base it directly on the Akaike weights, interpret the relative probabilities of each model being the actual best model; (2) use a cutoff Δ_i motivated by the idea of the sampling distribution of the approximate pivotal, Δ_p (using, say, the 95th percentile of the distribution as the cutoff Δ); or (3) think in terms of relative likelihood. Models with $\mathcal{L}(g_i|x)/\mathcal{L}(g_{min}|x) \equiv \exp(-\frac{1}{2}\Delta_i)$ is small, say 0.135 ($\Delta_i = 1$) or 0.050 ($\Delta_i = 6$). In general we favor this third approach.

The use of intervals based purely on relative likelihood is not justified by statistical theory (cf. Berger and Wolpert 1984, Edwards 1996, Royall 1997), but rarely taught or used. Rather, most people have been taught to think of confidence intervals in terms of coverage probability; hence they might feel more at home with methods of which are motivated by the sampling theory idea of a confidence interval on a parameter. The approach based on simple evidence is quite useful. More needs to be known about the properties of the confidence sets to construct a confidence set of models before we would be recommending just one approach.

4.5.2 Δ_i , Model Selection Probabilities, and the

For a given set of data we can estimate the sampling distribution of the selection frequencies and the distribution of $\Delta_p = \text{AIC}_p$ (see formula 4.12) using the bootstrap method. In this method the (unknown) K-L best model is played by the model selected by the data analysis; denote that model by model g_{best} . For example,

$\Delta_p^* = \text{AIC}_{\text{best}}^* - \text{AIC}_{\min}^*$, Δ_p^* does not change over bootstrap samples. The model producing AIC_{\min}^* varies by bootstrap sample. However, the best model g_{best} is the same for all bootstrap samples. If the best model is the model g_{best} in which case $\Delta_p^* = 0$. When it is not the best model, then $\Delta_p^* > 0$.

The B bootstrap samples provide B values of Δ_p^* that are conditional on the data. The percentiles of the empirical distribution of Δ_p^* provide the estimate of the percentile of the distribution of Δ_p^* , and hence provide a basis for a confidence interval for the best model for the actual data. For a $(1 - \alpha)100\%$ confidence interval, order the $\Delta_{p,(b)}^*$ (smallest to largest) and find the $b = [(1 - \alpha)B]$. For the actual data analysis results, the sum of the g_i having $\Delta_i \leq \Delta_{p,(1-\alpha)B}^*$ is the desired confidence set. For the upper tail percentiles of Δ_p^* , B needs to be 10,000.

Other information can be gained from these bootstrap samples. For example, selection uncertainty, in particular, the frequency of selection of each model. Let b_i be the number of samples in which model g_i is selected. Then an estimator of the relative frequency of selection of model g_i in the given situation is $\hat{\pi}_i = b_i/B$. These estimated selection probabilities are useful for assessing how much sampling variation there is in the selection of the best model: they directly quantify model selection uncertainty. If the estimated selection probabilities are similar to, but not identical to, the Akaike weights, which also quantify strength of evidence for each model, then selection uncertainty is small.

Also, for each bootstrap sample we can compute the Akaike weights

$$w_i^* = \frac{\exp(-\frac{1}{2}\Delta_i^*)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta_r^*)}$$

and then average these over the B samples to get \overline{w}_i^* . Coherence of the \overline{w}_i^* , and $\hat{\pi}_i$ is informative as to the coherence of these measures. The $\hat{\pi}_i$ provides information about the sampling uncertainty in the selection of the best model. The \overline{w}_i^* provides information about the sampling uncertainty in the estimation of the Akaike weights. The $\hat{\pi}_i$ is a theoretical measure of model selection sampling uncertainty. The \overline{w}_i^* is a measure of model selection uncertainty based on unknown selection probabilities, π_1, \dots, π_R . Either the $\hat{\pi}_i$ or the Akaike weights, w_i (we often prefer the latter because they require less computer-intensive calculations and they relate more directly to the strength of evidence based on the data at hand), may be taken as the measure of model selection uncertainty. Note that in this usage of the term "model selection uncertainty" we mean the structural form of the model (such as which variables are included vs. excluded) without consideration of the specific parameter values for each model. Parameter-estimation uncertainty is conceptually distinct from model selection uncertainty (but influenced by) model selection uncertainty.

Consider a set of three models, in which models g_2 and g_3 are 100% redundant in the set of models; the model set of three models we get $\mathcal{L}(g_1)/\mathcal{L}(g_2) = \mathcal{L}(g_1)/\mathcal{L}(g_3) = 7.4$. The redundancy has not affected the Δ_i nor the likelihood evidence. However, the (normalized) Akaike weights (Section 2.9) for the set of two models, $w_1 = 0.881$ and $w_2 = 0.119$; when model redundancy, $w_1 = 0.787$ and $w_2 = w_3 = 0.106$. In the model set we still have $w_1/w_2 = 7.4$ ($= w_1/w_3$); likelihood is affected by model redundancy.

The difference between a w_1 of 0.881 and one of 0.787 is not our point is that this, clearly erroneous, model redundancy calculation has affected the Akaike weights. The weights for the redundant model included are not correct because the weights are twice (one time “disguised” as Δ_3). The effect on the weights here (but it could be), but they are wrong, and this could adversely affect calculations using the w_i (as $\hat{\pi}_i$), as for example in conditional and unconditional variance calculations.

If the model redundancy was recognized, and we want to correct (should not), we could correct the situation by considering the model set having two subsets: Model g_1 is one subset; a second subset is g_2 and g_3 . Given that we know that models g_2 and g_3 are 100% redundant, we allocate prior weights, about which model is the expected model, 1/2 to each subset, and the 1/2 is further divided equally between g_2 and g_3 . Thus, $\tau_1 = 0.5$, $\tau_2 = 0.25$, and $\tau_3 = 0.25$. Now we can

$$w_i = \frac{\mathcal{L}(g_i|\underline{x})\tau_i}{\sum_{r=1}^R \mathcal{L}(g_r|\underline{x})\tau_r}$$

from Section 2.9 to compute correct Akaike weights for the model set, thus

$$w_1 \propto 1.0 \cdot \frac{1}{2}, \quad w_2 \propto 0.135335 \cdot \frac{1}{4}, \quad w_3 \propto 0.135335 \cdot \frac{1}{4}$$

The normalized (to add to 1) weights are 0.8808, 0.0592, and 0.0592. The sum of the weights for models g_2 and g_3 correctly add up to 1.0. This is the correct result, and it serves to introduce the concept, and issue, of model redundancy. It is possible to have actual model redundancy, and it is not careful in constructing the set of models considered.

models structured into two or more subsets. Different subsets can be specified in such a way that they have one key model in common. Then adjustment terms are applied to get a sequence of models. This can be used with different types of adjustment terms. So we can have this situation: The full set of models is given as two subsets, $\{g_1, g_2, g_3, g_4\}$ and $\{g_1, g_5, g_6, g_7\}$. If the full set of models contains 8 different models, then the redundancy of model g_1 is not zero. The situation can easily be rectified if it is recognized: either delete g_1 to 7, or compute the w_i from the Δ_i , for the models labeled g_i , with different priors, τ_i , as $\{\frac{1}{14}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}\}$ and $\{\frac{1}{14}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}\}$. It is clear there might be a partial model redundancy in the models used.

Failure to completely understand the models used can lead to model redundancy. For example, a logistic regression may be used as a model of the probability of success,

$$p(x) = \frac{\exp(a + bx)}{1 + \exp(a + bx)},$$

where the parameters are a and b . However, this model can also be expressed as,

$$p(x) = \frac{1}{1 + \exp[-\{(x - d)/c\}]},$$

where $a \equiv d/c$ and $b = -1/c$. The second model is just a different model representation and a 1-to-1 reparametrization. In forms where included in the model set, total redundancy may be avoided such model redundancy.

To further illustrate model redundancy we consider capture-recapture data, obtained on k capture occasions for a population of size N . The parameters of such models are population size, capture probabilities (denoted by p), by occasion, animal-specific effects. One possible type of model is model g_b under which there are different capture probabilities: for first capture or for recapture. This is the case where animals have a behavioral response to first capture. The number of factors affect capture probability ($K = 3$).

A different model (g_t) allows capture probabilities to vary by occasion; so we have p_1, p_2, \dots, p_k ($K = k + 1$). Thus, we have different models. However, the model under which capture probabilities vary by time allows for many submodels ($2^k - k$ possible submodels, the most general case). Some example (sub) models are

$$g_{t1} : p_1 = p_2, \text{ other } p_i \text{ all differ } (K = k),$$

$$g_{t2} : p_1 = p_2 = p_3, \text{ other } p_i \text{ all differ } (K = k + 1),$$

$$g_{t3} : \text{all } p_i = p \text{ } (K = 1),$$

$$g_{t4} : \text{all } p_i \text{ are different } (K = k + 1).$$

model redundancy becomes irrelevant because g_b is over model: The usual Akaike weight for that model is here claim, however, that the correct weights here should be based as $\{\frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{2}\}$ not $\{\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}\}$. If the Δ_i are $\{2, 0, 0, 0, 0\}$, model redundancy is irrelevant (redundancy is only between others; there is no redundancy in models 1 to 4 if model 5 is the result for the Δ_i is $\{2, 2, 2, 2, 0\}$, then model redundancy deal as regards the proper w_i . For wrong priors $\{\frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}, \frac{1}{5}\}$ under correct priors $\{\frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{8}, \frac{1}{2}\}$, $w_5 = 0.73$.

By adding submodels of the general time-specific model we dilute the absolute strength of evidence for model g_b as we add more models and their weights; and we must use such absolute weights in certain forms of averaging). Inasmuch as these added models deal only with capture probabilities, they are all of a type (hence, redundancy in evidence *against* model g_b), so they unfairly “gang up on g_b ” which is a totally different type of model.

The appropriateness of unequal priors if submodels of model, g_t , are included is justified here on a theoretical basis mentioned in the capture–recapture literature that there is no redundancy as regards estimating N , of considering constrained versus unconstrained time-specific model. Thus, the original set of two models can be augmented as above. Hence, in the last example we should have these two models, and they have Δ_1 and Δ_2 as $\{2, 0\}$. $w_2 = 0.73$. A key point here is that when we did have model redundancy the use of the unequal priors did produce the correct Akaike weights. I think that model redundancy can be coped with analytically by a modification of the otherwise equal model priors, τ_i .

Even more important than accepting model redundancy is accepting model priors, is to construct the set of models to be considered if there is no model redundancy. As the above example illustrates, knowledge about the correct formulation and use of models at hand should be utilized in defining the a priori set of models. Another point worth repeating is that neither the Δ_i nor the weights of the models will be affected by model redundancy. Thus, the models based on all models with Δ_i less than some cutoff value is the safest type to use. Our ideas on the cutoff value to use come from the distribution of Δ_p , but only for situations with no model redundancy. Recommendations already made on this matter were so difficult.

The concept and issue of model redundancy was brought up by S. T. Buckland (personal communication); the above example. Professor Buckland suggested that the bootstrap would be a solution to model redundancy as long as for a given bootstrap sample forced to select one best model. This seems reasonable; but

we must be able to recognize model redundancy in our data. If we can do that (we can and should), we think that redundancy can be useful. If model redundancy operates at a more subtle level than the bootstrap would have an advantage. We are currently discussing this issue, but there will be a model redundancy problem as regards Akaike weights if the set of models considered is carefully constructed. (The bootstrap issue would be helpful.)

4.7 Recommendations

If data analysis relies on model selection, then inferences must account for model selection uncertainty. If the goal is to get the best model, then model selection uncertainty is recommended. If the models have definite, and different, inferences as regards understanding relationships among variables, then the inferences that is sought, then one wants to identify the best model. Hence, reported parameters should be from the selected model (not model averaged values). When selecting a best model, also note the competing models and their Akaike weights. Restricting detailed comparisons to a small confidence set on models should often suffice. If a single model is supported, $w_{min} \geq 0.9$, and competing models give alternative inferences, they should be reported. It may occur that the basic inference is from all good models. However, this is not always the case. Inference based on a single best model may not be sound if support for the best model is weak (in all-subsets selection when $R > 1,000$ or small, e.g., < 0.01).

We recommend that investigators compute and report measures of precision based on (4.9) when inference is based on model selection unless the Akaike weight w_i for the selected model is large. If an unconditional confidence interval, often the form $\hat{\theta} \pm 2\hat{\sigma}$, is an interval of this type back-transformed from a function of the parameter via the link function in general linear models. If such a model has clear deficiencies, or in general if the computation can be avoided, profile likelihood intervals based on formulas (4.10) and (4.11) are recommended.

If interest is really just on some parameters in common to all models, we recommend using model-averaged parameter estimates. The sampling variance estimate to use is then (4.9). Again, often the bootstrap will suffice for a confidence interval.

We think that these analytical procedures can suffice, but sometimes it is better to use the bootstrap to evaluate model selection uncertainty.

The bootstrap provides direct, robust estimates of model selection probabilities rather than use of the Akaike weights, superior unconditional sampling variances or model-averaging estimators. The primary purpose of the bootstrap is to assess inferences; therefore, we recommend that the point estimates be MLEs from the selected model (not the bootstrap means). The problem is very complex, where there may be no suitable analytical expressions for factors of conditional (on model) sampling variances, the bootstrap is used to get conditional and unconditional measures of precision. It is recommended that more bootstrap samples be used than is commonly recommended, but for really reliable results, but even 400 would be better than 100. The model selection uncertainty (no assessment has often been made).

Be mindful of possible model redundancy. A careful selection of a priori models should eliminate model redundancy prior to the bootstrap. A central part of a sound strategy for obtaining reliable inferences is to use statistical tests of post hoc null hypotheses, the associated p-values, and decisions concerning supposed “significance” are not valid.

The theory here applies if the set of models is a priori to the bootstrap. If any models considered have been included after some analysis, then any model(s) are suggested by the data, then theoretical results (e.g., formulas) might fail to properly apply (in principle, the bootstrap is not used). Even for such data-driven model selection strategies, it is recommended to assess model selection uncertainty rather than ignoring it.

4.8 Cement Data

We return to the cement data of Section 3.2 to compare bootstrap results with model selection frequencies (π_i), Δ_i values, Akaike weights, and conditional estimation of sampling variances. These quantities are shown in Table 4.1; the AIC_c -selected model is shown there in bold. The seven models are not shown in Table 4.1 because they were not selected in the 10,000 bootstrap samples (also, they have virtually zero weight). The three simple approaches shown in Table 4.1 provide estimates of model selection uncertainty for this very small ($n = 13$) data set. Model {12} is indicated as the best by all approaches. However, the model selection uncertainty is evident because that best model has a weight of only 0.57 and a bootstrap selection probability of 0.57. These approaches cast substantial doubt concerning the utility of the results for models in Table 4.1. Model {34} is particularly unsupported, with $w_i = 0.0004$.

Δ_i and derived Akaike weights computed from the data.

Model	K	$\hat{\pi}_i$	Δ_i	w_i
{12}	4	0.5338	0.0000	0.5670
{124}	5	0.0124	3.1368	0.1182
{123}	5	0.1120	3.1720	0.1161
{14}	4	0.2140	3.3318	0.1072
{134}	5	0.0136	3.8897	0.0811
{234}	5	0.0766	8.7440	0.0072
{1234}	6	0.0337	10.5301	0.0029
{34}	4	0.0039	14.4465	0.0004

Evidence for the importance of each variable can be obtained by bootstrap and tallying the percentage of times that each variable was in the AIC_c selected model (Section 4.4). For the 10,000 bootstrap models, x_1 occurred in 93% of the models, followed by x_2 (76%), x_3 (71%), and x_4 . Again, this simple approach indicates the importance of x_1 , x_2 , and x_4 . Similar evidence can be obtained by summing the Akaike weights of those models with a particular variable present. Using the Akaike weights, the relative support of the four variables is as follows: x_1 (93%), x_2 (76%), x_3 (71%), and x_4 (32%). Considering the small sample size, the bootstrap and Akaike weights seem to give similar results.

Using the idea of the pivotal Δ_p (Section 4.5) to obtain a confidence distribution as Δ_p^* , we find that an approximate 90% confidence set is achieved if $\Delta_i < 8.75$, while a 95% set is achieved if $\Delta_i < 13.8$. The percentile values of Δ_i are quite extreme here because the sample size example is so small ($n = 13$).

Using the bootstrap selection frequencies ($\hat{\pi}_i$), models {12}, {124}, {123}, and {14} represent an approximate 86% confidence set, while adding models {134} and {234} reflects an approximate 94% confidence set of models. Using the Akaike weights (w_i), an approximate 90% confidence set includes model {12}, {124}, and {14}. The Δ_i values suggest that the final three models have little utility. These types of ranking and calibration measures are not available under a hypothesis testing approach or cross-validation.

We now illustrate the computation of unconditional estimates of Y_0 first for a parameter in common to all models. What if one is interested in the value $\hat{E}(Y_0)$, denoted for simplicity by \hat{Y}_0 , given the values of $x_3 = 10$, and $x_4 = 20$ (cf. Table 3.1)? The prediction unconditional on the models of Table 4.1 is shown in Table 4.2; we used PROC GENMOD in SAS to easily compute predicted values and their conditional standard errors ($\widehat{\text{se}}(\hat{Y}_0 | g_i)$). Clearly, \hat{Y}_0 is high for model {234}, relative to the other models. The estimated standard error for model {1234} is very high, and

model measures of precision are given for Y_0 ; \hat{Y} denotes a model-averaged prediction and $\hat{Y}_0 - \hat{\bar{Y}}$ is the estimated bias in using a given model to estimate

Model	K	\hat{Y}_0	$\widehat{\text{se}}(\hat{Y}_0 g_i)$	$\widehat{\text{var}}(\hat{Y}_0 g_i)$	$(\hat{Y}_0 - \hat{\bar{Y}}_{\text{boot}})$
{12}	4	100.4	0.732	0.536	4.26
{124}	5	102.2	1.539	2.368	0.07
{123}	5	100.5	0.709	0.503	3.86
{14}	4	105.2	0.923	0.852	7.48
{134}	5	105.2	0.802	0.643	7.48
{234}	5	111.9	2.220	4.928	89.01
{1234}	6	101.6	5.291	27.995	0.74
{34}	4	104.8	1.404	1.971	5.45

because the X matrix is nearly singular. Both of these models have little support, as reflected by the small relative weights, so the prediction under these fitted models is of little credibility.

The predicted value for the AIC_c -selected model is 100.4, with a conditional standard error of 0.73. However, this measure tends to underestimate because the variance component due to model uncertainty has not been incorporated. Model averaging (4.2) results in a value of 102.5 using the bootstrap estimated weights ($\hat{\pi}_i$) or the Akaike weights (w_i) (4.1). The corresponding estimated conditional standard errors are 3.0 using the bootstrap-based weights and 2.4 using the Akaike weights. These unconditional standard errors are substantially larger than the conditional standard error of 0.73. In Monte Carlo studies, we find that the unconditional standard errors better reflect the true conditional standard errors, and conditional confidence intervals cover near the nominal level (Chapter 5).

Study of the final three columns in Table 4.2 above shows that the variance in the model-specific predictions (i.e., the \hat{Y}_0) from the various models is large. The conditional variance $(\hat{Y}_0 - \hat{\bar{Y}}_{\text{bootstrap}})^2$ or $(\hat{Y}_0 - \hat{\bar{Y}})^2$ is substantial relative to the conditional variance $\widehat{\text{var}}(\hat{Y}_0 | g_i)$. Models {1234} and {34} are exceptions because they overfit the data (i.e., more parameters than observations). The Akaike weights are relatively easy to compute compared to the bootstrap estimates $\hat{\pi}_i$; w_i seem preferable in this example. That is, we perceive no advantage here for the bootstrap-based results, compared to the Akaike-weight-based results. This is true for the computational cost of the bootstrap (we do not claim that bootstrap-based results are any worse, just not better).

The investigator has the choice as to whether to use the predicted value from the AIC_c -selected model (100.4) or a model-averaged prediction

errors (5.6 for the bootstrap and 1.9 for Akaike weights), of weights makes no great difference. However, there is uncertainty associated with this data set, and we would model-averaged predictions (when prediction is the objective) using Akaike weights. Thus, we would use 101.6 as the predicted value, with an unconditional standard error of 1.9. If the AIC_c -selected model is strongly supported by the data, then we might suggest using the estimate based on that (best) model (i.e., $\hat{Y}_0 = 100.4$) combined with the standard error of the unconditional standard error (1.9), based on the Akaike weights.

The selected model includes only regressor variables X_1 and X_2 . In the model the estimated partial regression coefficients and their standard errors are $\hat{\beta}_1 = 1.4683$ (conditional $\widehat{se} = 0.1213$) and $\hat{\beta}_2 = 0.0122$ (conditional $\widehat{se} = 0.0459$). Each of these parameters appears in the model. To compute the estimate of unconditional sampling variance of $\hat{\beta}_1$, we find each model containing β_1 , its estimate and conditional standard error, and the model's Akaike weight:

Model	$\hat{\beta}_1$	$\widehat{se}(\hat{\beta}_1 g_i)$	w_i
{12}	1.4683	0.1213	0.5670
{124}	1.4519	0.1170	0.1182
{123}	1.6959	0.2046	0.1161
{14}	1.4400	0.1384	0.1072
{134}	1.0519	0.2237	0.0811
{1234}	1.5511	0.7448	0.0029
{1}	1.8687	0.5264	0.0000
{13}	2.3125	0.9598	0.0000

The first step is to renormalize the w_i so they sum to 1 across models. Here that sum is 0.9925 before renormalizing, so we divide the renormalized w_i , but they are the weights to use in applying the model-averaged estimate of β_1 is 1.4561 (from 4.1).

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{i=1}^8 w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]$$

For example, the first term in the needed sum is $0.5670 \sqrt{(0.1213)^2 + (0.0122)^2}$. Completing the calculation gives 0.01755 , or an estimated unconditional standard error of 0.01755 compared to the conditional standard error given the selected model of 0.1213.

{12}	0.6623	0.0459	0.6988
{124}	0.4161	0.1856	0.1457
{123}	0.6569	0.0442	0.1431
{1234}	0.5102	0.7238	0.0035
{234}	−0.9234	0.2619	0.0089
{23}	0.7313	0.1207	0.0000
{24}	0.3109	0.7486	0.0000
{2}	0.7891	0.1684	0.0000

When all 16 models are considered, the Akaike weights of the models above add to 0.8114. However, to compute resampling errors for these eight models we must renormalize the relevant Akaike weights to 1. Those renormalization Akaike weights are what are used in the model-averaged estimator of β_2 is 0.6110, and the unconditional standard error of $\hat{\beta}_2$ is 0.1206 (compared to the conditional standard error of 0.0459). It is important here to compute and use unconditional standard errors in all inferences after data-based model selection. The (conservative) confidence intervals on β_1 and β_2 , using resampling errors, should be constructed based on a t -statistic with 10 df (t_{10} for a two-sided 95% confidence interval). Such intervals here are well away from 0; for example, the 95% confidence interval for β_2 is 0.93.

We generated 10,000 bootstrap samples of these data. For each selection to all 16 models fit to each bootstrap sample. There is a common (and common belief about the bootstrap) it should be accurate estimate of the unconditional standard error of an estimated partial regression coefficient based on the standard deviation of the set of realized estimates of β_2 over all bootstrap samples, b , wherein the selected model is x_1 . The results are given below, along with the average value of β_2 estimate over all relevant bootstrap samples:

parameter	bootstrap results	
	average	st. error
β_1	1.461	0.760
β_2	0.453	0.958
β_3	−0.420	1.750
β_4	−0.875	1.237

From the selected model, {12}, we get $\hat{\beta}_1 = 1.47$ and $\hat{\beta}_2 = 0.45$ with unconditional standard errors of 0.18 and 0.12, respectively. These results are in good agreement with analytical methods using Akaike weights.

Based on the above, and other comparisons not given, we conclude that the bootstrap failed here when all 16 models were allowed to be considered. In contrast, noted in Section 3.2, the full design matrix, X , for this data set

eigenvalues sum to 95.5% of the total. Also, the pair-wise correlation of x_2 and x_4 is $r = -0.973$. This information, to us, strongly justified one to drop model {1234} from consideration and to drop models in which both x_2 and x_4 appear. Thus without any model selection, should, reduce the 16 possible models to 12 by eliminating {234}, and {1234}. These sorts of considerations should be done, but do not compromise an a priori (as opposed to exploratory) strategy.

With the reduced set of 12 models we computed the Δ_i and 10,000 new bootstrap samples (to get $\hat{\pi}_i$ and bootstrap estimates of the results below (models not shown were never selected in the bootstrap samples):

Model	K	$\hat{\pi}_i$	Δ_i	w_i
{12}	4	0.5804	0.0000	0.6504
{123}	5	0.1315	3.1720	0.1332
{14}	4	0.2340	3.3318	0.1229
{134}	5	0.0465	3.8897	0.0930
{34}	4	0.0076	14.4465	0.0005

Applying here the method of Section 4.4 based on the bootstrap estimates of the weights, we get the relative importance for the four variables: x_1 (0.9995), x_2 (0.7836), x_3 (0.2267), and x_4 (0.2164). Using the method of Section 4.3.2, especially (4.7) with the above Akaike weights, we get the unconditional standard errors as $\widehat{\text{se}}(\hat{\beta}_1) = 0.18$ (for $\hat{\beta}_1 = 0.12$) and $\widehat{\text{se}}(\hat{\beta}_2) = 0.046$ (for $\hat{\beta}_2 = 0.66$, $\widehat{\text{se}}(\hat{\beta}_2 | g_{\{12\}}) = 0.12$). The bootstrap estimates of unconditional standard errors are 0.34 and 0.046 for $\hat{\beta}_1$ and $\hat{\beta}_2$, respectively.

The two different methods (analytical Akaike weights and bootstrap) agree for $\widehat{\text{se}}(\hat{\beta}_2)$ but not for $\widehat{\text{se}}(\hat{\beta}_1)$. The resolution of this discrepancy is on two items. First, the correlation in the data of x_1 and x_2 is high. Second, the sample size is only $n = 13$. As a result, fitted values for {134} are very unstable over bootstrap samples as regards $\hat{\beta}_1$. For example, the sampling standard deviation (this estimate is based on the 1,315 bootstrap values of $\hat{\beta}_{1,b}$ that resulted when model {134} was selected by AIC_c was 0.65 (the average of the bootstrap estimates of the theory-based estimate is $\widehat{\text{se}}(\hat{\beta}_1 | g_{\{12\}})) = 0.12$.

There are several points we wish to make with this example. **sensitive to having demonstrably poor models in the selection process; thus it is very important to exclude models that are demonstrably poor.** The analytical method (vs. the bootstrap method) of assessing the standard errors seems more stable, as regards having or excluding models that are demonstrably poor.

seemed reasonable even with all 18 models considered (unprecise when only the 12 models were used). With the results of the bootstrap results are still suspect, but now only because of the small ($n = 13$). Monte Carlo evaluation and comparison of the results are needed before definitive statements about reliability will be possible.

4.9 Pine Wood Data

We consider here an example of only two simple linear models, neither one nested in the other. This example has been used by Carlin and Chib (1995) on Bayesian model choice using Markov chain Monte Carlo. The data also appear elsewhere, such as in Efron (1984). The response vector (y, x, z) can be considered a trivariate response vector (y, x, z), $n = 42$. Variable y is the measured strength of a piece of wood, x is the measured density of that wood, and z is the measured density for the measured resin content of the wood. The scientific question is whether x or z is a better predictor of the wood strength y , based on either $y = a + bx + \epsilon$ or $y = c + dz + \delta$ (ϵ or δ are normally distributed from the expected linear model structure). Residuals are normally distributed and homogeneous under either model. Scientists believe that wood density adjusted for resin content should be a better predictor of wood strength, but it takes more time and cost to measure z than x .

TABLE 4.3. Pine wood strength data y , wood density x , and wood resin content z (from Carlin and Chib (1995)); $n = 42$.

y	x	z	y	x	z	y
3040	29.2	25.4	2250	27.5	23.8	1670
2470	24.7	22.2	2650	25.6	25.3	3310
3610	32.3	32.2	4970	34.5	34.2	3450
3480	31.3	31.0	2620	26.2	25.7	3600
3810	31.5	30.9	2900	26.7	26.4	2850
2330	24.5	23.9	1670	21.1	20.0	1590
1800	19.9	19.2	2540	24.1	23.9	3770
3110	27.3	27.2	3840	30.7	30.7	3850
3160	27.1	26.3	3800	32.7	32.6	2480
2310	24.0	23.9	4600	32.6	32.5	3570
4360	33.8	33.2	1900	22.1	20.8	2620
1880	21.5	21.0	2530	25.3	23.1	1890
3670	32.2	29.0	2920	30.8	29.8	3030
1740	22.5	22.0	4990	38.9	38.1	3030

for model g_x (i.e., y vs. x) and 0.0002 for model g_z (i.e., y vs. z). The AIC_c values are 493.97 for model g_x and 494.00 for model g_z . The latter model being the estimated K-L best model, we select g_z as the better predictor. The two Δ_i are 0 and 17.01 for models g_z and g_x , respectively. The corresponding Akaike weights are 0.9998 and 0.0002.

The context in which we developed guidelines about model selection is having more complexity and more models than here. There is no overwhelming evidence for the superiority of one model over the other. The matter of interpretation of the strength of evidence is uncertain. For this model selection problem we recommend applying the guidelines as using the above analytical results. The bootstrap is quite useful, more so than in more complex model selection situations.

Based on 10,000 bootstrap samples of the data in Table 1, the AIC model selection we obtained results as follows. The distribution of Δ_p^* gives us estimated percentiles for the distribution of Δ_p as 1.29 (95th percentile), 2.56 (96th percentile), 4.00 (97th percentile), 9.84 (99th percentile), and 17 is at about percentile 99.85. These results are thus consistent here with the bootstrap sampling distribution. This gives us a basis to interpret Δ_i as regards the plausibility of model i actually being the K-L best model for the data. However, when we compare the selection relative frequencies from the bootstrap we find that 93.8% of the 10,000 bootstrap samples. This is still strong evidence for z as the better predictor.

It is clear that we select, based on strong evidence, $E(y) = c + dz$ as the better model. The estimates of the slope of this model, and their conditional standard errors, are $\hat{c} = 252.9$, $\hat{d} = 183.3$ ($\hat{se} = 9.3$). We do not have a way we would like to compute unconditional standard errors when a parameter is fixed in a model (which d definitely is). However, when the evidence for the selected model, such as here, it is reasonable to act as if that model, hence act as if that model would always be the best model for the data, in which case conditional standard errors apply. So here we use the conditional standard errors as a measure of estimator variability. As a rule of thumb we will hazard the suggestion that if the selected model has an Akaike weight ≥ 0.90 , it is acceptable to use the conditional standard errors. The exact value (i.e., 0.90) is not critical; the concept is that if the evidence supports the selected model strongly enough ($w_{min} \geq 0.9$) it is acceptable as a rule of thumb; also see Royall (1997), where a similar rule of thumb is suggested. Conditional and unconditional standard errors will be nearly identical if one is bothered by the issue here of using conditional standard errors in data-based model selection, a modest simulation study can be conducted on that issue, as well as other matters.

introduce one use of simulation here, namely generating data that closely mimic the apparent nature of the real data. This is a model selection in a case like these pine wood data when we have generated the data. Generally, we would be against such a use as having relevant applicability to AIC model selection is the issue is clearly one of just deciding between two linear complex data simulation models than used below do not.

We proceed by considering $(y, x, z)'$ as a trivariate normal with mean vector $\underline{\mu}$ and variance-covariance matrix Σ . The diagonal matrix \underline{D} has as its diagonal the marginal standard deviations of x , and z , and C is the matrix of correlations

$$\begin{bmatrix} 1 & \rho_{yx} & \rho_{yz} \\ \rho_{xy} & 1 & \rho_{xz} \\ \rho_{yz} & \rho_{zx} & 1 \end{bmatrix}.$$

From the data we obtain $\hat{\underline{\mu}}' = (2992, 27.86, 26.79)$; the standard deviations are 894.60 ($= \hat{\sigma}_y$), 4.4946 and 4.6475. The estimates are $\hat{\rho}_{yx} = 0.9272$, $\hat{\rho}_{yz} = 0.9521$, and $\hat{\rho}_{xz} = 0.9999$. A simulated observation mimicking the data we generated is a standard normal random variables (i.e., normal(0,1)), say

$$(y, x, z)' = \underline{\mu} + \underline{D}\underline{C}^{0.5}\underline{v}$$

for some parameter choices “near” the estimated parameter values. There are software packages that will find the needed “square root” of a positive definite matrix (MATLAB, Anonymous 1994).

The best model here is the one that has the smaller residual variance of y , given the predictor. Those true residual variances are $\sigma_{y|x}^2$ and $\sigma_{y|z}^2 = \sigma_y^2(1 - \rho_{yz}^2)$. Therefore, in the simulation, the model predicting y is the one with the biggest correlation coefficient. What we cannot determine without simulation is performance of the model selection method.

We can tell from theory that only the values of ρ_{yx} , ρ_{yz} , and ρ_{xz} affect selection performance, including the distribution of Δ_p , the distribution of selection frequencies, hence Akaike weights. Hence, parameter values $\underline{\mu}$ and \underline{D} are irrelevant to that aspect of the problem (we might be interested about values for $\underline{\mu}$ and \underline{D} if we wanted a realistic evaluation of the bias on parameter estimators). Therefore, in the simulation we set $\underline{\mu} = \underline{0}$ and $\underline{D} = \underline{I}$. Using these values results in $a = 0$ and $\bar{d} = \rho_{yz}$, and knowing these as truth, we can infer relative selection bias that might occur.

We generated 10,000 simulated observations for the trivariate normal, the coefficients being at their estimated values (as truth) and the

of ρ_{yz} and $\rho_{yx} = 0.925$ as in the real data and a final case where the correlation coefficients used in the simulation were set to 0.95. The main objective was to determine the relative frequency of model selection (this is without loss of generality since $\pi_z + \pi_x = 1$) and the weight $E(w_z)$ (also without loss of generality since $E(w_z) + E(w_x) = 1$) the 95th and 99th percentiles of Δ_p (denoted below by $\Delta_{p,0.95}$ and $\Delta_{p,0.99}$). Results are given below by assumed sets of correlation coefficients. The proportions have coefficients of variation of about 1%; the Δ_p have coefficients of variation more like 2.5%, and this for 10,000 simulations.

ρ_{yx}	ρ_{yz}	ρ_{xz}	π_z	$E(w_z)$	$\Delta_{p,0.95}$
0.927	0.952	0.958	0.97	0.96	0.0
0.927	0.952	0.900	0.90	0.89	4.6
0.927	0.952	0.980	0.99	0.99	0.0
0.900	0.925	0.958	0.92	0.90	1.8
0.900	0.925	0.900	0.84	0.83	7.6
0.950	0.950	0.950	0.50	0.50	18.3

In these cases it is clear that $\pi_z \approx E(w_z)$, and that the sample size of Δ_p is quite variable. This is a worse case as regards variability of Δ_p (only two models, and they are nonnested). We can see that in the case above were reality, we would expect to select the correct model in 96% of all samples (for $n = 42$). These Monte Carlo results give us some faith in the usefulness of the bootstrap results based on the bootstrap. We have faith in the strength of evidence deduced here from the data and from $w_z = 0.9998$.

By looking at the more detailed results (not given here) for the estimated parameters, $\hat{\theta}$, and their averaged estimated standard errors, given the selected model we can assess model selection bias. If case 1 above were truth, the suggested there would be little model selection bias here with respect to the standard errors applied (to be expected if the model is selected 96% of the time). When model g_x was selected, no strong model selection estimators were suggested, but the sample size for this was small ($m = 314$).

There is another interesting question we can explore with these results. When a model is selected, right (g_z) or wrong (g_x) is selected, we then judge the weight of evidence to be strongly in favor of the right model. Our interest in such a question is mostly focused on when the model is the right choice (we will not know this to be the case): Having picked the right model, not the K-L best model, will the data appear strongly to support the right model as being best, or will the evidence be weak? For case 1 above, the wrong choice with sampling probability only about 4%

evidence in favor of the wrong model. Thus only 1% of the 9,688 samples wherein model g_z was selected, 9,223 produced $(w_z > 0.9)$ in favor of the selected model. Hence we expect that in all samples (in this particular scenario) we would select the correct model, do so with convincing evidence. Note however, that we cannot know the actual data whether it is one of the “1%” strongly misleading, say, again just for this simulated scenario, that the *estimate* that we have reached a correct conclusion for these pine

4.10 The Durban Storm Data

Linhart and Zucchini (1986:176–182) apply AIC to storm data at the Botanical Gardens in Durban, South Africa. The details are in their Table 10.1. By seven-day periods in the year (“week i ” starting 1 January, they obtained the frequency of weeks with at least one storm occurring. For example, in 47 consecutive years of data, there were 6 years with at least one storm event. The data are defined by the definition of a storm: “a rainfall event of at least 30 mm in 24 hours” (and Zucchini 1986:176). We use here their period I data (years 1932 to 52), wherein for the first 22 weeks the sample size of years is 22, and for weeks $23 \leq i \leq 52$, $n_i = 48$. Thus, under a fixed effects model, the total sample size here is 2,474. The data are from January 1932 to December 1985, as did Linhart and Zucchini (1986), the minor matter of a year. To have 8 days (such as 26 February to 4 March when there was a storm). Listed in order $i = 1$ to 52, the data y_i are

6, 8, 7, 6, 9, 15, 6, 12, 16, 7, 9, 6, 8, 2, 7, 4, 4, 3, 3, 10, 4, 0, 2, 0, 3, 1, 1, 5, 4, 3, 6, 1, 8, 3, 4, 6, 9, 5, 8, 6, 5, 7

Conceptually, there exists a probability p_i of a storm at the Botanical Gardens in week i . Based on these data, what is a “good” estimate of p_i ? That was the analysis objective of Linhart and Zucchini, and it is one of our objectives. Our other objective is to reliably assess the uncertainty of \hat{p}_i . A simple estimator is $\hat{p}_i = y_i/n_i$; it is very nonparsimonious and (most seriously) fails to be a smooth, hence informative, function of trends in the true p_i . We expect that anyone considering the data should strongly believe that the p_i would have a considerable dependence on i as a function over the 52 weeks. Therefore, we want to fit a model for a not-large number of parameters represented by $\underline{\theta} =$

they construct a likelihood by treating the y_i as a set of independent random variables on sample sizes n_i for parameters p_i , and model as

$$\text{logit}(p_i) = \log(p_i/(1 - p_i)) = \sum_{j=1}^K \theta_j z_{ji}$$

being some suitable linear model on $\underline{\theta}$, for known “covariates” z_{ji} . This is a type of logistic regression (we consider theory for this in this situation in Section 7.6.6). Linhart and Zucchini used a series model for the z_{ji} and used TIC for model selection (which is essentially the same as AIC). We extend their example to model averaging; also, we compute unconditional confidence intervals for p_i . Here, K is the number of structural parameters in the model.

The structure of the simplest model, model g_1 , is given by

$$\text{logit}(p_i) = \theta_1, \quad i = 1, \dots, 52.$$

For model g_2 :

$$\text{logit}(p_i) = \theta_1 + \theta_2 \cos\left(\frac{2\pi(i-1)}{52}\right) + \theta_3 \sin\left(\frac{2\pi(i-1)}{52}\right)$$

For model g_3 :

$$\begin{aligned} \text{logit}(p_i) = & \theta_1 + \theta_2 \cos\left(\frac{2\pi(i-1)}{52}\right) + \theta_3 \sin\left(\frac{2\pi(i-1)}{52}\right) \\ & + \theta_4 \cos\left(\frac{4\pi(i-1)}{52}\right) + \theta_5 \sin\left(\frac{4\pi(i-1)}{52}\right) \end{aligned}$$

In general, the structure for model g_r (wherein $K = 2r$) is

$$\begin{aligned} \text{logit}(p_i) = & \theta_1 + \sum_{j=1}^{r-1} \left[\theta_{2j} \cos\left(\frac{2j\pi(i-1)}{52}\right) \right. \\ & \left. + \theta_{2j+1} \sin\left(\frac{2j\pi(i-1)}{52}\right) \right], \quad i = 1, \dots, 52. \end{aligned}$$

Assuming conditional binomial variation and independent likelihood for any model is

$$\mathcal{L}(\underline{\theta}) \propto \prod_{i=1}^{52} (p_i)^{y_i} (1 - p_i)^{n_i - y_i}.$$

Given the model for $\text{logit}(p_i)$ as a function of $\underline{\theta}$, say $h_i(\underline{\theta})$, compute p_i as

$$p_i = \frac{1}{1 + \exp[-h_i(\underline{\theta})]}.$$

Model	K^a	$\log(\mathcal{L})$	$\Delta\text{-AIC}$	χ^2	P	
1	2	-863.24	62.66	131.4	0.000	2
2	4	-833.83	7.85	76.5	0.007	1
3	6	-829.17	2.53	69.3	0.019	1
4	8	-826.37	0.93	61.2	0.054	1
5	10	-823.91	0.00	55.6	0.094	1
6	12	-823.89	3.95	55.6	0.064	1
7	14	-823.40	7.04	54.7	0.049	1
8	16	-822.76	9.70	54.0	0.035	1
9	18	-822.47	13.11	53.8	0.022	1

^aThe number of structural parameters plus 1 for \hat{c} for QAIC; the number of parameters in the model plus 1 for \hat{c} for AIC. In this example is $K - 1$.

The independence assumption may not be true, but it is not so badly wrong. Similarly, the count y_i may not be the sum of n_i independent Bernoulli events over the n_i years. Truth may correspond to having varying year-to-year weekly probabilities of a success. To cope with these types of model inadequacies is to use a variance inflation factor in the likelihood theory, hence to use a variance inflation factor computed from the global model goodness-of-fit chi-square divided by the freedom df. Then we use QAIC, rather than AIC; also, the variances based on assumed models are multiplied by \hat{c} (Linhart and Zucchini 1986).

Following Linhart and Zucchini (1986) we consider several models over which model uncertainty and model averaging are considered. For g_7 , $K = 14$. We obtained MLEs for these models by using PROC NLIN (SAS Version 6.12); it is easy to adapt PROC NLIN to produce MLEs (see, e.g., Burnham 1989). In this example, it is not clear what should serve as the basis for estimation of the variance inflation factor. Several models were explored (Table 4.4), and estimates of the variance inflation factor were stable at about 1.40. For each fitted model we also computed the square goodness-of-fit statistic, its significance level (P -value), for the purpose of a more thorough consideration of model fit with the global model. For g_8 ($K = 16$) and g_9 ($K = 18$). Table 4.4 gives basic results for these models: K , $\log(\mathcal{L})$, $\Delta\text{-AIC}$, χ^2 goodness-of-fit, and corrected χ^2 , \hat{c} , and $\Delta\text{-QAIC}$. The values of $\Delta\text{-QAIC}$ are for when model 5 is the global model.

The $\log(\mathcal{L})$ values in Table 4.4 for models 1 through 9 are from Linhart and Zucchini (1986) in their Table 10.3 (they denote models 1 and g_9). The AIC-selected model has 9 parameters, and the results of Linhart and Zucchini (on their page 182). The

nor their estimated conditional standard errors (standard errors applied to obtain the large-sample variance–covariance matrix). Rather, our goal is to estimate well the set of p_1 to p_{52} , parameters in common to all models.

4.10.2 *Consideration of Model Fit*

Before we accept the AIC-selected model, we must consider model fits. Based on the results in Table 4.4, the global model fit to the data: $P = 0.049$. More importantly, $\hat{c} = 1.4$ on a scale greater than 1 that we should not accept results of AIC that require $c \approx 1$. Even the AIC-selected model has $\hat{c} = 1.1$, even though this model is deliberately selected to fit well. To issue further we fit two more models; models g_8 and g_9 also. If the problem was an inadequate structural model, we would fit g_8 and g_9 , compared to model g_7 , to improve. The results strongly suggest that there is extrabinomial variation in the data; a result is common for real count data such as these, as is the case $1 < \hat{c} < \approx 2$).

However, before automatically resorting here to QAIC, it is worth noting. The expected counts from the models fitted to the data (i.e., the data are sparse in the sense of being small counts) for model g_5 , $\hat{E}(y_{26})$ to $\hat{E}(y_{32})$ are about 1.5; these are the expected count values here; the largest estimated expected count is 10. Perhaps even if the global model is structurally true, the count values will invalidate the usual central chi-square null hypothesis goodness-of-fit statistic.

We explored this matter by Monte Carlo methods (also called bootstrap method). We generated data based on truth being the selected model g_5 . That is, independent y_i^* were generated based on fitted model g_5 . For each such data set we then computed the chi-square goodness-of-fit statistic to see whether it was noticeably different from that of a central chi-square distribution. Only a sample of 100 such generated data sets because of the big effect: for this situation is $c = 1$ or 1.29?

The answer was clear: If the model truly fits, then on average $\hat{c} = 1$; i.e., the usual null distribution holds well here despite the average of the 100 χ^2 goodness-of-fit values was 41.9 (theoretical value is 40). The largest and smallest of the 100 values were 71.6 ($P = 0.997$); these are not unusual for a sample of 100. When each test statistic was converted to a P -value, the P -values fit a uniform (0, 1) distribution. Finally, the average of the \hat{c} was 0.98 ($\hat{se} = 0.21$). While the possibility remains that

based general statistical wisdom for real data supports the accept that extra binomial variation often exists in count data. We will use QAIC, not AIC, with $\hat{c} = 1.4$ as our basis for model selection.

The number of estimable parameters must be augmented by 1 to need to estimate the variance inflation factor (c); if more than one inflation factor is estimated, then the number of such estimable parameters must be included in K . If the estimate of c is close to 1, then no adjustment is necessary and K should not be increased. If one believes there is concern about overdispersion and, therefore takes $c = 1.4$, then K must not be incremented.

When sufficient precision is used in the calculations, we find that model g_3 is the QAIC best model, although for practical purposes model g_5 is tied for best (and model g_5 is almost as good, based on QAIC). Figure 4.1 shows a plot of the fitted \hat{p}_i for both models g_3 and g_4 . Also shown are the 95% confidence bands on p_i based on $\hat{p}_{L,i}$ and $\hat{p}_{U,i}$ for each model. To explain the calculation of these confidence intervals.

4.10.3 Confidence Intervals on Predicted Storm

Basically, $\hat{p}_{L,i}$ and $\hat{p}_{U,i}$ arise as back-transformed lower and upper limits on $\text{logit}(p_i)$. However, we used SAS PROC NLIN to compute the estimated MLE-based theoretical $\widehat{\text{se}}_t(\hat{p}_i | g)$ that is computed from the estimated MLE-based theoretical $\widehat{\text{se}}_t(\hat{p}_i | g)$. The first step is then to form the correct (inflated) estimated standard error $\sqrt{\hat{c}} \cdot \widehat{\text{se}}_t(\hat{p}_i | g) = 1.183 \widehat{\text{se}}_t(\hat{p}_i | g) = \widehat{\text{se}}(\hat{p}_i | g)$. The interval could be used. However, it is better to use here what is based on the back-transformed logit-based interval (Burnham et al. 1999).

$$\hat{p}_{L,i} = \frac{\hat{p}_i}{\hat{p}_i + (1 - \hat{p}_i)C},$$

$$\hat{p}_{U,i} = \frac{\hat{p}_i}{\hat{p}_i + (1 - \hat{p}_i)/C},$$

where

$$C = \exp \left[\frac{t_{\alpha/2, df} \widehat{\text{se}}(\hat{p}_i | g)}{\hat{p}_i(1 - \hat{p}_i)} \right]$$

(acceptable as long as \hat{p}_i does not get too close to 0 or 1). The confidence bands in Figure 4.1 were computed in this manner and are based on the model. We used $t_{\alpha/2, df} = 1.96$, because the df that apply is 39; thus here $df = 39$.

Estimates of unconditional standard errors require the use of the bootstrap, in this case based on Δ_i from QAIC (Table 2.6). w_1, \dots, w_7 are (from Section 2.6).

0.0000, 0.0833, 0.3149, 0.3149, 0.2465, 0.0342, 0.0000

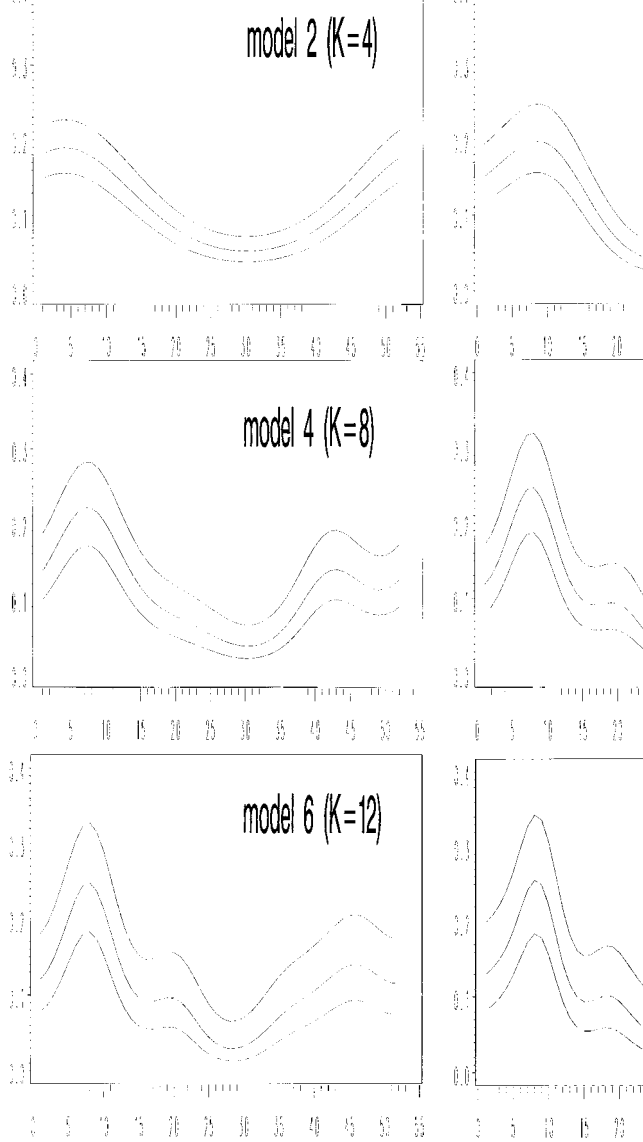


FIGURE 4.1. Plots of the predicted probability of one or more storms from models g_2 to g_7 fitted to the Durban storm data from Linhart and Linhart (1997) (see text for details). Also shown are approximate 95% confidence bands conditional on the model (see text for details).

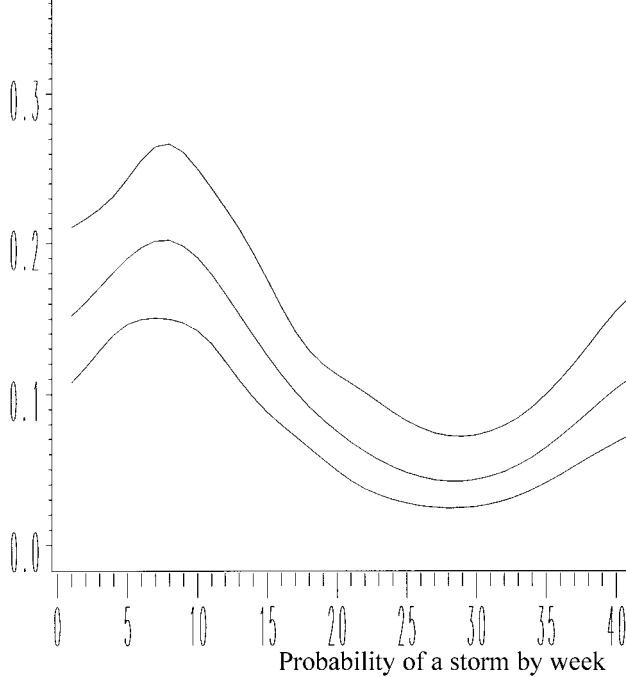


FIGURE 4.2. Plot of the model-averaged (from models g_2 to g_7) probability of one or more storms per week, $\hat{p}_{a,i}$, from the Durban storm data from 1986 (see text for details). Also shown are approximate 95% confidence bands; these bands include model selection uncertainty.

Then for each week we find \hat{p}_i under models g_1 to g_7 and average them to get the model-averaged $\hat{p}_{a,i}$. Next we apply (4.9) to obtain the unconditional standard error of \hat{p}_i . This unconditional standard error also applies to \hat{p}_i from model g_4 as well as to $\hat{p}_{a,i}$. In cases like this where the $\hat{\theta}$ are not available, we recommend that the \hat{p}_i to use are the model-averaged values. Here we use the unconditional $\widehat{\text{se}}(\hat{p}_{a,i})$, and we base confidence intervals on the above formula for $\hat{p}_{L,i}$ and $\hat{p}_{U,i}$, however, for C based on $\widehat{\text{se}}(\hat{p}_{a,i})$, $\hat{p}_{a,i}$, $\hat{p}_{L,i}$, and $\hat{p}_{U,i}$ are shown in Figure 4.2.

4.10.4 Comparisons of Estimator Precision

We now give some considerations about (estimated) standard errors for different versions of the \hat{p}_i . Common practice would be to use the standard errors conditional on the model. In this case, we would use model g_4 (by the slimness of margins). We compute

$$r_i(\text{se}) = \frac{\widehat{\text{se}}(\hat{p}_i)}{\widehat{\text{se}}(\hat{p}_{a,i})}, \quad i = 1, \dots, 52$$

and

$$r(\overline{\text{se}}) = \frac{\sum \widehat{\text{se}}(\hat{p}_i | g_4)}{\sum \widehat{\text{se}}(\hat{p}_{a,i})}.$$

These ratios are less than 1 if the unconditional is larger standard error (the notation used for these ratios has no s just need to represent them somehow).

We obtained $0.78 \leq r_i(\text{se}) \leq 1.02$ and $r(\overline{\text{se}}) = 0.99$. The unconditional standard errors are on average 1.11 times larger than those that are conditional on the model, and hence ignore model selection. We note that the average of the 52 values of $\widehat{\text{se}}(\hat{p}_{a,i})$ was $\widehat{\text{se}}(\hat{p}_{a,i}) \leq 0.035$. This is good absolute precision; the average width of the 52 confidence intervals was 0.084.

An alternative that avoids model selection is to use a parameter-saturated model is not very useful. Estimated standard errors under this model are given by $\widehat{\text{se}}_s(\hat{p}_i) = \sqrt{\hat{p}_i(1 - \hat{p}_i)/n}$, where any \hat{c} is used here since there is no basis on which to choose a variance inflation factor given this model). We computed $\widehat{\text{se}}_s(\hat{p}_i) / [\sum \widehat{\text{se}}(\hat{p}_{a,i})]$ and considered the separate ratios. We obtained

$$\frac{\sum \widehat{\text{se}}_s(\hat{p}_i)}{\sum \widehat{\text{se}}(\hat{p}_{a,i})} = 2.31;$$

so on average the unconditional standard errors of the \hat{p}_i were more precise by a multiplicative factor $0.433 = 1/2.31$ than the much less useful parameter-saturated model estimates. All $0 \leq \widehat{\text{se}}_s(\hat{p}_i) \leq 0.082$ (and a variance estimate of 0 is quite possible) and $0.012 \leq \widehat{\text{se}}(\hat{p}_{a,i}) \leq 0.035$; thus the standard errors for the $\hat{p}_{a,i}$ are much more stable than is the case for the parameter-saturated model.

Linhart and Zucchini (1986) used TIC for model selection. The only difference between the two methods is the use of $\text{tr}(J)$ for estimating this trace term, which Linhart and Zucchini do not. In fact, in Section 7.6.6 we demonstrate by theory and example that $\text{tr}(J)$ is very near K unless the structural model is truly terrible. In fact, it is better for count data (and generally simpler) to use $\text{tr}(J)$ than estimate this quantity. Linhart and Zucchini (1986:103) give in their Table 10.3. For models g_1 to g_7 the ratios of the estimated K are 0.95, 0.97, 0.97, 0.98, 0.98, 0.98, 0.98. Not only are they all close to 1, but they are all less than 1 (see Section 7.6).

Young and Young (1998, 510–514) give as an example the mortality of flour beetles (*Tribolium confusum*) caused by 24-hour exposure to gaseous carbon disulfide (CS_2); the data are from Bliss (1935). Table 4.5 gives the basic data, as dose level tested, and the number that died as an immediate causal response. We take sample size here as total beetles tested, hence n_i . The mortality rate (Table 4.5) increases with dosage in a roughly smooth, sigmoid form. It is typical to fit a parametric model to smooth such data, hence to get a simple estimated dose-response model, confidence bounds, and to allow predictions outside the observed range; extrapolation beyond the range of applied dose levels requires a response model; classically only one model was used. As we shall examine multimodel prediction of mortality at dose-levels.

A generalized linear models approach may easily, and appropriately, model the probability of mortality, π_i , as a function of dose level x_i . The probability distribution assumed for the data is binomial: at dose level x_i the response random variable y_i (number killed) can be distributed as $\text{binomial}(n_i, \pi_i)$. Therefore, the likelihood

$$\mathcal{L} = \prod_{i=1}^8 (\pi_i)^{y_i} (1 - \pi_i)^{n_i - y_i}.$$

Within this setting a model means some parametric form for $\pi(x)$ that $0 < \pi(x) < 1$ is maintained. Moreover, in the context of generalized linear models there must be a nonlinear transformation (i.e., link function) to give a linear structural model in the parameters. There is a natural choice here: as dose increases, modeled mortality must not decrease. There are several commonly used forms for such a

TABLE 4.5. Flour beetle mortality at eight dose levels of CS_2 (from Young and Young, 1998, their Table 14.1).

Dose (mg/L)	Number of Beetles		Observed mortality
	tested	killed	
49.06	49	6	0.12
52.99	60	13	0.22
56.91	62	18	0.29
60.84	56	28	0.50
64.76	63	52	0.83
68.69	59	53	0.90
72.61	62	61	0.98
76.54	60	60	1.00

We consider three commonly used generalized linear models: logistic, hazard, and probit (all are implemented in the LOGISTIC, SAS Institute Inc., 1985). The logistic model

$$\pi(x) = \frac{1}{1 + e^{-(\alpha + \beta x)}},$$

with link function

$$\log\left(\frac{\pi(x)}{1 - \pi(x)}\right) = \text{logit}(\pi(x)) = \alpha + \beta x$$

The hazard model and associated complementary log-log

$$\pi(x) = 1 - e^{-e^{(\alpha + \beta x)}},$$

and

$$\log[-\log(1 - \pi(x))] = \text{cloglog}(\pi(x)) = \alpha + \beta x$$

The cumulative normal model and associated probit link function (NORMIT) are

$$\pi(x) = \int_{-\infty}^{\alpha + \beta x} \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}z^2} \right] dz \equiv \Phi(\alpha + \beta x)$$

and

$$\Phi^{-1}(\pi(x)) = \text{probit}(\pi(x)) = \alpha + \beta x$$

Here, $\Phi(\cdot)$ denotes the standard normal cumulative probability function, which does not exist in closed form.

Traditionally, the predictor variable x can be either dose or log(dose). If dosages range over an order of magnitude or more, hence x is not too small, relatively near 0 compared to the largest dose, then log(dose) is preferred. It is used because then $\hat{\pi}(0) = 0$ will apply. Otherwise, with x = dose, above one risks getting $\hat{\pi}(0) > 0$, when in fact $\pi(0) = 0$. This is not recording acute mortality over a time interval so short that no deaths will occur. However, dose levels here are tightly clustered around zero, so either form of x is plausibly appropriate. One approach is to use models with both forms of x , thus six models.

A priori we do not know whether the simple model suffices to fit the data. Thus one might be motivated to extend the simple model to be, say, $\text{link}(\pi(x)) = \alpha + \beta x + \gamma f(x)$ where $f(x)$ is x^2 (or log(dose) if $x = \text{dose}$). Clearly, there is no unique model. So we consider six more models: the three link functions combined with $\alpha + \beta x + \gamma x^2$, where x can be either dose or log(dose). One might also consider without further thought. We do not recommend doing so. We present the results for these 12 fitted models and then note further similarities.

AIC	Δ	weight	K	predictors	
366.536	0.000	0.19356	2	dose	
366.641	0.105	0.18366	2	logdose	
367.608	1.072	0.11325	3	dose,	dose ²
367.698	1.162	0.10826	3	logdose,	logdo
367.804	1.268	0.10267	3	logdose,	logdo
367.998	1.462	0.09318	3	dose,	dose ²
368.420	1.884	0.07546	3	dose,	dose ²
368.430	1.894	0.07508	3	logdose,	logdo
370.246	3.710	0.03028	2	dose	
371.313	4.777	0.01776	2	dose	
374.138	7.602	0.00433	2	logdose	
375.222	8.686	0.00252	2	logdose	

that, in fact, allow an a priori restriction to just the three link($\pi(x)$) = $\alpha + \beta x$, for $x = \text{dose}$.

Table 4.6 gives AIC results for the 12 generalized linear models above. The models are identifiable by the predictors they use and the function used. There is considerable model uncertainty, as shown by the Akaike weights. However, a first thing to note is that the six models ($\Delta < 2$) are “paired” on dose and log(dose), by link function and predictor, in the sense that the two models of a pair have the same AIC value. There is a logical reason for this, and it should be noted that the models which would have led to consideration of only 6 models are log(dose) or dose.

For these data (Table 4.5) dosage is between 49 and 70, and the restricted interval log(dose) is almost perfectly linearly correlated with dose, $\log_e(\text{dose}) = 1.35376 + 0.007001 \cdot \text{dose}$, with $r^2 = 0.99$. This correlation justifies using either dose or log(dose), but not both. For our models. In fact, to include here the six models based on dose and log(dose) corresponding six models based on log(dose) is a form of overfitting (Section 4.6), and we recommend against it. Thus, a priori selection of six models; it is our choice to keep things simple and just stick with dose.

Of the remaining six models, a further a priori consideration is to eliminate the quadratic models such as $\text{logit}(\pi(\text{dose})) = \alpha + \beta \text{dose} + \gamma \text{dose}^2$. These three models cannot be monotonic increasing in dose if $\gamma \neq 0$. They might fit the data well (they do), but they might increase at lower doses or decrease at high doses. Because the quadratic models cannot be monotonic increasing in dose (unless $\hat{\gamma} = 0$ occurs), they have not been used. A plot of a fitted model will reveal this fact, at which point we feel justified in eliminating that model.

366.536	0.000	0.80114	2	dose	0
370.246	3.710	0.12534	2	dose	0
371.313	4.777	0.07352	2	dose	0

Table 4.7 shows the three models, of the 12, that can be compared on a priori basis. The associated Akaike weights change, relative to the model evidence ratios are invariant. For example, $0.1090 = 0.19356/0.01776$ for the hazard vs. logistic model. Figure 4.3 shows plots of the three fitted models.

Common analysis practice for these data would be to fit all three models of Table 4.7: historically probit, but in recent decades the logistic model. A recent analysis might have looked at all three models of Table 4.6) and used AIC to select the best model. This would have based inferences conditional on just that model. This ignores model uncertainty. Such uncertainty can be greater for extrapolation than the range of doses used than at the actual doses. To illustrate model inference (model averaging, here) we predict mortality at $x = 40$ ($\hat{\pi}(40)$).

Because there is a causal relationship of dose and mortality, extrapolation beyond the data, while risky, is both reasonable and often necessary. When extrapolation is done model deficiencies can be important. Figure 4.4 gives plots of four fitted models. The first model in Table 4.7 and the third model in Table 4.6, the logit link on the x-axis. The latter model predicts increasing mortality as dose goes beyond the range even though it provides a quite acceptable fit within the range of doses. Even if only discovered after the fact, we would use this model as a basis for eliminating this model.

Table 4.8 gives the basic conditional and unconditional model evidence ratios for the models considered. The estimated best model has a weight of 0.1090, not overwhelming (the evidence ratio for the best vs. second best is 6.39). The point estimates of $\pi(40)$ from the three models are 0.0257, 0.0308 and 0.0274, and produce, at these extremes, nonoverlapping intervals (Table 4.8). This information should be reported if known. It is not critical. If a single best-point estimate is acceptable, use the best model. $\hat{\pi}(40) = 0.0257$ with unconditional 95% confidence interval of 0.01274. Note that the unconditional standard error is 35% larger than the conditional standard error for the AIC best model. For a confidence interval based on the model-averaged $\hat{\pi}(40)$ ($= 0.0257$) and its unconditional standard error (0.01274), we used the logit-based interval of Burnham and Anderson (see also Section 4.10.3).

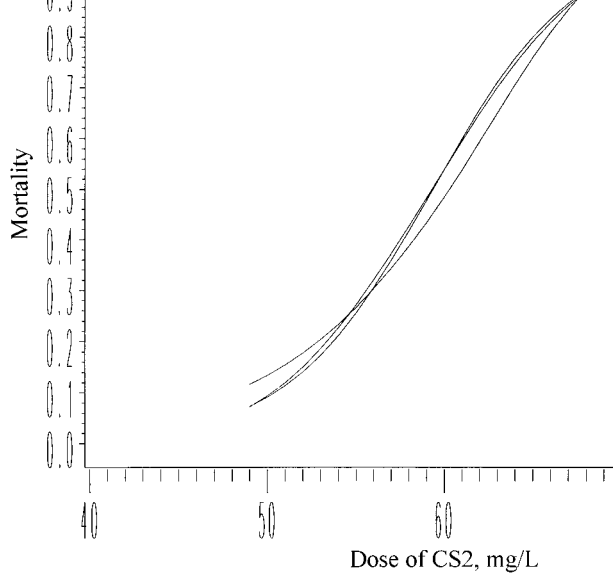


FIGURE 4.3. Fitted models of Table 4.7 plotted only over the range of dose 40 to 70 mg/L. The probit and logit models are nearly identical compared to the slightly different log-log model.

TABLE 4.8. Predicted mortality probability at dose 40mg/L; shown on each of the three models and the unconditional model-averaged prediction.

Link function	Akaike weight	$\hat{\pi}(40)$	\widehat{se}	95% Confidence interval lower
cloglog	0.8011	0.0308	0.00945	0.0168
probit	0.1253	0.0031	0.00226	0.0007
logit	0.0735	0.0085	0.00382	0.0035
model averaged:		0.0257	0.01274	0.0094

All three models used for Tables 4.7 and 4.8 fit the data well. There is no need for any overdispersion adjustment. A simple Pearson chi-square comparison suffices:

$$\chi^2 = \sum \frac{(O_j - \hat{E}_j)^2}{\hat{E}_j},$$

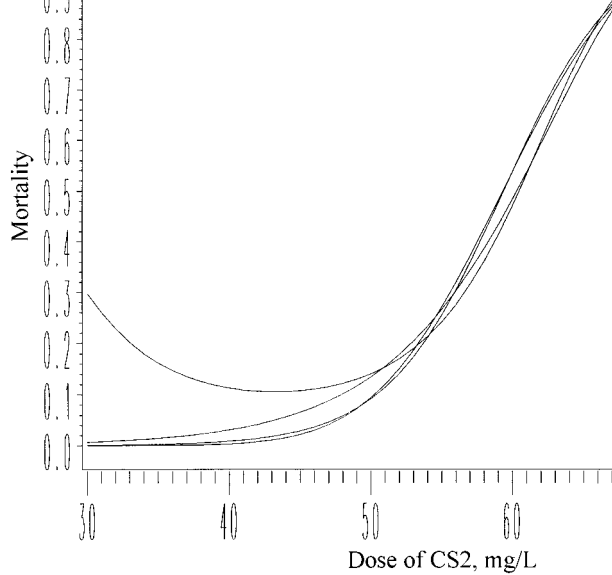


FIGURE 4.4. Four fitted models plotted beyond the range of dose for the tertiary log-log, probit and logit models of Table 4.7 (and Figure 4.3), and the dose and dose² (Table 4.6), which wrongly predicts increasing mortality at 44 mg/L.

on 6 degrees of freedom (= 8−2 since each model has 2 estimated parameters). For these binomially distributed data this chi-square statistic is

$$\chi^2 = \sum_{i=1}^8 \frac{(y_i - n_i \hat{\pi}_i)^2}{n_i \hat{\pi}_i (1 - \hat{\pi}_i)}.$$

Goodness-of-fit results are as follows

model	χ^2	P
cloglog	3.49	0.74
probit	7.06	0.32
logit	7.65	0.26.

A possible outcome, not observed here, is the best model fitting the data, say with $P = 0.01$ and $\hat{c} = 2$. If the worst-fitting model is not a global model, but is just an alternative model of the same dimension as the other models, we should attribute the failure to structural inadequacies, not overdispersion. Therefore, we would not need overdispersion adjustments). Overdispersion would result in rejecting all three models, i.e., even the best model would not fit well.

dose caused the observed mortality, whereas for the storm data prediction. For that reason it sufficed with the storm data a function plus a series of increasingly general nested models to find a global model for that example. For the beetle data we concluded that (1) the only predictor needed, or useful, is dose and that the expected response must be imposed by any model. The issue was thus reduced to one of an appropriate functional form, hence the linear models framework, to what is the appropriate link function. As a result, we have no global model, but rather several alternatives for a best causal-predictive model (this situation is common in observational studies). We saw that this has implications for model selection and overdispersion evaluation.

Moreover, an important issue here is prediction outside the range of data applied. It is well known that many models might fit the observed data, but give quite different extrapolated predictions. It is thus not a sound idea to pick a single model and use its base extrapolated predictions on it when there is model uncertainty. This example illustrates how easy it is to compute model-averaged predictions from the information from the fitted models, and that so doing leads to more realistically cautious predictions.

4.12 Publication of Research Results

We provide an outline of how results under the information-theoretic approach might be presented in papers submitted for publication (see Anderson et al. 2001d). The *Introduction* is the place to state the study objectives and note the degree to which the paper is exploratory or confirmatory.

Chamberlin's (1965) concept of "multiple working hypotheses" underlies the *Methods* section. This is the place to describe the hypotheses and models in the set and how these relate to the study objectives. Ideally, one should be able to justify why a particular model is in the set, as well as support the decision to exclude another model. Avoid the routine inclusion of a trivial hypothesis or model. All the models considered should have some reasonable support. The *Methods* section should always provide sufficient detail so the reader can understand what was done.

A common mistake is the use of AIC rather than the corrected AIC, AIC_c . Use AIC_c unless the number of observations is at least 10 times the number of explanatory variables (i.e., $n/K > 40$) for the most parsimonious model in the set. If using count data, some detail should

rounds, the log-likelihood should be computed as $\log(L)/n$, used in model selection, and the covariance matrix should be $\hat{\sigma}^2$. When using QAIC or QAIC_c, one parameter must be estimated: the number of estimated parameters in the model (K). When the estimation of the overdispersion parameter c is not required, the estimation has been identified (AIC, AIC_c, or QAIC_c), it should be reported for all models in the set. The adequacy of the global model should be assessed when the response variable is continuous (e.g., regression), and the use of procedures to aid in this task and many general ways to assess variation.

Discuss or reference the use of other aspects of the information approach, such as model averaging, confidence sets on model coefficients, or the relative importance of variables. Define or reference the use of these (e.g., K , Δ_i and w_i). Ideally, the variance component due to uncertainty should be included in estimates of precision (e.g., vs. conditional standard errors), unless there is strong evidence for the best model, say an Akaike weight (w_i) > 0.9.

For well-designed, strict experiments in which the number of factors is small and factors are orthogonal, use of the full model is appropriate (rather than considering more parsimonious models). The objective is to assess the relative importance of variables, then select models to provide a balance among variables (Section 4). The relative importance of variables can then be based on the Akaike weights for each variable, across models, and these sums can be used. Avoid the implication that variables not in the selected model are unimportant. Give estimates of the important variables (effect size) and measures of precision (preferably a confidence interval). Evidence can be assessed using the Akaike weights (w_i). Provide quantities of interest from the best model, or other models (e.g., $\hat{\sigma}^2$, coefficient of determination, estimates of model parameters, standard errors, insights from an analysis of the residuals). In general, evidence should be gathered, interpreted, and presented clearly.

The *Results* section should be easy to report if the *Methods* section convincingly the hypotheses and associated models of interest. Report the value of the maximized log-likelihood function ($\log L$), the number of estimated parameters (K), the appropriate selection criteria (AIC, QAIC_c), the simple differences (Δ_i), and the Akaike weights (w_i) in the set (or at least the models with some reasonable likelihood, where $\Delta_i < 10$). It is often helpful to report the models in the set based on the differences (Δ_i); this makes it easy to see which is the second best, and so on.

Do not include test statistics and *P*-values when using the information theoretic approach as this inappropriately mixes

The classical tests one would use for this purpose are in pairs have been ordered by an information criterion. Do information-theoretic approaches are a test in any sense **terms such as significant and not significant, or reject** instead, view the results in a strength of evidence context.

If some data dredging and modeling were done after then be sure that this is clearly explained when such re in the *Discussion* section. It is important to separate analytical questions and hypotheses formed before examining the data sequentially from examining the results of data analyses tends to be more confirmatory, while the second approach is exploratory. In particular, if the analysis of data suggests leading to an interesting hypothesis, then at this midway, tests or measures of precision remain valid. That is, an interesting patterns or hypotheses as being an actual feature of the population of interest are not well supported. Conclusions reached after re of the results of prior analyses, while interesting, cannot be of the degree of confidence as those from the more confirmatory post hoc results, while somewhat likely to be spurious, may hypotheses to be readdressed with a new, independent set an important part of good science.

4.13 Summary

Model selection should not be considered just the search. Rather, the basic idea ought to be to make more reliable the entire set of models that are considered a priori. This and scale the set of models and determine, perhaps, a cross the R models, for the K-L best model. Parameter estimation use of all the models when appropriate (e.g., model averaging) and attempt to use unconditional variances unless the set is strongly supported (say its $w_{min} > 0.9$).

In general there is a substantial amount of model selection in many practical problems (but see the simulated starling and pine wood data, for exceptions). Such uncertainty about variances (and associated parameter values) is the K-L best approximation whether one uses hypothesis testing, information-theoretic consistent criteria, cross-validation, or various Bayesian methods. There is a nonnegligible variance component for estimated parameters (prediction) due to uncertainty about what model to use, which should be included in estimates of precision.

bootstrap method. The simple computation of the relative Akaike weights, w_i , is often effective and easy to understand. (w_i provide a basis for a measure of relative support of the model, which is what we can call model scaling. The use of evidence ratios is effective and convenient.) The bootstrap, while computationally intensive, provides estimates of model selection probabilities, $\hat{\pi}_i$.

If there is substantial model selection uncertainty and the parameter estimates are estimated conditionally on the selected model, then the estimated parameters will likely be overestimated, and the coverage of the confidence interval will be below the nominal level (e.g., 95% instead of 95%). Estimates of unconditional standard errors can be computed using either Akaike weights or bootstrap selection probabilities. This is done by computing the R models considered over all bootstrap samples and then averaging the computed unconditional standard errors of estimated parameters.

Relatively little work has been done to understand and reduce the bias (Section 1.6). Shrinkage estimators such as

$$\tilde{\theta} = w_+ \hat{\theta}$$

deserve investigation. In particular, the estimation of the variance component in the precision of parameter estimates needs further work.

Conditional inference can be relatively poor. There is a need for research on general methodology to incorporate model selection uncertainty as a variance component in the precision of parameter estimates. This has potential in some applications where interest is concentrated on a small number of models that appear (implicitly, at least) in every model in the set. Problems can be approached as a problem in (weighted) model averaging. The weights can be either Akaike weights (w_i) or bootstrap selection probabilities.

The importance of a *small* number (R) of candidate models for a detailed analysis of the data, cannot be overstated. Small is in contrast to what is not small: If the number of models considered is large, then R is small. If the sample size is large, then R is small. This condition often occurs in the commonly used case of all possible linear models in variable selection. In that all-models case, if there are p variables, then $R = 2^p$; hence $R = 1024$ for $p = 10$ and $R = 1,048,576$ for $p = 20$. One should have R much smaller than n . Model selection is increasingly important in cases where there are many models to consider, the background science is lacking, so that needed a priori information is deemed impossible, then the analysis should probably be considered exploratory.

Finally, investigators should explain what was actually done in model selection (i.e., data analysis). Was it objective model selection?

introducing new models into consideration as data and the former case we recommend AIC (or its variants as mentioned above). In the latter case, if the strategy can be implemented as a computer algorithm, we recommend a bootstrap to assess model selection uncertainty for such searches. If the model selection strategy cannot be represented by a well-defined algorithm, one cannot determine model selection

Monte Carlo Insights and Examples

5.1 Introduction

This chapter gives results from some illustrative exploration of information-theoretic criteria for model selection and precision when there is model selection uncertainty. The Chapter 4 are illustrated and additional insights are provided on simulation and real data. Section 5.2 utilizes a chain binomial Monte Carlo evaluation of unconditional sampling confidence intervals, and model averaging. For this simulation process is known and can be of relatively high dimension. The model and the models used for data analysis in this chapter are easy to understand and have no nuisance parameters. Comparisons of AIC versus BIC selection and use of achieved coverage as an integrating metric to judge the success of Monte Carlo to inference.

Section 5.3 focuses on variable selection (equivalent to a test in multiple regression for observational data assuming normal homogeneous errors). A detailed example of AIC_c data analysis is given, and it is shown how to extend the example to a regression investigation of methodology. The same Monte Carlo method generate additional specific illustrative results. We also discuss how to reduce a priori the number of variables. A discussion of details of model selection bias and other ideas with emphasis on better approaches to, all-subsets selection.

One recent example uses distance sampling data on kangaroos. The reader is encouraged to read these examples, because they convey some general information about data analysis under the approach.

In the following material it is important to distinguish between variances, or standard errors, that are conditional on one parameter versus those that are unconditional, hence not based on just one parameter. Conditional measures of precision based on a restrictive model are noted by $\text{var}(\hat{\theta}_j | \theta_j)$ or $\text{se}(\hat{\theta}_j | \theta_j)$; however, they are more properly noted by $\text{var}(\hat{\theta}_j | g)$ or $\text{se}(\hat{\theta}_j | g)$. Corresponding unconditional variances are noted by $\text{var}(\hat{\theta}_j)$ or $\text{se}(\hat{\theta}_j)$, and have (as needed) an added variance to model selection uncertainty. Even such “unconditional” variances do depend on the full set of models considered.

The “as needed” phrase is appropriate because some models are design-based rather than model-based (see, e.g., Schreuder et al. 1996, and Edwards 1998). This applies, for example to the use of $\hat{\theta}_j$ from a random sample; or here to the use of just $\hat{S}_i = n_{i+1}/n_i$ for a restrictive model assumed for this inference. However, model selection is usually necessary and can be very effective for parameter estimation on complex data.

5.2 Survival Models

5.2.1 A Chain Binomial Survival Model

We consider here tracking a cohort of animals through the life span process of the cohort. Thus at time 1 there are a known number n_1 alive. Usually, we are thinking that these are all young animals, or born at time 1. One year later there are n_2 survivors; in general, at anniversary dates $i > 1$ there are n_i survivors. Eventually at the last anniversary year ℓ , $n_\ell > 0$ but $n_{\ell+1} = 0$. Given the information to date, the survival probability in year one (i.e., for the first year of life) is defined as $S_1 = E(n_2 | n_1)/n_1$. In general, for year i , $S_i = E(n_{i+1} | n_i)/n_i$. The obvious general estimator is $\hat{S}_i = n_{i+1}/n_i$; this estimator is valid for any assumed model structure imposed on the set of basic parameters. If n_i is large enough, then \hat{S}_i is an acceptable estimator in the absence of any assumed model structure. However, as the cohort dies out, n_i becomes small, and then the estimator is biased. In particular, in the year ℓ that the cohort dies out we always get $\hat{S}_\ell = 0$. This is a terrible estimator of age ℓ survival probability. The problem, and to improve age-specific survival estimates from the data.

The quantity n_1 is known at the study initiation, while n_2, n_3, n_4, \dots , being not then known, are treated sequentially as unknown variables and known ancillaries in the eventual data analysis. We will generate, in the Monte Carlo study, n_2 as binomial(n_1, S_1), where the survival process is sequential in time, once the survivors at age 1 are known, we can model n_3 , given n_2 , as binomial(n_2, S_2), where we know n_i ; hence we can generate n_{i+1} as a binomial(n_i, S_i). Thus we generate our probability model, hence likelihood, as a chain binomial model with conditional independence of n_{i+1} given the temporally preceding n_1 to n_i . The likelihood of the general possible model is

$$\mathcal{L}(\underline{S} | \underline{n}) = \prod_{i=1}^{\ell} \binom{n_i}{n_{i+1}} (S_i)^{n_{i+1}} (1 - S_i)^{n_i - n_{i+1}}$$

where \underline{S} and \underline{n} are vectors of the survival parameters and sample sizes, respectively. The underlying parameters explicitly in (5.1) are S_1, \dots, S_ℓ . In this context a (restricted) model for the parameters is some $S = (S_1, \dots, S_\ell)$, such as $S_i \equiv S$ for all ages $i \geq 1$ (model g_1), or perhaps $S_i = S$ only after age three; hence $S = S_4 = \dots = S_\ell$, and S_1, S_2, S_3 are free (restricted model g_4). If no structural restrictions are imposed, we have the general possible global model, g_i , for which the MLEs are

$$\hat{S}_i = \frac{n_{i+1}}{n_i}, \quad i \leq \ell,$$

with conditional sampling variances

$$\text{var}(\hat{S}_i | g_i) = \frac{S_i(1 - S_i)}{n_i}.$$

If such data were from moderately long-lived species such as wolves, one might expect substantial differences in survival probabilities across age classes, hence $S_1 < S_2$ and so on for several age classes. For salmon, be near-equal survival probabilities for adults until a decrease in survival as the surviving animals approach old age (i.e., senescence). If there are some large age-specific effects in survival following fishing or other harvesting effects (this conceptual model is not universally applicable, it does not apply to salmon). In real populations, the age-specific effects would be confounded with annual environmental effects and individual heterogeneity. We will not pursue this here and instead focus on some points related to conditional inference in model selection and some comparison of information approaches versus dimension-consistent criteria (e.g., BIC). We will illustrate the amount of variability in model selection as a function of sample size n_1 .



Ritei Shibata was born in Tokyo, Japan, in 1949. He received a B.S. in mathematics in 1973, and a Ph.D. in mathematics in 1981 from the Tokyo Institute of Technology. He is currently a professor in data science, Department of Mathematics, Faculty of Science, Yokohama National University, Yokohama, Japan. He has been a visiting professor at the Australian National University, University of Pittsburgh, Mathematical Science Research Institute (University of California), Victoria University of Wellington, New Zealand. He is interested in statistics, time series analysis, S software, and data and description (D&D).

Because of the paucity of data at older ages, some correction must be assumed, at least for the older ages, to get reliable estimates. One solution is to pool data beyond some particular age and assume that older ages have a constant survival probability. Thus, for animals, we assume that they all have the same survival probability after some age. It is not to say that truth is no longer age-specific, but rather that the *model* of the survival process is not age-specific after some age. The solution is to do some modeling (e.g., logistic models) of the survival process and do model selection. We will address both approaches in the next section.

Model	K	Parameters	
1	1	$S \equiv S_i,$	$i \geq$
2	2	$S_1, S \equiv S_i,$	$i \geq$
3	3	$S_1, S_2, S \equiv S_i,$	$i \geq$
\vdots	\vdots	\vdots	
R	R	$S_1, S_2, \dots, S_{R-1}, S \equiv S_i,$	$i \geq$

The pattern is obvious: Model g_r has r parameters with probabilities for ages 1 to $r - 1$, and constant annual survival probabilities for age r and older. The global model is g_R ; model g_r has K

This particular set of models is convenient for Monte Carlo simulation because all MLEs exist in closed form. Let the tail sum of survival probabilities alive at and after age r be denoted by $n_{r,+} = n_r + n_{r+1} + \dots$; the MLEs for model g_r ($r \leq \ell$) are

$$\hat{S}_i = \frac{n_{i+1}}{n_i}, \quad i = 1, \dots, r - 1,$$

$$\hat{S} = \frac{n_{r+1,+}}{n_{r,+}} = \frac{n_{r+1} + n_{r+2} + \dots}{n_r + n_{r+1} + \dots}.$$

Also, the likelihood for the r parameters given model g_r (constrained by model g_r)

$$\mathcal{L}_r(\underline{S} | \underline{n}) = \left[\prod_{i=1}^{r-1} \binom{n_i}{n_{i+1}} (S_i)^{n_{i+1}} (1 - S_i)^{n_i - n_{i+1}} \right] \left[\binom{n_{r,+}}{n_r} \right]$$

5.2.2 An Example

The first example compares the performance of AIC_c and BIC for sample size $n_1 = 150$. BIC is the Bayesian information criterion developed independently and from somewhat differing viewpoints by Schwarz (1978) and Akaike (1978, 1979). The form of BIC is

$$-2 \log(\mathcal{L}) + K \cdot \log(n).$$

This is a type of criterion that Bozdogan (1987) calls “dimensional selection criterion” since such criteria are often based on the assumption that the true model is in the set of candidate models. Then, BIC chooses the model with the smallest value. This is an estimate of the dimension of this true model with probability 1 as sample size increases asymptotically. Other interpretations weaker than this are “quasi-true models” (see Cavanaugh and Neath (1999)) and that BIC provides an estimate of the true model’s dimension.

0 in some detail. Here, we only present some comparisons of exposure to BIC or SIC.

The true survival probabilities S_1, S_2 , and so forth were 0.8, and then beyond S_5 , the age-specific annual survival of year (e.g., 0.784, 0.768, 0.753, 0.738, and so forth for S_6 and the data-generating model the survival parameters vary slightly to a maximum at ages 4 and 5, then simulating senescence decrease in S_j . For each Monte Carlo sample the cohort is simulated until all animals are dead. For example, the simulated animal counts ($n_i, i = 1, 2, \dots, 16 = \ell + 1$) were

150 (fixed), 74, 49, 34, 27, 21, 13, 9, 6, 6, 6, 4,

The set of models considered here is g_1 to g_{10} ; thus $R = 10$ models within the global model g_{10} . Hence, the global model provides separate estimates of S_1, S_2, \dots, S_9 but a single “pooled” estimate of S_{10} at age 9. The results for 10,000 Monte Carlo samples are given in Table 5.1. The two selection approaches yield substantial differences in model selection. On average, AIC_c selected an approximating model with $K = 2.1$ parameters while BIC selected a model with an average of 2.1 parameters.

The thinking underlying BIC is that the true model exists among the set of candidate models; this condition is not met here: True survival probabilities of ten models (the model closest to truth is model g_{10}). BIC attempts to estimate the dimension of the true model, a concept that is not well-defined here because there are nominally an unbounded number of candidate models.

TABLE 5.1. Comparison of model selection relative frequencies for 10 models g_1 to g_{10} and data from the chain binomial generating model with S_1 to S_{10} as 0.5, 0.7, 0.75, 0.8, 0.8, 0.784, 0.768, 0.753, 0.738, and 0.718 ($S_{j+1} = 0.98 \cdot S_j$); results are for sample size $n_1 = 150$.

Model	K	AIC _c model selection		BIC model selection
		percent	cumul. %	percent
1	1	0	0	0.2
2	2	52.0	52.0	88.6
3	3	22.4	74.4	10.0
4	4	5.6	80.0	0.5
5	5	5.5	85.5	0.5
6	6	5.0	90.5	0.1
7	7	3.4	93.9	0
8	8	2.4	96.3	0
9	9	2.2	98.5	0
10	10	1.5	100.0	0

models used were g_1 to g_{10} (see text), these results are based on 10,000 simulations for sample size $n_1 = 150$.

Age i	S_i	$E(\hat{S}_i)$		$E(\widehat{se}), AIC$		$E(\widehat{se}), BIC$	
		AIC_c	BIC	cond.	unc.	cond.	unc.
1	0.500	0.500	0.501	0.041	0.041	0.041	0.041
2	0.700	0.711	0.732	0.039	0.053	0.029	0.041
3	0.750	0.751	0.749	0.035	0.053	0.026	0.041
4	0.800	0.772	0.751	0.033	0.052	0.026	0.041
5	0.800	0.769	0.750	0.032	0.050	0.026	0.041
6	0.784	0.760	0.749	0.032	0.048	0.026	0.041
7	0.768	0.754	0.749	0.032	0.048	0.026	0.041
8	0.753	0.750	0.749	0.032	0.048	0.026	0.041
9	0.738	0.745	0.749	0.032	0.047	0.026	0.041
10	0.723	0.740	0.749	0.032	0.047	0.026	0.041
Average							

specific survival rates. However, there will not be enough information to estimate all those parameters, and increasing n_1 does not change that.

In sharp contrast, AIC_c attempts to select a parsimonious model as a basis for inference about the population sample. It does not assume that full truth exists as a model, nor does it assume that the true model is in the set of candidates. AIC_c estimates relative expected Kullback-Leibler distance and then selects the approximating model that is closest to the truth (i.e., the model with the smallest value of AIC_c). Based on $E(\Delta_i)$ (good to about two significant digits), it is model g_2 that is selected here, on average, but only by a minute winning margin. Therefore, without loss of generality we can set $E(\Delta_2) = 0$. For the other values of $E(\Delta_i)$ compared to this minimum value. In Table 5.2 we have the $E(\Delta_i)$ as 25.0, 0.0, 0.09, 1.2, 2.0, 2.7, 3.6, 4.1.

Further results from the Monte Carlo study are presented in Table 5.2. Formally, we look at the properties of the estimated age-specific survival rates under AIC_c and BIC model selection strategies. In Table 5.2 we give the estimated expected value of \hat{S}_i , where \hat{S}_i , by sample, is the estimated survival rate based on whatever model was selected as best. This estimator is computed for both AIC_c and BIC model selection. For example, if model g_4 is selected, then $\hat{S}_i = n_{i+1}/n_{i+1,+}$ for $i \geq 4$. For $i = 5, \dots, 10$, $\hat{S}_i = n_{5,+}/n_{4,+}$ for $i \geq 4$.

From Table 5.2 the estimators of age-specific survival rates are nearly unbiased; however, bias was slightly smaller for the AIC_c -selected models than for the BIC-selected models for 8 of the 9 ages. The expected estimated standard errors, conditional on the selected model, are also given in Table 5.2.

under both AIC_c and BIC selection are too small, since for any model selection uncertainty.

Various comparisons are possible from these simulation empirical standard errors versus expected estimated standard errors under the theory. However, confidence-interval coverage is an interesting way to see how well the methodology is performing. Therefore, we report the empirical confidence interval coverage for nominal 95% intervals, as $\hat{S} \pm 2 \hat{se}$. Use of this simple form may have resulted in a 1 to 3% failure in coverage of S_i for older ages, but we judge this to be in the contrast of conditional versus unconditional coverage.

Conditional confidence interval coverage in Table 5.2 for AIC_c model selection is generally below the nominal level, ranging from 0.959 over S_1 to S_{10} . Adjusting the conditional standard errors to unconditional (unc.) using (4.9) provides much improved coverage (95.0%), ranging from 0.906 to 0.977 (Table 5.2).

Achieved coverage using conditional standard errors for BIC model selection averaged 77.8% across the 10 age classes (Table 5.2). BIC model selection (there is no theoretical basis to justify this) improved the achieved coverage of true S_i for BIC model selection (range 60.1 to 98.7%); coverage generally remained below the nominal (95%). Buckland et al. (1997) present results on a survival analysis to the one used here.

This example shows the large amount of uncertainty associated with model selection when sample size is small ($n_1 = 150$). BIC model selection has tapering effects. This simple simulation exercise shows that BIC selection cannot be recommended. It requires very large samples to achieve consistency; and typically, BIC results in a model that is underfit (e.g., biased parameter estimates, overestimation of precision, achieved confidence interval coverage below that achieved by AIC models). Conditional estimates of precision, under either model selection, include model selection uncertainty, and this is often an important component. Incorporating model selection uncertainty can bring achieved confidence interval coverage up to approximately the nominal level (e.g., under AIC_c (using (4.7)) (this applies to AIC also).

The results in Table 5.2 are typical of many similar simulations we have examined (with different sample size, various sets of models R). To make more use of this particular example, we obtained confidence interval coverage on S_i , $i = 1$ to R , for AIC_c model selection for some additional sample sizes. We also tabulated the 90th, 95th, and 99th percentiles of Δ_i . This is useful for interpretation of the Δ_i . As sample size increases,

number of models considered, achieved confidence interval coverage, and the number of models selected. The results have been averaged over conditional (cond.) and unconditional (unc.) in

Sample size n_1	R	Best model K	Confidence interval		
			coverage, %		percentage
			cond.	unc.	0.90
100	10	2	84.4	95.5	4.0
250	10	3	77.8	93.3	4.8
500	10	3	75.5	92.2	7.6
1,000	10	8	78.6	92.8	5.9
1,000	12	8	76.0	90.6	6.4
10,000	10	10	86.2	93.0	1.0
10,000	20	13	78.3	88.9	5.4

model, should include more parameters. Thus one needs a larger size of the model set, i.e., R .

Unconditional confidence interval coverage is much smaller than conditional coverage (Table 5.3). When R is large, such as 20 for $n_1 = 10,000$ (over S_1 to S_{20}) confidence interval coverage suffers as much as the extrapolated estimates of, say, S_{20} based on a fitted model g_{13} . The S_i have a shallow peak at $i = 4$ and 5, and this feature is not a large sample size to detect reliably. Thus the theoretical coverage, though increasing as sample size increases, stalls at moderate sample size is passed, after which that theoretically best model is not responsive to increasing sample size. It is this feature that causes the percentiles of Δ_p to be so big (e.g., for $n_1 = 500$, a 95th percentile of 16.8 and a 99th percentile of 16.8).

The other notable feature in Table 5.3 is found when $n_1 = 10,000$. The theoretically best AIC_c model is then model g_{10} , the 10th model in the set considered. Model selection probabilities here in the set considered. Model selection probabilities here in the set considered. $\pi_9 = 0.18$; hence, Δ_p is 0 in 80% of the samples and small in the rest of the samples.

Use of BIC is theoretically inappropriate for model selection. This is evident from Table 5.3 because the set of models used does not include a model with fixed, small K and with all effects being fixed. As n_1 increases. However, BIC performs best at large sample size. It is not clear whether BIC performs well here when $n_1 = 10,000$ and the coverage is achieved (for nominally 95% intervals) for BIC. The results for these two cases is given below, along with average (over

confidence interval coverage (nominally 95%) has been averaged over 100 replications. The conditional (cond.) and unconditional (unc.) intervals on S_1 to S_R .

Sample size n_1	R	Best model K	Confidence interval coverage, %		
			cond.	unc.	perce 0.90
100	10	3	78.0	93.4	2.7
250	10	3	74.6	92.2	3.6
500	10	3	76.4	92.3	4.8
1,000	10	4	76.3	92.6	3.5
1,000	12	4	72.0	91.5	3.6
10,000	10	6	80.1	92.4	3.5
10,000	20	6	59.6	85.2	3.7

both BIC and AIC (AIC coverage is in Table 5.3):

n_1	R	coverage, BIC		average MS	
		cond.	unc.	BIC	AIC
10,000	10	77.8	83.7	0.000245	0.000245
10,000	20	45.9	51.5	0.003077	0.003077

Average mean square error is lower for AIC, and confidence interval coverage is much better. Using BIC, model g_8 is selected with probability 0.96 ($R = 10$) or 0.36 ($R = 20$), and the probability of selecting one of the other models is 0.96 ($R = 10$) or 0.93 ($R = 20$).

For some additional results and comparison to Table 5.3, we present the 5.3 results for a model wherein true $S_i = 0.5+0.3/i$. Table 5.3 presents the interval coverage results and percentiles of Δ_p for this case. Due to the monotonicity of the true S_i , the percentiles of Δ_p are approximately 4 (90th percentile), 5 (95th percentile), and 6 (99th percentile). Again, for the last case ($n_1 = 10,000$ and $R = 20$) the probability of selecting model g_{10} to averaging coverage over the intervals on all of S_1 to S_R is 0.96, a considerable extrapolations. Generally, AIC_c does not select model g_{10} . If for this last case we look only at intervals on S_1 to S_R , the levels under AIC_c model selection are 80.7% (cond.) and 85.2% (unc.).

5.2.3 An Extended Survival Model

A second Monte Carlo example is given here to add additional realism in the set of approximating models. This example uses a binomial generating model and associated parameters p and r (see 5.2.2; only the set of approximating models is different. Results for model g_r the constraint $S \equiv S_i$ for all $i \geq r$, the a

Model 1	$\text{logit}(S_i) = \alpha + \beta \cdot i,$
Model 2	$S_1, \text{logit}(S_i) = \alpha + \beta \cdot i,$
Model 3	$S_1, S_2, \text{logit}(S_i) = \alpha + \beta \cdot i,$
\vdots	\vdots
Model 10	$S_1, S_2, \dots, S_9, \text{and } \text{logit}(S_i) = \alpha + \beta \cdot i,$

Each model has an intercept and slope parameter for the fitting of age-specific survival rates S_i for ages $i \geq r$. Model g_r fits unconstrained age-specific survival rates for ages 1 to r . Model g_r has $K = r + 1$ parameters.

The initial population size was $n_1 = 150$ animals, as before. The animals were followed until the last one died. The true parameter values were $S_2 = 0.7$, $S_3 = 0.75$, $S_4 = 0.8$, $S_5 = 0.8$, and for $r > 5$, $S_r = 0.9$. The computer generated 10,000 independent repetitions of the simulation. In the analysis using the 10 approximating models, in some samples no model could be fit because all animals died before reaching age 1. In other samples, the last age at which $n_\ell > 0$ was $\ell - 2$. If $\ell - 2 \geq 10$, then all ten models were fit. Otherwise, only $\ell - 2$ models were fit (if $\ell - 2 = 9$, then 9 models). The number of models that could be fit to a sample varied somewhat, but this did not matter to the overall strategy of either selecting a best model or averaging, because for each fitted model it was always possible to compute \hat{S}_i (and its estimated conditional variance) for any age i .

Generally, the logistic approximating models were selected more often than when r was set as compared to the models used in Section 5.2.2. The best model (under AIC_c selection) was g_3 ($K = 4$) with model selection error $E(\Delta_4) = 0.17$ relative to setting $E(\Delta_3) = 0$). The model selection frequencies are shown in Table 5.5. In 84.5% of samples, model g_3 or g_4 was selected. Model selection uncertainty in this example is less than that of the example in Section 5.2.2 (compare Tables 5.5 and 5.6). In this example, the best model to use has 4 parameters, rather than 2, as was the case for the models previously considered. **Use of a better set of models can lead to more informative inferences from the data.**

For the simulation underlying Table 5.5 results, all 10 models were fitted in 9,491 samples. In 18 samples, only models g_1 to g_7 could be fitted (see the protocol noted above). Only models g_1 to g_8 were fitted in 363 samples. Only up to model g_9 were fitted in 363 samples.

Table 5.6 summarizes the estimated expected values of S_i for the 10 models and under model averaging (MA) using (4.2), and under model selection using the Akaike weights. Also shown in Table 5.6 is the estimated conditional variance of confidence intervals based on estimated unconditional survival rates (i.e., using (4.9)). Model averaging provided a slightly better estimate of the true conditional survival. However, both methods performed well.

to increase to a maximum at ages 7 and 9 and then decrease slowly (see Section 5.2.2 (see text for details). Model g_3 is the theoretically best model for sample size $n_1 = 150$.

Model	Percent	Cumul. %
1	00.0	00.0
2	35.2	35.2
3	29.5	64.7
4	19.8	84.5
5	6.3	90.8
6	3.5	94.3
7	2.2	96.4
8	1.5	98.0
9	1.2	99.2
10	0.8	100.0

TABLE 5.6. Summary of Monte Carlo results, based on 10,000 samples, for survival estimation under AIC_c selection and model averaging (MA). The true S_i had parameters S_i ; 10 logistic models (g_1, \dots, g_{10}) were fit to the data; results are for sample size $n_1 = 150$.

Age i	S_i	$E(\hat{S})$		Coverage	
		AIC_c	MA	AIC_c	MA
1	0.500	0.499	0.499	0.957	0.957
2	0.700	0.708	0.707	0.907	0.907
3	0.750	0.758	0.757	0.920	0.920
4	0.800	0.789	0.790	0.926	0.926
5	0.800	0.785	0.786	0.917	0.917
6	0.784	0.772	0.773	0.951	0.951
7	0.768	0.759	0.760	0.964	0.964
8	0.753	0.744	0.745	0.962	0.962
9	0.738	0.727	0.728	0.962	0.962
10	0.723	0.709	0.709	0.957	0.957
Average				0.942	0.942

The achieved confidence interval coverage for these intervals, based on model selection uncertainty, is very close to the nominal 0.95 coverage; both approaches use the same estimate of uncertainty, so the error is the same.

From Table 5.6 the unconditional confidence interval coverage is close to 0.95 over the intervals on S_1 to S_{10} , and over all 10,000 samples. This suggests that the strategy of inference based on the selected AIC_c best model is

0.000. In the worst case, which was at age 5 (i.e., for S_5), the MSE was only 0.717, versus unconditional of 0.917.

In this example the simulation program also computed the error (MSE) for estimators of S_i at each age for AIC selection, model averaging, and for the simple (almost model selection) age-dependent survival $\hat{S}_i = n_{i+1}/n_i$, $i = 1, 2, \dots$. In this example, the model-averaged estimator was smaller than that of the simple estimator (ranging from 11% to 25% smaller, except for age 5, where all approaches had the same MSE). In this example, at least, using Akaike weights, has important advantages. Both AIC selection and model averaging are quite superior to the use of the simple estimator of S_i . For example, at age nine, the MSEs for \hat{S}_9 were 0.0186 for selection, model averaging, and the simple estimator were 0.0186, respectively. This illustrates the advantage of *modeling* of the underlying survival probabilities (this is what was learned from the starling experiment in Section 3.4).

We computed the pivotal quantities $\Delta_p = \text{AIC}_{\text{best}} - \text{AIC}_p$ for each of the 10,000 samples to better understand these values for Δ_p of this type. Some percentiles of the sampling distribution of Δ_p are shown below:

Percentile	Δ_p
50.0	1.20
75.0	2.09
80.0	2.29
85.0	3.11
90.0	4.21
95.0	6.19
97.5	8.33
98.0	8.80
99.0	10.63

Thus, in approximately 90% of the simulated data sets the age, between AIC_c for the selected model, for that sample, the age of the theoretically best model was ≤ 4.21 . If we adopted the idea of an approximate 90% confidence set on what is the actual AIC_c theoretical best model, then in this example, fitted models that had Δ_i values below 4.21 would be in that set, and hence might be candidates for some further model making inferences from an individual data set.

To make further use of the generating model and set of ages, we determined several quantities for differing sample sizes. For each age, we determined the conditional and unconditional AIC_c for model selection (as opposed to model averaging), and 50th, 75th, 80th, 85th, 90th, 95th, 97.5th, 98th, and 99th percentiles of Δ_p . Regardless of this variation in sample size,

of the (approximate) pivotal Δ_p . These same results hold for generating models and fitted models that we have examined by other methods. Often the conditional coverage averages are less than nominal (more like 70% to 85%).

The true S_i do not exactly fit any of the approximations. However, model g_4 turns out to be an excellent approximation at least for $i \leq 15$ (we did not plan this to be the case). Here we know S_1 to S_{15} exactly, but we needed a mathematical model for larger i . If we knew the true S_i numbers, we would likely fit and use model g_4 . We would not use the fitted model then provides a concise, yet almost exact, summary of the data. The lack of fit is both trivial and statistically “significant” given the sample sizes. We do not generally have such huge sample sizes, but if we did, the concerns of small-sample-size statistics would not apply (perhaps we would then replaces statistical significance), and we might very well prefer g_4 rather than model g_8 as a preferred approximation to truth. This is a philosophical point: Even if we knew truth, we would of course not say said truth by a low-order parsimonious fitted model because the model itself confers about understanding the basic structure of the data.

5.2.4 Model Selection if Sample Size Is Huge, or

The results in Table 5.7 provide a motivation for us to re-examine philosophical issues about model selection when truth is essentially known, or equivalently in statistical terms, when we have a huge sample size.

TABLE 5.7. Some Monte Carlo results (10,000 repetitions per sample size) for model selection, for the same generating model and set of models g_i as used in Table 5.6. n_1 is the sample size; R is the maximum number of models considered; K is the number of models; AIC_c best model g_k ; achieved confidence interval coverage (nominal coverage) averaged over conditional (cond.) and unconditional (unc.) intervals.

Sample size n_1	R	Best model K	Confidence interval coverage, %		percentage of correct 0.90
			cond.	unc.	
50	10	2	84.8	92.6	3.8
100	10	3	86.2	94.1	3.7
150	10	3	86.0	94.2	4.2
200	10	4	86.1	94.2	3.6
500	10	4	88.3	94.6	2.7
1,000	12	4	90.4	95.4	2.6
10,000	15	4	90.3	95.5	4.9
100,000	15	8	86.4	93.3	5.0

could) be based on statistical criteria. There is no theory that says a deterministic model must fit deterministic truth in order for that model to be useful. Similarly, if truth as a probability distribution is known, a model approximating that truth is useful. If truth is unknown, then K-L can be used, but there is no reason why much information loss is tolerable. That is a subject-matter judgment.

5.2.5 A Further Chain Binomial Model

Here we consider a generating model with damped oscillations as age increases. It is not intended to be considered as a biological model, just a case of a complex truth against which we can examine model selection. We let the initial population size be $n_1 = 10,000$ independent samples using the chain binomial model with true survival probabilities as

$$S_i = 0.7 + (-1)^{i-1}(0.2/i),$$

where i is year of life (“age” for short). The data ($n_1, n_2, \dots, n_{18} = 0$ but $n_{17} > 0$) for one repetition were

1000 (fixed), 890, 520, 414, 270, 202, 130, 93, 67, 50, 35, 25, 15, 10, 5, 3, 2, 1, 0.

Animals were followed until all were dead. The models used were the logistic models g_i of Section 5.2.3. The model set was $\{g_i : i = 1, \dots, R = 15\}$. It is because sample size was 1,000 rather than 10,000 that the size of the model set used. It is a general principle that as sample size increases, one can expect to reliably estimate more parameters. The set of models considered should depend weakly on sample size.

Results for AIC_c and BIC model selection are very different. On average, AIC_c selected a model with 7.6 parameters, while BIC selected a model on average with 5.1 parameters. The theoretically expected K-L criterion, in this example is g_{15} with 15 parameters. Under BIC selection 95% of the models selected had 5 or 6 parameters. In contrast, under AIC_c selection approximately 95% of models selected had between 5 and 13 parameters, including g_{15} .

That BIC selection produces a more concentrated distribution of models certainly seems to be an advantage. However, useful models are not necessarily those that focus on bias and precision of parameter estimates, hence the importance of interval coverage (given that we have essentially the shortest intervals under the different model selection strategies). From Table 5.9, AIC_c selection produces average conditional and unconditional coverage of 92.7%, respectively. Under BIC selection we had average conditional coverage of 73.9% and 78.5%, respectively. AIC_c (but not always) did achieve a smaller MSE in these examples.

Of Table 5.9, MSEs averaged over S_1 to S_{15} were 0.0014

(-1) = (0.2/4), and the model set for data analysis as g_1 to g_{15} , size $n_1 = 1,000$.

Model <i>i</i>	<i>K</i>	AIC _c model selection		BIC model
		percent	cumul. %	percent
1	2	0	0	0
2	3	0	0	0
3	4	0.6	0.6	23.0
4	5	11.5	12.1	50.3
5	6	26.3	38.4	21.7
6	7	20.5	58.9	4.1
7	8	16.0	75.0	0.8
8	9	7.8	82.8	0.1
9	10	6.0	88.8	0
10	11	3.0	91.8	0
11	12	2.5	94.4	0
12	13	1.6	96.0	0
13	14	1.6	97.6	0
14	15	1.1	98.7	0
15	16	1.3	100.0	0

selection, whereas for AIC that average MSE was 0.002. At this cost of smaller MSEs (by BIC) for the \hat{S}_i is poor coverage; for S_1 to S_{15} , unconditional BIC coverage varied from 0.1 to 1.0 (see Table 5.9 for AIC coverage).

Comparison of the expected estimated standard errors in the magnitude of the variance component due to model selection uncertainty. For example, for age 10, the expected conditional standard error was 0.037, versus the expected unconditional standard error of 0.037. Table 5.9 shows results for ages 1 to 15.

We computed Δ_p values for each of the 10,000 repetitions. There were damped oscillations in the S_i parameters. Some of the sampling distribution of Δ_p are shown below:

Percentile	Δ_p
50.0	1.7
75.0	3.0
80.0	3.5
85.0	4.2
90.0	5.7
95.0	7.9
97.5	10.1
98.0	10.9
99.0	13.3

the generating model allowed damped oscillations in survival as for details). Conditional (cond.) and unconditional (unc.) standard intervals were also evaluated.

Age i	S_i	$E(\hat{S}_i)$	$E(\widehat{se}(\hat{S}_i))$		Cov cond.
			cond.	unc.	
1	0.900	0.900	0.009	0.009	0.950
2	0.600	0.600	0.016	0.016	0.955
3	0.767	0.767	0.018	0.018	0.951
4	0.650	0.651	0.023	0.024	0.904
5	0.740	0.734	0.023	0.027	0.821
6	0.667	0.678	0.025	0.031	0.701
7	0.729	0.716	0.024	0.032	0.671
8	0.675	0.690	0.024	0.034	0.711
9	0.722	0.706	0.025	0.035	0.767
10	0.680	0.693	0.028	0.037	0.838
11	0.718	0.699	0.030	0.039	0.871
12	0.683	0.691	0.035	0.043	0.887
13	0.715	0.692	0.040	0.048	0.900
14	0.686	0.685	0.046	0.054	0.902
15	0.713	0.681	0.052	0.061	0.908
Average					0.849

These results are consistent with percentile values for Δ_P for other examples.

Table 5.10 gives expected Akaike weights $E(w_i)$ and model probabilities π_i from the 10,000 Monte Carlo samples for this example. For results from this case of a large sample size ($n_1 = 1,000$ is the best), Table 5.10 also gives $E(w_i)$ and π_i for $n_1 = 50$ (g_3 is the best) under the same data-generating model. For sample size $n_1 = 50$, Table 5.10 that on average models g_4 to at least g_9 must be considered in inference about the population under AIC_c model selection. Akaike weights give some support to models g_{10} and g_{11} , while π_i gives somewhat less support to these models.

There is a high degree of model selection uncertainty, particularly for the larger sample size. The underlying process generating the data (the S_i) is complicated, and the set of logistic approximations was relatively poor. Had some science been brought to bear on this special problem, hopefully the set of approximating models would include some models with at least some oscillating features. The uncertainty would likely have been greatly reduced.

Use of the bootstrap to estimate model selection probabilities is computationally intensive. Akaike weights can be easily calculated.

allowing damped oscillations in survival as age increases, models g_1 to g_{15} ; two sample sizes are considered.

Model i	$n_1 = 50$		$n_1 = 1,000$	
	$E(w_i)$	π_i	$E(w_i)$	π_i
1	0.050	0.038	0.000	0.000
2	0.293	0.431	0.000	0.000
3	0.283	0.325	0.011	0.006
4	0.160	0.096	0.091	0.115
5	0.104	0.064	0.177	0.263
6	0.056	0.026	0.176	0.205
7	0.033	0.014	0.153	0.160
8	0.015	0.004	0.109	0.078
9	0.006	0.002	0.084	0.060
10	0.000	0.000	0.058	0.030
11	0.000	0.000	0.044	0.025
12	0.000	0.000	0.032	0.016
13	0.000	0.000	0.026	0.016
14	0.000	0.000	0.020	0.011
15	0.000	0.000	0.019	0.013

simple and effective alternative. From Table 5.10 we can see a close agreement between the expected Akaike weights and the model probabilities π_i . Note, however, that these are not estimates of the true probabilities that exact agreement is not expected. Akaike weights, not model probabilities, reflect the relative likelihood of each fitted model in the set. They can provide information about the relative support of the data for the different models. Finally, as all the examples in this section (and all the others) show, the weights w_i are very useful in model averaging and computing unconditional sampling variances, hence obtaining unconditional confidence intervals that do substantially improve coverage after model selection.

5.3 Examples and Ideas Illustrated with Linear Regression

The model selection literature emphasizes applications in time series, often as selection of variables in regression; this is not the focus here, as model selection. McQuarrie and Tsai (1998) is devoted to this topic. *Regression and Time Series Model Selection*. Our exploration is very much with those of McQuarrie and Tsai (1998), because we focus on model selection uncertainty and, in general, multimodel inference. We focus on all-subsets model selection by presenting an example.

5.3.1 All-Subsets Selection: A GPA Example

We use an example based on four regressors (Table 5.11) 20 (a larger number of regressors makes it too demanding for the full results). The example of Table 5.11 comes from (1994). They use these example data extensively to illustrate including all-subsets selection based on several criteria or the other ideas we use herein. Also, we note that the “realistic but not real” (H. Iyer, personal communication).

The full model to fit is

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 +$$

This model, g_{15} , is also denoted in tables by $\{1234\}$ because four variables are used (see, for example, Table 5.12). As another notation, model g_6 uses only the predictors x_1 and x_3 and is denoted by $\{1 \cdot 3 \cdot \cdot\}$:

$$y = \beta_0 + \beta_1 x_1 + \beta_3 x_3 + \epsilon.$$

Including the intercept-only model g_{16} , $\{\cdot \cdot \cdot \cdot\}$, there are 16 standard all-subsets approaches. We use AIC_c for selection of models, and we also apply the bootstrap (10,000 sample) full AIC_c model selection applied to all 16 models for each

The selected best model (g_{11}) includes predictors x_1, x_2 . However, support for g_{11} as the only useful model to use here is weak; its Akaike weight is just 0.454. Model g_5 , $\{12 \cdot \cdot\}$, is also competitive. Model g_6 , $\{1 \cdot 3 \cdot \cdot\}$. The Akaike weights w_i and bootstrap estimates $\hat{\pi}_i$ agree to a useful extent, and both demonstrate uncertainty in this example. For a confidence set on models, we are compelled in practice to provide one in this sort of application. We would use a likelihood ratio (evidence) criterion: the set of models for which $w_{11}/w_i \leq 8$ (cf. Royall 1997). Thus, we have the set of models for which $w_{11}/w_i \leq 8$. This criterion is identical to using as a cutoff models for which $w_{11}/w_i \leq 8$ would suffice).

We can explore for this example the coherence (with the selection of such a rule by determining the bootstrap distribution of $\Delta_p^* = AIC_{c,11}^* - AIC_{c,min}^*$. Here, $AIC_{c,min}^*$ is the minimum value of AIC_c over the given bootstrap sample, and $AIC_{c,11}^*$ is the AIC_c value in the bootstrap sample for model g_{11} . We computed the 10,000 bootstrap

First-year GPA	SAT		High School	
	math	verbal	math	Engli
y	x_1	x_2	x_3	x_4
1.97	321	247	2.30	2.63
2.74	718	436	3.80	3.57
2.19	358	578	2.98	2.57
2.60	403	447	3.58	2.21
2.98	640	563	3.38	3.48
1.65	237	342	1.48	2.14
1.89	270	472	1.67	2.64
2.38	418	356	3.73	2.52
2.66	443	327	3.09	3.20
1.96	359	385	1.54	3.46
3.14	669	664	3.21	3.37
1.96	409	518	2.77	2.60
2.20	582	364	1.47	2.90
3.90	750	632	3.14	3.49
2.02	451	435	1.54	3.20
3.61	645	704	3.50	3.74
3.07	791	341	3.20	2.93
2.63	521	483	3.59	3.32
3.11	594	665	3.42	2.70
3.20	653	606	3.69	3.52

give below some percentiles of this random variable:

percentile	Δ_p^*
50	1.00
80	3.21
90	3.62
95	5.18
98	7.00
99	8.79

Thus, we have a good basis to claim that the confidence on the expected K-L best model is about a 95% confidence. In statistical inference, this result depends on sample size, and if n increased.

Because model selection in regression is often thought of as identifying the importance of each predictor variable, we next consider the importance weights of each variable as the sum of the Δ_p^* for each model in which the predictor variable appears (C

indicated by the set notation code such as {123·} for the full model. The model selection probabilities $\hat{\pi}_i$ are based on 10,000 bootstrap samples.

Model <i>i</i>	Predictors used	AIC _c	Δ _{<i>i</i>}	<i>w_i</i>
11	{123·}	−43.74	0.00	0.454
5	{12··}	−42.69	1.05	0.268
6	{1·3·}	−40.77	2.97	0.103
15	{1234}	−39.89	3.85	0.066
12	{12·4}	−39.07	4.67	0.044
13	{1·34}	−38.15	5.58	0.028
1	{1···}	−38.11	5.62	0.027
8	{1··4}	−35.19	8.55	0.006
14	{·234}	−32.43	11.30	0.002
10	{··34}	−31.59	12.15	0.001
7	{·23·}	−29.52	14.21	0.000
9	{·2·4}	−25.87	17.86	0.000
3	{··3·}	−25.71	18.02	0.000
2	{·2··}	−23.66	20.07	0.000
4	{···4}	−21.70	22.03	0.000
16	{····}	−15.33	28.41	0.000

TABLE 5.13. Evidence for the importance of each regressor variable based on sums of Akaike weights $w_+(j)$ (denoted by $\sum w_i$ here) of models in which variable occurs and based on relative frequency of occurrence of each model based on 10,000 bootstrap samples.

Predictor	Σw_i	$\Sigma \hat{\pi}_i$
x_1	0.997	0.965
x_2	0.834	0.722
x_3	0.654	0.658
x_4	0.147	0.148

computed the relative frequency of models selected containing x_i (hence, the same sum of the bootstrap-based $\hat{\pi}_i$). The results show that at this sample size, x_4 is not important, whereas x_1 is very important. x_2 and x_3 are at least moderately important. Indeed, the most important model was g_{11} , {123·}.

Standard theory gives us the MLEs for each β_j including the estimated conditional standard error $\widehat{\text{se}}(\hat{\beta}_j \mid g_i)$. We give the MLEs and compute the model-averaged estimate of each regression coefficient.

corresponding estimated unconditional standard error $\widehat{\text{se}}$ use as the estimated unconditional standard error of any use the AIC_c -selected best model. The needed formulas are given especially (4.9) with $\hat{\pi}_i = w_i$. The numerical inputs to the given in Table 5.14 for each of the four parameters: Note that a model-averaged estimate of a parameter independently of a single best model. In Table 5.14 we show for each β_j eight models it appears in, and the corresponding conditional and Akaike weight. These weights must be renormalized using (4.1) and (4.9). In Table 5.14 the normalizing constant is the model-averaged estimate and its estimated unconditional standard error that for a given regression coefficient β_j , its estimate can vary by model. Part of this variation is due to model bias, as discussed in 5.3.5. For example, $\hat{\beta}_3$ is 0.18 ($\widehat{\text{se}} = 0.09$) for model g_{11} , model g_3 , $\{\cdot \cdot 3 \cdot\}$, we have $\hat{\beta}_3 = 0.51$ ($\widehat{\text{se}} = 0.12$). It is because are correlated here that β_3 actually varies over the different inclusion/exclusion of other predictor variables), and this contributes to model uncertainty about β_3 under model averaging, of the set of $\hat{\beta}_3$ values in Table 5.14.

We also used the 10,000 bootstrap samples to estimate parameter $\bar{\beta}_j$ and its unconditional sampling variation $\widehat{\text{se}}(\bar{\beta}_j)$. 4. The results from the analytical approach using model selection weights, and the bootstrap results, are shown in Table 5.14. The estimated unconditional standard errors based on Akaike model averaging (i.e., from (4.9)) are each larger than the corresponding standard error for the selected best model, g_{11} , $\{123 \cdot\}$. In the case for β_1 where based on analytical methods $\widehat{\text{se}}(\hat{\beta}_1) = 0.000455$, and from the bootstrap, $\widehat{\text{se}}(\hat{\beta}_1) = 0.000455$, and the other bootstrap results are precise to essentially two decimal places.

The three different point estimates (bootstrap, AIC_c -based, and best model) for each of β_1 , β_2 , and β_3 are quite similar, and the standard errors are quite similar. However, this is not true for the two estimates of β_1 . Whether the bootstrap result (point estimate 0.29, $\widehat{\text{se}} = 0.09$) or the AIC_c -based result (point estimate 0.09, $\widehat{\text{se}} = 0.20$) is better in this example is not clear.

Note also that here the bootstrap-estimated unconditional standard error for $\hat{\beta}_2$ is less than the estimated conditional standard errors for $\hat{\beta}_2$ and $\bar{\beta}_2$. This example constitutes only a sample of size 1 as regards model selection and information-theoretic analytical methods. It will take a Monte Carlo study to make a reliable general comparison of the different methods for assessing model selection uncertainty.

averaged β_j (4.1), its unconditional standard error (Section 4.3.2), weights over the relevant subset of models.

Model	Predictors	Results by model	
i	used	$\hat{\beta}_1$	$\widehat{\text{se}}(\hat{\beta}_1 \mid g_i)$
11	{123·}	0.002185	0.0004553
5	{12··}	0.002606	0.0004432
6	{1·3·}	0.002510	0.0004992
15	{1234}	0.002010	0.0005844
12	{12·4}	0.002586	0.0005631
13	{1·34}	0.002129	0.0006533
1	{1···}	0.003178	0.0004652
8	{1··4}	0.002987	0.0006357
$\hat{\beta}_1$ and unc. se:		0.002368	0.0005350

Model	Predictors	Results by model	
i	used	$\hat{\beta}_2$	$\widehat{\text{se}}(\hat{\beta}_2 \mid g_i)$
11	{123·}	0.001312	0.0005252
5	{12··}	0.001574	0.0005555
14	{·234}	0.001423	0.0007113
15	{1234}	0.001252	0.0005515
12	{12·4}	0.001568	0.0005811
7	{·23·}	0.002032	0.0007627
9	{·2·4}	0.002273	0.0008280
2	{·2··}	0.003063	0.0008367
$\hat{\beta}_2$ and unc. se:		0.001405	0.0005558

5.3.2 A Monte Carlo Extension of the GPA Exam

Simulation is a very useful way to gain insights into co-
tion issues. In particular, here we can assume that the five
 $(y, x_1, x_2, x_3, x_4)'$ is multivariate normal, $MVN(\mu, \Sigma)$ (
generate 10,000 independent simulated sets of data un
model, and do full model selection to learn about select
sues. Given the matrix Σ we can determine the true reg
(and their approximate true conditional standard errors)
regression models. The β_j (and other needed quantities
model depend only on elements of Σ , so it suffices to set
zero: $\mu = \underline{0}$. The needed 5×5 variance–covariance matri
sample variance–covariance matrix from the GPA data. T
will be under an assumed truth that is close enough to

i	used	$\hat{\beta}_3$	$\widehat{\text{se}}(\hat{\beta}_3 g_i)$
11	{123·}	0.1799	0.0877
7	{·23·}	0.3694	0.1186
6	{1·3·}	0.2331	0.0973
15	{1234}	0.1894	0.0919
3	{··3·}	0.5066	0.1236
13	{1·34}	0.2474	0.0990
14	{·234}	0.3405	0.1045
10	{··34}	0.4171	0.1054
$\hat{\beta}_3$ and unc. se:		0.1930	0.0932

Model	Predictors	Results by model	
i	used	$\hat{\beta}_4$	$\widehat{\text{se}}(\hat{\beta}_4 g_i)$
15	{1234}	0.08756	0.1765
12	{12·4}	0.01115	0.1893
13	{1·34}	0.17560	0.1932
8	{1··4}	0.09893	0.2182
14	{·234}	0.45333	0.1824
10	{··34}	0.57902	0.1857
9	{·2·4}	0.51947	0.2207
4	{···4}	0.77896	0.2407
$\hat{\beta}_4$ and unc. se:		0.09024	0.1989

TABLE 5.15. Bootstrap (10,000 samples) and Akaike weight-based example for model-averaged estimated regression coefficients and unconditional standard errors, which include model selection uncertainty. The estimate from the selected model g_{11} and its estimated conditional

j	Bootstrap results		w_i model-averaged		AIC
	$\hat{\beta}_j$	$\widehat{\text{se}}(\hat{\beta}_j)$	$\hat{\beta}_j$	$\widehat{\text{se}}(\hat{\beta}_j)$	$\hat{\beta}_j$
1	0.00236	0.000652	0.00237	0.000535	0.00237
2	0.00156	0.000508	0.00141	0.000556	0.0013
3	0.2296	0.0684	0.1930	0.0932	0.1799
4	0.2938	0.3056	0.0902	0.1989	

these GPA data to provide useful results and insights about the GPA example.

This particular use of Monte Carlo simulation is also called parametric bootstrap: We use as the generating model the param-

the (nonparametric) bootstrap. The advantage of this procedure is that we can specify true values of parameters and hence obtain interval coverage.

Symbolically, the full variance–covariance matrix is partitioned for model g_i :

$$\Sigma_i = \begin{bmatrix} \sigma_y^2 & \underline{c}' \\ \underline{c} & \Sigma_x \end{bmatrix}.$$

The marginal variance of the response variable y is σ_y^2 . If the predictors x_j are in the regression model, the vector \underline{c} (m elements) gives their covariances with y ($\text{cov}(y, x_j)$). The variance of just the predictors considered (i.e., for any of the 15 models, the intercept-only model, $\{\cdot \cdot \cdot\}$) is given by matrix Σ_x (m elements). The true regression parameters, other than the intercept β_0 , is

$$\underline{\beta}' = \underline{c}'(\Sigma_x)^{-1}$$

(we ignore β_0). The approximate sampling variance–covariance matrix is given by $\sigma_{y|x}^2(\Sigma_x)^{-1}$, where

$$\sigma_{y|x}^2 = \sigma_y^2 - \underline{c}'(\Sigma_x)^{-1}\underline{c}$$

is the true residual variance in the regression (a good discussion of multivariate theory is Seber 1984).

The actual simulation process generates the rows of the design matrix random, but then we condition on them in the regression. Conditionally on X , $\hat{\underline{\beta}}$ is unbiased; so it is also unconditionally unbiased, $E(\hat{\underline{\beta}}) = \underline{c}'(\Sigma_x)^{-1}$. This same argument applies to $\hat{\sigma}_{y|x}^2$. Conditionally (by sample), the variance–covariance matrix of $\hat{\underline{\beta}}$ is $E(X'X)^{-1} = (\Sigma_x)^{-1}$ holds only asymptotically as sample size increases. For the simulations, $\sigma_{y|x}^2(\Sigma_x)^{-1}$ is only an approximation to the variance–covariance matrix of $\hat{\underline{\beta}}$.

Rather than show Σ , we show the derived correlation matrix R only:

	x_1	x_2	x_3	x_4
y	0.850	0.653	0.695	0.606
x_1		0.456	0.559	0.663
x_2			0.434	0.417
x_3				0.272

(the ordered x_i are SATmath, SATverbal, HSmath, HSEnglish). The x_i, x_j correlations are so high that we would need to eliminate

as g_{15} (as detailed in the text).

Model $g_{15}, \{1234\}: \sigma_{y x}^2 = 0.05692$				
j	β_j	$\text{se}(\hat{\beta}_j g_{15})$	$\text{cv}(\hat{\beta}_j g_{15})$	
1	0.002010	0.0005061	0.252	
2	0.001252	0.0004776	0.381	
3	0.1895	0.0796	0.420	
4	0.0875	0.1528	1.745	
Model $g_{11}, \{123 \cdot\}: \sigma_{y x}^2 = 0.05785$				
j	β_j	$\text{se}(\hat{\beta}_j g_{11})$	$\text{cv}(\hat{\beta}_j g_{11})$	
1	0.002185	0.0004072	0.186	
2	0.001312	0.0004698	0.358	
3	0.1799	0.0784	0.436	
Model $g_5, \{12 \cdot \cdot\}: \sigma_{y x}^2 = 0.07307$				
j	β_j	$\text{se}(\hat{\beta}_j g_5)$	$\text{cv}(\hat{\beta}_j g_5)$	
1	0.002606	0.0004086	0.157	
2	0.001574	0.0005121	0.325	
Model $g_1, \{1 \cdot \cdot \cdot\}: \sigma_{y x}^2 = 0.10759$				
j	β_j	$\text{se}(\hat{\beta}_j g_1)$	$\text{cv}(\hat{\beta}_j g_1)$	
1	0.003178	0.0004413	0.139	
Model $g_3, \{\cdot \cdot 3 \cdot\}: \sigma_{y x}^2 = 0.20001$				
j	β_j	$\text{se}(\hat{\beta}_j g_3)$	$\text{cv}(\hat{\beta}_j g_3)$	
3	0.5066	0.1173	0.117	

is an important consideration. In addition, a principal-component analysis of the covariance matrix of the predictor variables is a reasonable choice.

The data-generating model used here ($g_{15}, \{1234\}$), which generated the data, has the pairwise correlations given above and the conditional variance $\sigma_{y|x}^2 = 0.05692$ (from (5.5)). From (5.4) we compute the conditional standard errors of the $\hat{\beta}_j$ under model g_{15} . We also compute the conditional coefficient of variation of each $\hat{\beta}_j$ (Table 5.16). These quantities are computable for any submodel fitted to the generated data. Table 5.16 shows these theoretical values for models $\{1234\}$, $\{123 \cdot\}$, $\{12 \cdot \cdot\}$, and $\{\cdot \cdot 3 \cdot\}$.

Table 5.16 if we use model g_3 , then $E(\hat{\beta}_3) = 0.5000$ (cv = 0.436). For model g_{11} , $E(\hat{\beta}_3) = 0.1799$ (cv = 0.436). Results when the specified model is always fit to the data; hence no model selection occurs. When the inference strategy is to first select a model on the data, then use it for inference, the properties of the inferences are affected (model selection biases and uncertainty).

Examination of results in Table 5.16 demonstrates that the variance measured by its coefficient of variation, precision of a parameter estimate, as the number of other parameters in the models decreases. The variance varies by model is because the predictors are correlated. The variation in $E(\hat{\beta}_j)$ will get more pronounced if correlations among predictors do not occur if all predictors are uncorrelated with each other. The effect of fitted model g_i has little effect if that predictor is unimportant, as by a large $cv(\hat{\beta}_j | g_i)$. For example, x_4 (HSenglish) can be dropped from model {1234}; hence one uses {123·}. Indeed, in this Monte Carlo study the expected K-L best model is g_{11} (based on 10,000 simulations).

What we want to illustrate with simulation here (and in the next section) are some results under all-subsets AIC_c model selection in regression. We focus on unconditional vs. conditional confidence interval coverage, i.e., the value of β_j in the generating model g_{15} in Table 5.16. For the intervals we used $\hat{\beta} \pm 2 \hat{se}$; hence, we ignored the issue of a coverage multiplier. This affects coverage a little, but the focus is really on the difference between conditional and unconditional coverage.

Second, we look at induced model selection bias in $\hat{\sigma}_{y|x}^2$. For model g_i , out of the 16 models fitted (especially at a small sample size), AIC_c tends to have a better fit for that data set, hence a smaller sum of squares, than would occur on average if model g_{15} were selected. Thus, data-based selection in regression will tend to result in a model that fits a little too well; as a result, we get $E(\hat{\sigma}_{y|x}^2) < \sigma_{y|x}^2$ (the bias on confidence interval coverage that is not correctable by using multiplier intervals). Confidence intervals depend on $\hat{\sigma}_{y|x}$; however, $\sigma_{y|x}$ varies by model, so what we report to assess selection bias is

$$RB = \frac{E(\hat{\sigma}_{y|x}) - \sigma_{y|x}}{\sigma_{y|x}}.$$

Other quantities of interest include the expected value of the estimator $\hat{\beta}_j$ and unconditional interval coverage based on model selection variation, and percentiles of Δ_p . From the

50	0.9
80	3.1
90	3.6
95	5.1
98	7.3
99	9.4

When a parameter (hence $\hat{\beta}_j$) appeared in the selected several quantities: the model-averaged estimate $\hat{\hat{\beta}}_j$, the unconditional standard error $\widehat{\text{se}}(\hat{\beta}_j)$, and three confidence intervals (nominally 95% conditional interval (cond.) is based on $\hat{\beta}_j$ and its estimated conditional standard error given the selected model. The unconditional interval (uncond.) is based on $\hat{\hat{\beta}}_j$ and $\widehat{\text{se}}(\hat{\beta}_j)$. The interval based on model averaging (MA) is based on $\hat{\hat{\beta}}_j$ and $\widehat{\text{se}}(\hat{\hat{\beta}}_j)$. Finally, it needs to be clearly understood that the coverage of the unconditional interval for the true parameter from the actual data-generating model is not 95%.

In this example, the achieved coverage of the unconditional interval is 92.1%, less than that of the conditional intervals (Table 5.17), especially for β_1 (90.8% vs. 80%; both coverage percentages increase by about 0.02 if a model-averaged interval is used). A source of lowered coverage comes from model selection bias. For example, here model selection results in the unconditional standard error, which is here $\text{se}(\hat{\beta}_1) = 0.00229$ and the conditional standard error is $0.00201 - 0.00201 = 0.00028$. This bias is important on the unconditional standard error, which is here $\text{se}(\hat{\beta}_1) = 0.00229$. The bias/se ratio is $\delta = 0.54$; this value of δ will result in a coverage of 92.1% if coverage would be 95% at $\delta = 0$ (see Cochran 1957).

An unexpected result in this example is that the conditional coverage for β_2 and β_3 was as high as (about) 0.9. The unconditional coverage for β_2 and β_3 improve and without exceeding 95% coverage.

The model selection bias induced in $\hat{\sigma}_{y|x}$ is negative, i.e., the estimated probability that the model will be selected (Table 5.18). Since the true probability is 0.5, the bias is -0.05.

TABLE 5.17. Expected values of estimators of β_j and confidence intervals for the true β_j under AIC_c-based all-subsets model selection from the Monte Carlo simulation (10,000 samples) mimicking the GPA example; occurrence frequency of the 10,000 samples in which the selected model included the indicated β_j .

Occur. freq.	j	β_j	$E(\hat{\beta}_j)$	$E(\hat{\hat{\beta}}_j)$	Achieved
					cond.
9544	1	0.00201	0.00229	0.00232	0.801
7506	2	0.00125	0.00156	0.00157	0.906
6506	3	0.190	0.248	0.253	0.887

of the 10,000 values of $\hat{\sigma}_{y|x}^2$ when the model is fit to every generated set (occurs).

Model i	Predictors used	$E(\hat{\sigma}_{y x}^2 \mid g_i)$ no selection	π_i	S
11	{123·}	0.0576	0.3786	—
5	{12··}	0.0728	0.2730	—
6	{1·3·}	0.0803	0.1451	—
13	{1·34}	0.0762	0.0458	—
15	{1234}	0.0566	0.0389	—
1	{1···}	0.1076	0.0351	—
12	{12·4}	0.0728	0.0301	—
14	{·234}	0.1017	0.0216	—
10	{··34}	0.1272	0.0161	—
8	{1··4}	0.1063	0.0078	—
7	{·23·}	0.1420	0.0064	—
9	{·2·4}	0.1695	0.0012	—
2	{·2··}	0.2227	0.0001	—
3	{··3·}	0.2012	0.0001	—
4	{···4}	0.2446	0.0001	—

is the average value of $\hat{\sigma}_{y|x}^2$ when the model is fit to all 10 sets (compare to theoretical results in Table 5.16). The RB under model selection is given by RB from (5.6). Good models (paradigm) do not correspond to very bad levels of RB. As less acceptable (in terms of expected K-L value), it is selected if the data are an unusually good fit to that model.

While general in their qualitative nature, these numerical results are extreme for a sample size of 20 than would be the case at $n = 50$. In fact, for this generating model the selection bias in $\hat{\sigma}_{y|x}^2$ at sample size $n = 50$, and confidence interval coverage is nearly 90% interval. Even though model selection can induce biases in the theoretic selection and associated unconditional inference, the results are good and certainly better (for the sample size) than using a high-dimensional global model that includes all predictors.

5.3.3 An Improved Set of GPA Prediction Models

An even better way to improve on the all-predictors global model is to reduce one's models to a smaller set of a priori meaningful predictors by subject matter or logical considerations. Basically, this involves transformations of the predictors (consideration of meaningful interactions also important) and dropping predictors that are very unlikely to be

different SAT prediction models based on simple derived variables in some sense.

The original four predictors are each just indices to general ability and they are measured with error. That is, a person's test score varies by circumstances (and luck), such as if they had a good or bad exam. Viewing these predictor variables as just semicircular random variables, just compute a single averaged index? In so doing we average out the verbal (English) ability, but grades in many courses depend on many other things anyway. With a large sample size (say $n > 1,000$) it may be possible to use regression fit to calibrate the relative importance of the four predictors, but with only 20 observations some combining of indices may be necessary.

The SAT and HS scores are on very different scales. There are many ways to allow for this, such as first to normalize each predictor to have a mean of 0 and a standard deviation of 1 and then to combine the adjusted predictors to get a total (*tot*) predictor index. To control for the nuisance we used geometric means to cope with the scale differences. The new variables that replace the original four variables are

$$\begin{aligned} sat &= (x_1 \times x_2)^{0.5}, \\ hs &= (x_3 \times x_4)^{0.5}, \\ math &= (x_1 \times x_3)^{0.5}, \\ engl &= (x_2 \times x_4)^{0.5}, \\ tot &= (sat \times hs)^{0.5} = (math \times engl)^{0.5} = (x_1 \times x_2 \times x_3 \times x_4)^{0.25}. \end{aligned}$$

These variables are interpretable and seem just as adequate as the original variables as indices to first-year college GPA.

Next, we would not use, in an a priori analysis, any of the models that use the derived variables. The only linear regression models we would (did) consider are those where the predictors are given below, in terms of predictor variables. The models have an intercept and σ^2 . We numbered these as g_{17} through g_{21} in order to compare results to the original 16 models:

model	K	variables included
g_{17}	3	<i>tot</i>
g_{18}	4	<i>sat</i> <i>hs</i>
g_{19}	4	<i>math</i> <i>engl</i>
g_{20}	3	<i>sat</i>
g_{21}	3	<i>hs</i>

We conceptualized these models before examining fit of the models to the GPA data, and no other derived models were considered.

The AIC_c best model of the above is g_{17} (Table 5.19). It is the best in the full set of all 21 models (results are not shown in Table 5.19). In adding new models to an existing set, no

models originally considered (see text for details of the new prediction

Model i	Predictors used	AIC _c	Δ_i	w_i
17	<i>tot</i>	-48.20	0.00	0.3
18	<i>sat hs</i>	-45.97	2.23	0.1
20	<i>sat</i>	-43.96	4.24	0.0
11	{123·}	-43.74	4.47	0.0
5	{12··}	-42.69	5.52	0.0
6	{1·3·}	-40.77	7.44	0.0
15	{1234}	-39.89	8.32	0.0
19	<i>math engl</i>	-39.23	8.97	0.0
...
21	<i>hs</i>	-33.98	14.22	0.0
...

need to be recomputed: Just reorder the full set from sma
The full set of Δ_i values may need to be recomputed if the
(as here, from g_{11} to g_{17}). Given the new set of Δ_i , re
weights w_i .

The results in Table 5.19 illustrate that the best model
relative only to that set of models. Kullback–Leibler mod
provide an absolute measure of how good a fitted model
only in the set of 16 all-subsets models. Compared to m
(and the entire original set of 16 models) can almost be dis
for expected K-L best model for these data. Correspondi
that any model-based inference is conditional on the set of
The specifics of inferences and computable uncertainties
the models formally considered.

The models used here are useful only for prediction;
any causal process. Hence, we illustrate inclusion of mo
prediction based on models g_{17} – g_{21} (standard aspects of
given a fitted linear model are assumed here; see, e.g., Gra
As computed in Graybill and Iyer (1994:244), under mo
prediction of expected GPA at $x_1 = 730$, $x_2 = 570$, $x_3 = 2.7$ is $\hat{E}(y) = 3.185$ with conditional (on model) stand
comparison we note that given model g_{11} , {123·}, the co
are $\hat{E}(y) = 3.253$, $\hat{se} = 0.102$.

The model-averaged predicted expected GPA is $\hat{\bar{E}}(y) =$
unconditional standard error of 0.11 (Table 5.20). These r
using (4.1) and (4.9). To construct a confidence interval he
small degrees of freedom of $\hat{\sigma}_{y|x}^2$ we suggest that it suffic

Model	Predictors				
i	used	Δ_i	w_i	$\hat{E}(y)$	\widehat{se}
17	<i>tot</i>	0.00	0.685	3.016	
18	<i>sat hs</i>	2.23	0.224	3.095	
20	<i>sat</i>	4.24	0.082	3.177	
19	<i>math engl</i>	8.97	0.008	3.271	
21	<i>hs</i>	14.22	0.001	2.632	
Weighted results:				3.056	

use $3.06 \pm t \times \widehat{se}(\hat{E}(y))$, where the multiplier $t = 2.10$ is from the t -distribution on 18 df. Model g_{17} has 18 df for $\hat{\sigma}_{y|x}^2$, and the weight on it is 0.685 (more sophisticated procedures will not make a difference here). The model-averaged result is distinctly more precise than the result based on the fitted global model (standard errors of 0.108 vs. 0.074, model-averaged versus global model-based). Also, the inclusion of g_{21} increases the standard error as compared to the result conditional on g_{17} ($\widehat{se} = 0.074$).

5.3.4 More Monte Carlo Results

The theory for Monte Carlo generation of regression data and the predictors, was presented in Section 5.3.2. Using that approach, we ran a few more simulations. Our motivation was firstly to see what biases occurred in using model averaging and unconditional confidence intervals (if they did not), and secondly to see what biases might result from model selection and what confidence interval coverage could be achieved.

This is far from a full-scale simulation study because of the many factors to consider in the design of an all-subset simulation study. For example, we used only $m = 4$ predictors here (many other values of m need to be explored). Thus, the global model is not included. Given a sample size n , one generates a sample from the $\text{MVN}(\mu, \Sigma)$ generating model. Without loss of generality, we can assume $\mu = 0$. However, there are still (in general) $(m + 1) \times (m + 2)/2$ parameters in Σ . To make this design problem tractable we used the following approach on the generating model: either

$$\underline{c}' = [0 \ 0 \ 0 \ 0] \quad \text{or} \quad \underline{c}' = [0.8 \ 0.6 \ 0.4 \ 0]$$

$$\Sigma_x = \begin{bmatrix} \rho & 1 & \rho & \rho \\ \rho & \rho & 1 & \rho \\ \rho & \rho & \rho & 1 \end{bmatrix}.$$

The regression parameters of the generating model are given by

Another design factor was taken to be the regression parameter \underline{c} . Given values for $\sigma_{y|x}^2$ (we used only the marginal variance of y , σ_y^2). These quantities (i.e., \underline{c} , Σ_x) compute the full 5×5 variance–covariance matrix of $(\hat{\beta}_1, \dots, \hat{\beta}_5)$. The variance–covariance of $\hat{\beta}_1$ is given by

$$\frac{\sigma_{y|x}^2}{n} (\Sigma_x)^{-1}.$$

We see that factors $\sigma_{y|x}^2$ and n are redundant in their effect on the conditional variance. Therefore, for not-small sample sizes it is much easier to fix n (say $n \geq 28 + K$, for the global model value of K) and vary $\sigma_{y|x}^2$ to gain precision, rather than to fix $\sigma_{y|x}^2$ and simulate by increasing sample size. We did not do so here; it is still of interest to use small actual n to explore small-sample-size effects.

Because of the choice of the form of Σ_x , the estimators $\hat{\beta}_i$ all have the same conditional variance. In fact, for any ρ and \underline{c} , the constant diagonal element v of $(\Sigma_x)^{-1}$, because

$$\text{var}(\hat{\beta}_i | g_{15}) = v \times \sigma_{y|x}^2 / n.$$

For $\underline{c}' = [0.8, 0.6, 0.4, 0.2]$ we simulated 10,000 samples for each combination of $\rho = 0, 0.2, 0.4, 0.6$, and 0.8 crossed with $n = 28, 32, 36, 40$. We focused on $\hat{\beta}_1$, but looked at other parameters in a few cases (in another set of 10,000 samples). To these cases 1 to 5 for $\rho = 0$ and 6 to 7, as noted in Table 5.21. In total we looked at 29 simulations.

We tabulated some basic results (Table 5.22) wherein the selection was applied to all 16 possible models (the labeling is the same as in Table 5.12). In particular, we tabulated the 5th, 50th, and 95th percentiles of Δ_p . There is one variation here; with no nonboundary effect than the generating global model (g_{15}) if sample size gets large, Δ_p converges on model g_{15} , and all percentiles of Δ_p go to the same value (i.e., some degree of reduction in percentiles of Δ_p occurs in real data analysis, so we flagged populations where a boundary effect is occurring. Our recommendations about interpreting the results: no boundary effect occurs. For the 15 populations where a boundary effect occurred, the mean percentiles of Δ_p in Table 5.22 are 4.7th, 50th, and 95th, respectively).

Case	ρ	n -set	$\sigma_{y x}^2$	v	β_1	β_2	
1	0	a	1	1.00	0.8	0.6	
2	0.2	a	1	1.09	0.69	0.44	
3	0.4	a	1	1.36	0.73	0.39	
4	0.6	a	1	1.96	0.93	0.43	−0.01
5	0.8	a	1	3.82	1.65	0.65	−0.01
6	0	a	1	1.00	0	0	0
7	0	b	25	1.00	0.8	0.6	

Sample size also has an effect on the distribution of Δ effect for n greater than about 20. For case 7 in Table 5.22 size ($20 \leq n \leq 500$) is about 2 units at the 90th and 95th p 3 units at the 99th percentile. These are typical of sample observed.

There is considerable model selection uncertainty in populations (Table 5.22), for example, as indexed by ho probability π_k is even for the expected AIC_c best model model selection uncertainty used in Table 5.22 is simply a of models, of the 16, that have selection probabilities \geq half of the possible models). For the all too typical app selection with 10 or more variables ($R \geq 1,024$) and no the selection process will be highly unstable (cf. Breim model is selected. That is, selection probabilities can be low, even for the actual K-L best model, and not exhibit result the selected model itself is not at all the basis for about the relative importance of the predictor variables, model provides reliable predictions.

A confidence interval for a parameter β_i was compu parameter was in the selected model, in which case point were the MLE $\hat{\beta}_i$ and the model-averaged $\hat{\hat{\beta}}_i$. Three typ computed: the classical conditional interval $\hat{\beta}_i \pm 2 \widehat{se}(\hat{\beta}_i | g$ unconditional interval $\hat{\beta}_i \pm 2 \widehat{se}(\hat{\hat{\beta}}_i)$, and the interval b averaged point estimate $\hat{\hat{\beta}}_i \pm 2 \widehat{se}(\hat{\hat{\beta}}_i)$. One result was th interval types two and three was barely different, but wa the interval $\hat{\hat{\beta}}_i \pm 2 \widehat{se}(\hat{\hat{\beta}}_i)$ (ratio of coverages: 0.995). Henc only for this latter interval (Table 5.23). We focused on (i.e., parameter index 1); Table 5.23 gives coverage result of a difference parameter (β_3 or β_4).

column six is the number of models for which $\pi_i \geq 0.01$, also given.

Case	ρ	n	Best	π_k	# π_i	Percentiles	
			Model	best	≥ 0.01	90	95
1	0	20	11	0.243	13	5.37	6.74
1	0	50	15	0.368	4	2.55	2.59
1	0	50	15	0.378	4	2.56	2.59
1	0	100	15	0.668	2	1.81	2.12
2	0.2	20	5	0.342	12	4.65	6.51
2	0.2	20	5	0.334	12	4.58	6.41
2	0.2	50	11	0.264	7	2.67	3.78
2	0.2	100	11	0.462	4	2.22	3.55
3	0.4	20	1	0.398	12	7.03	9.35
3	0.4	50	12	0.353	8	3.35	4.29
3	0.4	100	12	0.611	6	2.17	3.36
4	0.6	20	8	0.227	12	6.38	8.73
4	0.6	50	12	0.544	8	2.47	3.65
4	0.6	100	12	0.743	4	1.74	2.89
5	0.8	20	8	0.499	10	5.43	7.68
5	0.8	50	15	0.305	4	3.32	4.03
5	0.8	50	15	0.308	4	3.24	4.03
5	0.8	100	15	0.595	4	2.04	2.23
5	0.8	100	15	0.597	4	2.03	2.23
6	0	20	16	0.594	11	3.43	5.11
6	0	50	16	0.535	11	3.35	4.96
6	0	100	16	0.518	11	3.41	4.90
7	0	20	16	0.493	11	4.55	6.47
7	0	50	1	0.180	13	5.62	7.48
7	0	100	5	0.131	16	5.11	6.90
7	0	200	11	0.143	15	4.08	5.20
7	0	200	11	0.140	15	4.10	5.20
7	0	500	11	0.391	8	3.42	4.93
7	0	1000	15	0.380	4	2.00	2.04

Confidence interval coverage is affected by bias in either β_i or its standard error estimator. Therefore, we tabulated intervals and biases for the interval based on the model-averaged estimator. In this context bias is important only in relation to standard error. Table 5.23 shows,

$$\delta = \frac{E(\hat{\beta}_i) - \beta_i}{E(\widehat{\text{se}}(\hat{\beta}_i))},$$

coverage is trivial for $|\beta_i| \leq 0.25$ and ignorable for $|\beta_i| \geq 0.75$. Hence we show the ratio

$$se-r = \frac{E\left(\widehat{se}\left(\widehat{\beta}_i\right)\right)}{MC-se\left(\widehat{\beta}_i\right)}$$

in Table 5.23. Here, $MC-se(\widehat{\beta}_i)$ is the actual achieved standard error from the Monte Carlo samples (out of 10,000) wherein $\widehat{\beta}_i$ is chosen. A $se-r = 1$ is desirable. Coverage would be reduced (other than 0.95) if $se-r$ becomes much less than 1. The effect on coverage is small if $0.9 \leq se-r \leq 1.1$. For all these results, the relevant sample size is “Freq.” in Table 5.23: the number of samples wherein $\widehat{\beta}_i$ is chosen. The parameter is in the AIC_c -selected model.

The biggest surprise was the high achieved coverage of the conditional confidence interval (Table 5.23). When that coverage is less than 0.95, for example, case 2, $\rho = 0.2$, $n = 20$, unconditional coverage is less than 0.95. It was because of severe bias in either the point estimator or the conditional estimator. Moreover, these biases are clearly a form of model selection bias and they occurred when the reference parameter was infrequently selected (infrequent selection of any model containing β_i), which is due to the predictor variable x_i being unimportant at the given sample size. In this matter, however, not all instances of small frequency of selection including β_i resulted in deleterious effects on coverage (e.g., $\rho = 0.8$, $n = 50$, unconditional coverage of 0.925 on β_3). Coverage is good if a predictor variable was important (as judged by high ρ), and its unconditional (and conditional) coverage was always good.

In all 29 simulated populations the unconditional interval coverage was greater than or equal to the conditional coverage and provided good coverage when the conditional coverage was less than 0.95. For example, $\rho = 0.8$, $n = 20$, conditional and unconditional coverage were 0.937 and 0.937 on β_1 . For the 23 populations where the parameter was not selected to the model, the average conditional and unconditional coverage was 0.930 versus 0.947 (and $\bar{\delta} = 0.2$). The improvement in coverage is not dramatic, but is generally good.

For the other six populations, the bias in coverage is small. There are no strong biases in point estimates or their standard errors, and $se-r$ ($abs(\delta) = 0.64$ and $se-r = 0.63$). These biases are a direct result of model selection, i.e., they are model selection bias. In those cases where the selected model is not the expected AIC_c best (i.e., expected K-L best), it includes a variable x_s that is rarely included in the selected model.

of such selected models is denoted by \hat{r} . Freq., coverage is for the interval and the interval based on the model-averaged estimator $\hat{\beta}_i$; of δ and se- r .

Case	ρ	n	Parm. index	Freq.	Coverage	
					cond.	MA
1	0	20	1	8864	0.932	0.948
1	0	50	1	9994	0.943	0.948
1	0	50	4	4337	0.934	0.934
1	0	100	1	10000	0.948	0.950
2	0.2	20	1	8309	0.923	0.947
2	0.2	20	3	2106	0.694	0.755
2	0.2	50	1	9970	0.945	0.955
2	0.2	100	1	10000	0.944	0.950
3	0.4	20	1	8026	0.908	0.945
3	0.4	50	1	9912	0.918	0.947
3	0.4	100	1	9999	0.932	0.944
4	0.6	20	1	8178	0.905	0.957
4	0.6	50	1	9953	0.905	0.942
4	0.6	100	1	10000	0.937	0.945
5	0.8	20	1	9175	0.858	0.937
5	0.8	50	1	10000	0.918	0.945
5	0.8	50	3	3572	0.913	0.925
5	0.8	100	1	10000	0.930	0.946
5	0.8	100	3	6055	0.957	0.960
6	0	20	1	1220	0.396	0.495
6	0	50	1	1450	0.643	0.677
6	0	100	1	1501	0.658	0.676
7	0	20	1	2014	0.726	0.759
7	0	50	1	3624	0.904	0.913
7	0	100	1	5594	0.952	0.955
7	0	200	1	7911	0.967	0.969
7	0	200	4	2225	0.785	0.796
7	0	500	1	9828	0.969	0.969
7	0	1000	1	9998	0.956	0.956

inference on β_s (an unimportant variable) can be very much affected by the resultant model selection bias for $\hat{\beta}_s$. Fortunately, this selection bias is almost by definition, since it is a case of x_s being common to all models in the selected model. Also, even then inference on an important variable is generally sound in these simulations and others we have conducted. The inference on an unimportant β_s was slightly (but not strongly) affected. The unconditional confidence interval was considered.

reference to linear regression and so-called variable selection, the most used and misused application area of model selection. Every conceivable type of model (variable) selection method has been tried in the context of having m predictors and using linear regression models (see, for example, Hocking 1976, Draper and Henderson and Velleman 1981, Breiman and Freedman 1983, Miller 1990, Hjorth 1994, Breiman 1995, Tibshirani 1996, etc.). However, almost always the statistical literature approaches the problem as if it is only a matter of “just-the-numbers” data analysis method. In fact there is always a subject-matter scientific context, and a choice of sample size, that must be brought into the problem, and an enormous difference as compared to any naive model selection method that does not consider context, prior knowledge, and sample size.

Secondly, there is always a goal of either (1) selection of the best model (should include ranking competitor models) because one of the relationships (presumably causal) between \underline{x} (independent variable) and y , or (2) prediction of $E(y | \underline{x})$ at values of \underline{x} (predictor). (prediction of $E(y | \underline{x})$ for \underline{x} in the sample can be considered as a form of estimation). These goals really are different. That is, if there is model selection uncertainty, then selecting the best model for goal (1), and using it for goal (2), prediction, is not optimal.

We recommend that prior to any data analysis full consideration be given to how the problem (i.e., set of models) should be structured. This means dropping variables that cannot reasonably be related to the response, or cannot reasonably be causally related at detectable levels, given the sample size. From the literature and our experience, statisticians are too reluctant to drop clearly irrelevant variables and other considerations based on logic and theory. This is the “measure what is easy to measure and let the computer sort it out” syndrome. Even a good exploratory analysis needs input of informed theory to reach useful results. In this regard we quote Freedman et al. (1991):

A major part of the problem in applications is the curse of dimensionality. There is a lot of room in high-dimensional space. That is, there is a need for model specifications tightly derived from good theory. We expect statistical modeling to perform at all well in an era characterized by consisting of large, complicated data sets and weak theoretical foundations. At present that describes many applications.

An important a priori aspect is to consider reducing the number of independent variables by functionally combining them into a small number of variables. This may be as simple as computing, by observation, that some of the predictor variables (such as in the GPA exam

density as mass per unit volume in physics or animals per rates of all sorts, and so forth). Consider also any bounds on often college GPA is bounded on 0 to 4; hence, we do not be able to make a prediction of 4.2. We could model GPA with a link function and rescale predictions by 4.

As a rule of thumb, the maximum number of predictors to allow in a regression (or other univariate) model is not possible to reliably estimate anything like $n/2$ (or $n/3$) “noisy” data. Mistakenly, models of such size are often fitted selected based on an invalid criterion such as minimum variance, an inappropriate to the situation criterion like adjusted R^2 does not correctly adjusts for either small sample size or large K).

To illustrate some of this thinking we consider another example of variable selection (Hocking 1976; see also Hocking 1981 for the actual data; we did not read this latter paper, but below): automobile gas mileage (y) as MPG (miles per gallon) as a function of 10 independent variables (there is clearly causation involved, but not so we use the term predictor variable). Note that y is a response variable. The 10 x_i are:

- 1 Engine shape (straight or V)
- 2 Number of cylinders (4, 6, or 8)
- 3 Transmission type (manual or automatic)
- 4 Number of transmission speeds
- 5 Engine size (cubic inches)
- 6 Horsepower
- 7 Number of carburetor barrels
- 8 Final drive ratio
- 9 Weight
- 10 Quarter mile time

The data arose from testing 32 ($= n$) different types of cars under standardized conditions. We independently generated a predictor variable x_i (it would be better to get an automotive engineer involved). We then fit a model at the dependent variable (either as y alone, or as y versus x_i), given that one has decided that the analysis will be only to predict a response variable y based on x_1, \dots, x_m as predictor variables. If the models are all conditional, one next reduces the number of predictors as much as possible by logical and subject-matter considerations. From the resultant reduced set of predictors, we recommend looking at the correlation of those independent variables to be assured that there are no pairs having an extremely high pairwise correlation.

A more comprehensive examination would be a principal component value evaluation of the design matrix X . Such results are available (Hocking 1976) for the full set of 10 predictors (but the data are not

eigenvalues add to about 98% of the total of all 16 eigenvalues. The first 10 eigenvalues explain 96% of the explained variation is often all we can hope for without a sample size 32.

The x_i are highly intercorrelated in this observational study. Design: Big cars have bigger engines; are more likely to have 6 cylinders; are therefore more likely to have a V-engine design. For such observational studies if the issue of interest is causal, there are substantial inference problems (see Draper and Smith 1999). Some sage cautionary comments on such problems with observational data.

One of us (KPB, who is automotively challenged) proposes. Because $n = 32$, do not include more than three structural variables. Mileage is strongly dependent on car weight, so always include an intercept will be used here, this leaves room for only one other variable. As a first thought, then, consider the nine models

$$y = \beta_0 + \beta_1 x_9 + \beta_i x_i + \epsilon, \quad i \neq 9$$

plus $y = \beta_0 + \beta_1 x_9 + \epsilon$. However, bearing in mind the intercorrelation of these variables, consider dropping some on a priori grounds. Dropping x_i based just on a high correlation unless it is extreme such as $|r| > 0.9$ then there is a variable redundancy problem (near collinearity) (near) collinearity problems; and do eliminate variables based on statistical reasoning, and experience.

As a type of thought experiment (because the data do not support a causal experiment) consider whether engine shape (x_1) is causal. Do we really think that if all car features were held fixed except engine shape, engine is a straight or V8 that there would be any effect on mileage? Either not at all, or at a level we will never care about a statistically significant effect except with an experiment and a huge sample size. Conclusion: drop x_1 (we surmise that it was recorded because it is easy to record, not justification for including a variable). **Recommendation: drop experiments in conjunction with observational studies**

The same reasoning leads KPB to drop variable 2: number of cylinders. Again the thought is that if all else (horsepower, total car weight, etc.) were fixed, would just number of cylinders (as 4, 6, 8) affect MPG? And again, no; at least not in these data. Variable 10 might be retained on fundamental grounds. Quarter mile time is a derived variable; it might predict MPG well, but it is not a structural variable. Instead, variable 10 might itself be well predicted based on other variables. Conclusion: drop x_{10} . Thus KPB would consider only models with x_9 always included and at most one of variables 3–8 (seven models). This is different from an all-subsets selection over $10^m = 1,024$ models. The problem of thinking can, and should, always be brought to bear on model selection problem.

that horsepower itself is a derived variable but one that can to some extent horsepower might replace variables 5 and 7. More important than variables 3, 4, 6, and 8. Thus we have the impression that the best two-variable model might be based on x_9 and x_8 but not best, but nearly tied for best with a less-interpretive prior thinking would justify objective selection of these two as most important.

Another issue is the suitability of the linear model for a response bounded below by zero but weight can be unbounded, and can predict negative MPG. Surely, over a big enough weight range, the model is curvilinear, such as $E(y | x_9) = \beta_0 e^{-\beta_1 x_9}$ or $E(1/y | x_9) = \beta_0 + \beta_1 x_9$ for MPG only by convention; hence a priori KPB would fit the data better but based on an inverse link function to MPG. One could just fit several linear models and select a best model. However, AIC_c is a criterion for models fit to y versus models fit to $1/y$.

Less time was afforded to this exercise by DRA (who was not particularly motivatedly knowledgeable), who independently put forth two models. Both include weight (x_9);

$$y = \beta_0 + \beta_1(x_9)^2 + \beta_2 z + \epsilon$$

and

$$y = \beta_0 + \beta_1(x_9)^2 + \beta_2 x_{10} + \epsilon.$$

The variable z is a derived variable meant to reflect the effect of several variables on MPG:

$$z = \frac{x_2 \times x_5}{x_6}.$$

Similarly, x_{10} is used here as a predictor that summarizes the effect of several variables on the ability of the car to consume fuel.

Considerations like these based on reasoning and theory should be made about before data analysis if reliable uncertainty bounds are to be made on an inference made after model selection. One can always do exploratory analyses after the a priori analysis. We just need to separate the two processes, because results of exploratory analyses are not reliable inferences in the sense that the data cannot both confirm the model (the model) and then reliably affirm the inferential uncertainty (the same model). There is a saying from the USA western: "ask questions first, then create models later." This strategy often precludes confirmation. Similarly, "compute first, then create models later" (or "think later") is also not a strategy for making reliable inferences. It can be a model that describes the data very well (because it confirms the data) but is a poor model as an inference to independent data from

5.3.6 Discussion

A variety of comments and opinions are given here, some fit well elsewhere in Section 5.3. There is no particular order to the comments and opinions.

The Monte Carlo simulations of Sections 5.3.2 and 5.3.3 violate our general philosophy that the actual data-generating process (should) in reality be more general than the global model used for data analysis (an expanded vector of predictors \underline{x}_T would be the true generating model; the global model does not use all the predictors). This is only partly true. The part that is not true is thinking that the data generated under the global model, no more general than the global model, can be used to apply the global model. In fact, the residual variation of the global regression model is a confounding of average (with respect to f) model variation arising from the differences $E(y | \underline{x}, G) - E(y | \underline{x}_T, f)$, plus the unexplained residual variation σ_ϵ^2 (it might be 0) under f . Thus, the use of a conceptually more general data-generating model and the global model. Hence, it is more economical simply to use the assumed global model.

One way in which this lacks generality is that the number of components of the true parameter vector $\underline{\beta}_T$ that are in the global model $\underline{\beta}$ may not exactly equal their counterpart components. This would affect confidence interval coverage, which should be based on the appropriate components of $\underline{\beta}_T$, not to $\underline{\beta}$. This seems like a problem in initial Monte Carlo studies intended to explore basic model properties.

The more important lack of generality relates to how we select a model sequence of models as sample size increases. Classical model selection is a model or set of models fixed, independent of sample size. This is not true with reality, because as sample size grows we will include more predictors in the data and in our models (e.g., in the GPA example, we might include school type, university attended, major, student age, and so on). The reason we should simply have a larger global generating model is that we used in simulations here, and include more factors so that the selected model can grow without the arbitrary bound of only four predictors. In a sense the issue becomes one of "growing" up against a bound (i.e., large π_i for the generating model). As sample size increases, because this feature of data analysis is often uncontrolled in experimental studies. This problem is solved simply by having a larger global generating model, and it is then still acceptable (and more general) to have that generating model also as the global model for data analysis.

It is well known that selecting a best model from a set of models can lead to important biases in parameter estimates and in model predictions.

selecting a good model, because an selection criteria that is seeking a fitted model with a relatively small residual sum of squares, AIC_c will not entirely protect one from this possible bias (e.g., Table 5.18; the relative bias in $\hat{\sigma}_{y|x}^2$ is not high for this example and this example is for a small sample size). The bias is more infrequently a model is selected because for such models are selected only when they fit a sample unusually well. In the context of a study (hence, sample size is then a given) the best way to reduce bias in $\hat{\sigma}_{y|x}^2$ is to keep the candidate set of models small.

When the predictors are intercorrelated and model selection tends to induce a bias in the estimators of regression coefficients of the selected predictors. The less important a predictor x_i is, the less likely it is to be selected, and then when selection occurs, both of the adjusted estimators $\hat{\beta}_i$ or $\hat{\tilde{\beta}}_i$, conditional on the model, tend to be biased away from zero. Let $E(\hat{\beta}_i | g_r \text{ always})$ denote the expected value of $\hat{\beta}_i$ under the model g_r if model g_r is always fit to the data. Let $E(\hat{\beta}_i | g_r \text{ selected})$ denote the expectation of $\hat{\beta}_i$ when model g_r is selected, as by AIC_c . If $E(\hat{\beta}_i | g_r \text{ always}) > 0$, then we usually find that $E(\hat{\beta}_i | g_r \text{ selected}) > E(\hat{\beta}_i | g_r \text{ always})$. If $E(\hat{\beta}_i | g_r \text{ always}) < 0$, then we find that $E(\hat{\beta}_i | g_r \text{ selected}) < E(\hat{\beta}_i | g_r \text{ always})$. The strength of the bias depends mostly on the importance of the predictor, measured by its overall selection probability (and that probability increases with sample size and goes to 1 as n goes to infinity if $|\beta_i| > 0$).

Consider Table 5.17, which gives Monte Carlo results for a model mimicking the GPA data. Using that information, and excluding the model that computed the percent relative bias of $\hat{\beta}_i$, $PRBias(\hat{\beta}_i)$ relative to model selection. The reference value for computing bias is the value of β_i from the data-generating model, not the parameter value estimated conditionally to model g_r when model g_r is always fit to the data. The frequency of occurrence of the given parameter (i.e., predictor x_i) in the model is denoted by $Pr\{x_i\}$:

i	$Pr\{x_i\}$	$PRBias(\hat{\beta}_i)$
1	0.954	14%
2	0.751	25%
3	0.650	31%
4	0.162	267%

A 31% relative bias in conventional estimators due to model selection is of concern (let alone 267%, but x_4 is not in the K-L best model).

We have looked at this issue for other models and sample sizes, and selection, and it is quite clear that this aspect of model selection bias is, as above, strongly related to the importance of the predictor. Selection bias is less for a predictor always included in the model.

presence of better predictors they correlate with) and that a few of those unimportant predictors will end up in the set that happens, all the model selection biases operate in a direction that think that the selected variables are important (“significance testing terms”). The best way to reduce this risk of misleading is a small list of carefully considered candidate variables. (This is selecting one or more unimportant variables; just have a small set of conceived variables and a small sample size; see, e.g., Freedman et al. 1988, 1990).

Two undesirable, but mutually exclusive, properties of model selection strategies particularly relevant to all-subsets regression for model selection are worth noting here: overfitting the data or overfitting the model. The strategy is to always fit and use the global model, you will overfit the model (i.e., include unnecessary variables). This approach avoids avoid subjectively tailoring the model to the data, but you will overfit the model and inflate standard errors of all the $\hat{\beta}_i$. This loss of precision makes all the estimates are worthless. Thus, usually one is forced to use model selection with multivariable observational data (it is a priori considerations).

If you use a subjective selection procedure of first fitting a model and examining the results (e.g., residual plots, r -squares, leverage, etc.) in search of a better model based on a synthetic criterion of your own choosing, you probably will overfit the data. Thus, you will include in model structure what are really noise of the data, thereby possibly biasing $\hat{\sigma}_{y|x}^2$ quite low and making “noise” as real structure. The resultant model may become a poor representation of the particular data at hand than a valid inference for the population. Samples have their nearly unique peculiarities as well as common features that would show up in all, or most, samples you might get from the population at hand. Inference is about correctly identifying the common features of samples. When you overfit the data, you mistakenly include in the model structure uncommon data features that would not be found in a new sample that might arise.

When you have a large number of models for a much larger number of variables (like $R = 1,024$, $m = 10$, all-subsets) and a large number of models fit and considered for selection, you run a high risk that you will overfit the data. The use of AIC_c reduces this risk (because it looks at model fit penalized by a function of model size, n ; but there is no built-in “penalty” for having a huge number of models). It does not eliminate it for all-subsets selection: Some degree of overfitting remains. For this reason, and the instability of all-subsets selection, it is better to properly evaluate model selection uncertainty under a variety of conditions and use inference rather than just use the selected best model.

is, either the bootstrap failed (the comment data, Section 5.3.1), or it gave some peculiar results that the theoretical approach did not (the comment data, Section 5.3.1). Two small studies that evaluated the use of assessing aspects of model reliability after model selection (the comment data, Section 5.3.1) expressed pessimism that the bootstrap would always be a reliable way to assess model selection uncertainty for a task (Freedman et al. 1988, Dijkstra and Veldkamp 1988). However, we can now offer is that for even moderately high-dimensional data ($m \geq 7$, hence $R > 100$) one should not blithely think that the bootstrap is a reliable way to assess model selection uncertainty for a task. The method needs more study.

In fairness, it can also be said that AIC_c and associated methods here need more evaluation for their performance under a task. However, a more basic issue is whether or not ever to do model selection (especially when the number of predictors is large) when there is a single best model and ignoring all other models. The problem is the instability of what model is selected (cf. Brieman 1996) and the selection biases. Model instability arises when all model selection probabilities (i.e., the π_i) are low. For large R (hence if m is at all large) the expected K-L best model might have selection probability zero, so the set of supposed important regressor variables, as judged by the best model, can vary dramatically over samples (An illustration is given in Chapter 6 for $R = 8,191$ models).

At a fundamental level the question of variable selection is a question of the strength of evidence in the data for the inclusion of a predictor variable. If the problem is thought of this way, model selection as such is an illegitimate discretization of what ought to be estimating continuous parameters (the regression coefficients in a model). The flaw in using model selection is then just like the flaw in hypothesis testing that makes a problem a reject-or-not problem. It ought to be approached as an evaluation of strength of evidence. The use of P -values rather than strict reject-or-not procedures is also a flawed methodology; see, e.g., Harlow et al. 1997, Selwyn et al. 1997. Selwyn et al. strongly recommend against doing all-subsets selection with the goal of identifying a single “best” model: Promoting this practice is a failing of statistical science.

We believe that the only defensible reason for fitting all-subsets models should be to obtain the full set of Akaike weights, and to use them are based on the full set of models as mediated by their Akaike weights (i.e., model averaging). The selected best model is a subset (of R) of the predictor variables. Unless the Akaike weight of the best model is very high (say $w_k \geq 0.9$), we maintain that it is not reasonable to infer that one has found *the* important predictors, and that the other selected are unimportant. As noted above, it is not properly

value is only with respect to prediction, reliable causal inference cannot be made from just the data alone when those data are from an experiment.

We are led to believe that the only legitimate application of model fitting with purely observational data (and then only after careful selection of the number of predictors, as discussed in Section 5.3.) is prediction. In this context we recommend model averaging. Model averaging is made with each model, and the Akaike weights are used to produce a weighted average of these predictions. We do not know of any other averaging, but we have seen it only in the Bayesian literature. In the Bayesian-based model weights (see, e.g., Madigan and Raftery 1995, Hoeting et al. 1999, Raftery et al. 1997, Hoeting et al. 1999).

As noted by Breiman (1996), selection of a best model from a set of models is inherently unstable in its outcome. The solution proposed by Breiman to produce stabilized inferences is a type of model averaging. He suggests to perturb sets of the data (such as bootstrap samples cross-validation) to produce a model in each case, and produce some sort of averaged inference. This, for stabilized inference, is a sort of reverse strategy: Keep the data as is, but find for each fitted model its Akaike weight; then average the predictions as some form of weighted average over all the models.

Interest in regression parameter estimates in conjunction with model averaging all-subsets model fitting will no doubt continue. Perhaps the need for this (we are not convinced). Motivated by this recommendation to use model-averaged predictions, we decided to relate such prediction to parameter estimation. This led to the development of ideas and issues that we will outline here. These are issues that we will discuss in Chapter 6, but still need additional research.

The model-averaged prediction (estimate) of $E(y | \underline{x})$ is

$$\hat{E}(y | \underline{x}) = \sum_{r=1}^R w_r \hat{E}(y | \underline{x}, g_r).$$

We define an indicator function for when a predictor is in a model

$$I_i(g_r) = \begin{cases} 1 & \text{if predictor } x_i \text{ is in model } g_r \\ 0 & \text{otherwise.} \end{cases}$$

For model g_r the value of β_i is denoted here by $\beta_{i,r}$. One model-averaged parameter estimator is

$$\hat{\beta}_i = \frac{\sum_{r=1}^R w_r I_i(g_r) \hat{\beta}_{i,r}}{\sum_{r=1}^R w_r I_i(g_r)} = \frac{\sum_{r=1}^R w_r I_i(g_r) \hat{\beta}_{i,r}}{w_+(i)}$$

$$w_+(i) = \sum_{r=1}^R w_r I_i(g_r).$$

we simply set $\beta_{i,r} = 0$. Thus a new estimator, denoted by

$$\tilde{\beta}_i = w_+(i)\hat{\beta}_i.$$

This is just $\hat{\beta}_i$ shrunk toward zero by the amount $(1 - w_+(i))\hat{\beta}_i$.
found, based on empirical results, that we could also construct an estimator as

$$\tilde{\beta}_i \equiv \hat{\beta}_i - (1 - w_+(i))\hat{\beta}_i = \hat{\beta}_i - \widehat{\text{model selection bias}}$$

that is, (5.8) is our original model-averaged estimator adjusted for model selection bias. Certainly, the term $(1 - w_+(i))\hat{\beta}_i$ is an estimator of model selection bias, but it is a usable estimator.

Then we realized that $\tilde{\beta}_i$ is of fundamental importance. The averaged prediction can be expressed as

$$\hat{E}(y | \underline{x}) = \hat{\beta}_0 + \sum_{i=1}^m w_+(i)\hat{\beta}_i x_i = \hat{\beta}_0 + \sum_{i=1}^m \tilde{\beta}_i x_i$$

If we accept $\hat{\beta}_i$ as the appropriate naive estimate of β_i given the data, then heuristically, the above suggests that prediction should be shrunkage toward zero of each parameter's estimate by a factor of the parameter's unimportance ($= 1 - w_+(i)$). The value of shrinkage is a well-known concept in statistics (see, e.g., Copas 1983, Tibshirani 1996); hence, this line of thought worth pursuing.

Thus we have compelling reasons to want to replace, in the multiple regression, the conditional estimator of (5.7) by the unconditional estimator of (5.8). This would allow us to ignore the issue of what to do when we simply make inferences from the full set of models as regression or prediction. In our limited Monte Carlo evaluation of this estimator, we found that $\tilde{\beta}_i$ is less biased by model selection than is $\hat{\beta}_i$.

An unresolved matter is a simple, yet reliable, estimator of the variance of $\tilde{\beta}_i$ and an associated confidence interval for $\tilde{\beta}_i$. The theoretical sampling variance formula for a model averaged estimator is given here:

$$\text{var}(\tilde{\beta}_i) = \left[\sum_{r=1}^R \pi_r \sqrt{\text{var}(\hat{\beta}_{i,r} | g_r)} + (\beta_{i,r} - \tilde{\beta}_i)^2 \right]$$

However, what we need is an estimator of $\text{var}(\tilde{\beta}_i)$. When θ_r is of interest, say θ , appears in every model (hence, θ_r in model

$$\widehat{\text{var}}(\hat{\theta}) = \left[\sum_{r=1}^R w_r \sqrt{\widehat{\text{var}}(\hat{\theta}_r | g_r)} + (\hat{\theta}_r - \hat{\theta}) \right]$$

This variance estimator makes no allowance for the uncertainty inherent in $\hat{\theta}$ because the weights $w_r (= \hat{\pi}_r)$ are random. Because these weights are positive and sum to 1, a decrease is compensated for by an increase in one or more other weights so that the result is an acceptably stable variance estimator even in the presence of sampling variation of the weights.

Applying the above variance estimator to $\tilde{\beta}_i$ the result, $\widehat{\text{se}}(\tilde{\beta}_i)$, is

$$\widehat{\text{se}}(\tilde{\beta}_i) = \sum_{r=1}^R w_r \sqrt{\widehat{\text{var}}(\hat{\beta}_{i,r} | g_r) + (\hat{\beta}_{i,r} - \tilde{\beta}_i)^2 + |\tilde{\beta}_i|}$$

Now the sampling variation in the weights matters very much for the performance of this variance estimator. Another way to look at this is to note that $\tilde{\beta}_i = w_+(i) \hat{\beta}_i$, and the variance of the conditional distribution can be reliably estimated (because one renormalizes the Akaike weights over a subset of models that contain β), but now we need a formula for $\widehat{\text{var}}(\tilde{\beta}_i) = (w_+(i))^2 \widehat{\text{var}}(\hat{\beta}_i) + (\hat{\beta}_i)^2 \widehat{\text{var}}(w_+(i))$. We do not know $\widehat{\text{var}}(w_+(i))$. In simulations we know that it is not trivial (i.e., cannot be ignored).

Despite seeking one, we do not yet know a reliable and simple estimator of $\widehat{\text{se}}(\tilde{\beta}_i)$. Moreover, if we had one, there is a second issue as to how to construct an interval on β_i under this unconditional model averaging. The interval $\tilde{\beta}_i \pm 2 \widehat{\text{se}}(\tilde{\beta}_i)$ is not justified, in general, because the sampling distribution of $\tilde{\beta}_i$ can be *very* skewed (i.e., quite nonnormal) when $E(w_+(i))$ is small.

While finding a reliable estimator of $\text{var}(\tilde{\beta}_i)$ is worthwhile, it is also clear that we have the issue that the sampling distribution of $\tilde{\beta}_i$ can be quite nonnormal. Hence, the only reliable approach to frequentist inference seems to be the bootstrap. To apply the bootstrap to this problem, one must compute this estimator in each of the B bootstrap samples. One will obtain $\tilde{\beta}_{i,b}^*$, $b = 1, \dots, B$, and then determine the percentile of $\tilde{\beta}_i$ in the bootstrap distribution.

interval, and, if desired, the usual bootstrap-based estimate of the standard error.

We again consider the Monte Carlo evaluation of the variance estimator. Results in Table 5.17 were extended to β_4 and $\tilde{\beta}_i$ to compare the two estimators. We need to be clear on what was done in this new Monte Carlo study. The predictor variable x_i was in the selected AIC_c best model for each of the B bootstrap samples.

these estimates is $E(p_i)$ under model selection, the expected size is $10,000 \times \pi_+(i)$ for $\pi_+(i)$ the probability that x_i is model. In contrast, $\tilde{\beta}_i$ was computed for every Monte Carlo $E(\tilde{\beta}_i)$. The results are below; $\tilde{\beta}_i$ has the better performance

i	$\pi_+(i)$	β_i	$E(\hat{\beta}_i)$	$E(\tilde{\beta}_i)$
1	0.95	0.00201	0.00229	0.0021
2	0.75	0.00125	0.00156	0.0011
3	0.65	0.18945	0.248	0.1577
4	0.16	0.08752	0.321	0.0588

5.4 Estimation of Density from Line Transect Sampling

5.4.1 Density Estimation Background

Animal inventory and monitoring programs often focus on estimating population density (i.e., number per unit area). Buckland et al. (1999) provide the theory and application for field sampling and analysis using line transects. We will illustrate several aspects of this in the face of model selection uncertainty using line transect data by Southwell (1994) on the eastern grey kangaroo (*Macrotis lagotis*) at Wallaby Creek, in New South Wales, Australia. The paper by Southwell (Laake et al. 1994) was written for the analysis of line transect data without model selection, and has an option for bootstrapping the standard errors (Laake et al. 1997 for a similar example). Thus, line transect sampling with DISTANCE will be used to provide some deeper insight into model selection uncertainty and will serve as another comprehensive example.

In line transect sampling, the estimator of density (D) is

$$\hat{D} = \frac{n}{2wL\hat{P}},$$

where n ($= 196$ in this example) is the number of objects detected on 78) transects of total length L ($= 88.8$ km) and width w . The unconditional probability of detection, within the strips of width w , the focus of the estimation of animal density is on the probability of detection, this is defined as

$$P = \frac{\int_0^w g(x)dx}{w},$$

where $g(x)$ is the detection function (i.e., the probability of detecting an animal is at perpendicular distance x from the line).

in detectability among animals being surveyed. The detection probability is estimated from perpendicular distances taken from the transect to the first object detected. Assumptions required in line transect sampling and theory are given in Buckland et al. (2001).

Substituting the expression for P into the estimator of D and the w and $1/w$, gives

$$\hat{D} = \frac{n}{2L \int_0^w \hat{g}(x) dx}.$$

Thus, the essence of data analysis here is to find a good approximation for $g(x)$, the detection function. Buckland et al. (2001) review the general form

$$g(x) = \text{key}(x)[1 + \text{series}(x)].$$

The key function alone may be adequate for modeling $g(x)$ if the sample size is small or the distance data are easily described by a simple model. If one or more adjustment terms must be added to achieve a good fit to the data. For the purposes of this example we chose several models of the above form. Each of these models provides a unique, basis for data analysis in this example.

5.4.2 Line Transect Sampling of Kangaroos at Wallaby Creek

Eastern grey kangaroos often occur in family groups; thus, the animal density is the product of the estimated number of groups per unit area and group size. In this example we will focus only on estimating the density of groups of this species of kangaroo. We define the set of parameters for this example in Table 5.24. The analysis theory of line transect sampling has been the subject of a great deal of work since about 1970. The selection of candidate models is relatively well based in this example. The DISTANCE (Laake et al. 1994) was used to compute the maximum likelihood parameters in the key functions (σ , a , or b , in Table 5.24), and the expansions (the a_j in Table 5.24), and choose, using Akaike's Information Criterion among the four.

5.4.3 Analysis of Wallaby Creek Data

The results of the initial analysis of these data are given in Table 5.25. Model 3 was selected using AIC and provides an estimated density of 3.86 groups/km² (conditional se = 1.00 and conditional cv = 10.12). The other models produce relatively similar point estimates of D for these data. The estimated value was 3.86 (model 3). In this example, the estimated lo

half-normal key function.

Model	$g(x)$	K	Key function
1	$\left\{ \frac{1}{w} \right\} \left\{ 1 + \sum_{j=1}^2 a_j \cos \left(\frac{j\pi x}{w} \right) \right\}$	2	uniform
2	$\left\{ e^{-x^2/(2\sigma^2)} \right\} \left\{ 1 + \sum_{j=2}^3 a_j H_{2j} \left(\frac{x}{\sigma} \right) \right\}$	3	half-normal
3	$\left\{ 1 - e^{-(x/a)^{-b}} \right\} \left\{ 1 + \sum_{j=2}^3 a_j \left(\frac{x}{w} \right)^{2j} \right\}$	4	hazard
4	$\left\{ e^{-x^2/(2\sigma^2)} \right\} \left\{ 1 + \sum_{j=2}^3 a_j \cos \left(\frac{j\pi x}{w} \right) \right\}$	3	half-normal

TABLE 5.25. Summary statistics for the line transect data on eastern grey kangaroos at Wallaby Creek, New South Wales, Australia (from Southwell 1999). The best model is shown in bold; \hat{D}_i is used to clarify that \hat{D} is based on model averaging.

Model	K	$\log(\mathcal{L}(\hat{\theta}))$	AIC	Δ_i	$\exp(-\Delta_i/2)$	w_i
1	2	-1,021.725	2,047.449	0.000	1.000	0.495
2	3	-1,021.546	2,049.092	1.643	0.440	0.220
3	4	-1,021.654	2,051.307	3.858	0.145	0.072
4	3	-1,021.600	2,049.192	1.743	0.418	0.209

number of model parameters, and estimated density are similar. However, the estimated conditional sampling variances are almost 3. In this case, all four models contain the same information (about 99%, P); thus, model averaging (4.1) and (4.9) should be used. In either analysis, estimates of unconditional variances and associated confidence intervals should be used in making inferences about population density.

Using the Akaike weights w_i and the conditional variances $\widehat{\text{var}}(\hat{D}_i|g_i)$ for each model, we computed the model-averaged density $\hat{D} = 10.19$ and an estimate of its unconditional variance $\widehat{\text{var}}(\hat{D}) = 1.51$ (4.9). Hence, the (unconditional) standard deviation is 1.23 and its cv is 12.06%. This unconditional cv is slightly higher than the 10.12% conditional on the AIC-selected model. Inferentially the same here whether based on the model-averaged density estimate from the AIC-selected model but using the unconditional estimate of 1.51. In either case the achieved

5.4.4 Bootstrap Analysis

The most obvious advantages of using Akaike weights as estimated unconditional sampling variances (and $\hat{\bar{D}}$) are simple. However, the bootstrap method also can be used to make inferences; the bootstrap is especially useful in complex situations for analytical variances, even given the models, is lacking. Such analytical theory to compare to the bootstrap results is rare.

We used the program DISTANCE to draw and analyze 10,000 bootstrap samples based on transects as the sampling unit (thus, there were 700 bootstrap samples per model). The bootstrap method computes lower and upper confidence intervals as well as an estimated unconditional sampling variance for each model. We present first the results from 10,000 bootstrap samples; we present first the results from 10,000 bootstrap samples. Then we examine the variability inherent here in a “meta-bootstrap” of 1,000 samples, based on the 10 sets of 1,000 samples each from the 10,000 bootstrap samples, based on the 10 sets of 1,000 samples each from 1,001–2,000, and so forth.

The resultant density estimates, by model, and the model selection relative frequencies are shown in Table 5.26. The mean of the estimated density from the bootstrap samples, 10.39, is quite close to the estimate based on Akaike weights and (4.1) and (4.9) (10.19, $\hat{se} = 1.23$). Based on \hat{D}^* , the bootstrap estimate of the unconditional standard error (and $\hat{\bar{D}}^*$) is 1.48. The model selection relative frequencies from the bootstrap procedure are similar to, but do not exactly match, the Akaike weights (as expected). However, the results are close for the favored model (weight $w_1 = 0.50$ (Table 5.25) and from the bootstrap (Table 5.26).

5.4.5 Confidence Interval on D

There are several options for setting a confidence interval on the estimated density and its estimated unconditional sampling variance. There is the usual procedure that assumes that the sampling distribution of the estimator is approximately normal. Hence, an approximate confidence interval is based on

$$\hat{D} \pm 2 \hat{se}(\hat{D}),$$

where $\hat{se}(\hat{D})$ ($= 1.23$ from Section 5.4.3) is the estimated unconditional standard error. For this example, $9.88 \pm 2 \times 1.23$ (7.42, 12.34).

sampling distribution for an estimator θ in any context where θ is strictly positive, and for fixed sample size the $\text{cv}(\hat{\theta})$ tends to zero as the sample size increases. Then one computes lower and upper bounds (Burnham et al. 1987)

$$D_L = \hat{D}/C \quad \text{and} \quad D_U = \hat{D}C,$$

where

$$C = \exp \left[t_{\alpha/2,df} \sqrt{\log[1 + (\text{cv}(\hat{D}))^2]} \right]$$

The confidence level is $1 - \alpha$; $t_{\alpha/2,df}$ is the upper $1 - \alpha/2$ quantile of the t -distribution on df degrees of freedom. The degrees of freedom are the degrees of freedom of the estimated $\text{var}(\hat{D})$. For an approximate 95% interval, $\alpha = 0.05$, it suffices to use 2 in place of $t_{0.025,df}$.

For this example $\hat{D} = 9.88$ (from the AIC-selected model), and the conditional $\text{cv}(\hat{D}) = 0.124$, and thus $C = 1.28$. Therefore, if we base inference on the AIC-selected model, the approximate 95% confidence interval is 8.01 to 12.65. If we base inference on the model-averaged estimate (which increasingly strikes us as the preferred approach), $\hat{\bar{D}} = 10.19$, again with standard error estimate 1.23, hence the approximate 95% confidence interval 8.01 to 12.96. This method would provide a point estimate of $\hat{\bar{D}}$, hence the corresponding confidence interval is more comparable to the analytical results for model selection than the results based on the selected single best model.

A third option is to use the bootstrap to produce a robust estimate, for example, based on the percentile method (Efron and Tibshirani 1994 and Tu 1995). Here the 10,000 values of \hat{D}_b^* generated in previous

TABLE 5.26. Summary of results from 10,000 bootstrap samples for eastern grey kangaroos at Wallaby Creek, New South Wales, Australia (Burnham et al. 1994): Empirical means of the \hat{D}^* by selected model and overall, standard errors, and selection frequencies.

Model	$\hat{\bar{D}}^*$	Standard error estimate	Selection frequency
1	9.97	1.10	4,529
2	10.63	1.41	2,992
3	10.92	2.34	1,235
4	10.75	1.38	1,240
All	10.39	1.48	10,000

confidence interval endpoints on D : $D_{(250)} \leq D \leq D_{(9750)}$, were $\hat{D}_{(250)}^* = 7.88 \leq D \leq 13.78 = \hat{D}_{(9750)}^*$.

The interval lower bounds from the three methods are upper bounds. Results from the bootstrap in this example show less selection uncertainty than the results based on use of Akaike's (2.9); we rectify this matter in Section 5.4.6 below. In general, a good analytical approach or the bootstrap, achievement of nominal interval coverage is likely if a good model is selected, if model uncertainty has been incorporated into an estimate of the standard error, and if nonnormality has been accounted for.

It can be problematic to identify a correct unit of distance for bootstrap resampling. Aside from this fundamental issue, bootstrapping is conceptually simple and can effectively handle model selection, even if computer software exists or can be written. The program described below bootstrapping in the context of distance sampling. In contrast, bootstrapping the experimental starling data would have been nearly impossible. Specialized software development in this case would be prohibitive; and the computer time required would have been in weeks. In these cases, we recommend use of Akaike's estimate of an unconditional standard error, and then an analytical confidence interval procedure.

5.4.6 Bootstrap Samples: 1,000 Versus 10,000

The $B = 10,000$ bootstrap samples were partitioned, randomly generated, into 10 sets of 1,000 samples per set, and estimated on a per-set basis. The results are given in Table 5.27. Before these results we need to establish our goals for precision of estimation (estimate, actually) of quantities such as

$$\widehat{\text{se}}^*(\hat{D} \mid B) = \sqrt{\frac{\sum (\hat{D}_b^* - \bar{\hat{D}}^*)^2}{B - 1}}.$$

The true bootstrap estimate of the standard error of \hat{D} is actually the limit of $\widehat{\text{se}}^*(\hat{D} \mid B)$ as B goes to infinity. We denote this by $\widehat{\text{se}}(\hat{D})$; however, this bootstrap standard error need not converge to the same number as the analytically computed standard error of \hat{D} (using the same notation). For any value of B we have $\widehat{\text{se}}^*(\hat{D} \mid B) = \widehat{\text{se}}(\hat{D}) + \epsilon$, which goes to 0 quickly as B gets large and $\text{var}(\epsilon) = \phi/B$ (ϕ unknown). The goal in selecting B should be to ensure that $\sqrt{\phi/B}$ is small relative to the value of $\widehat{\text{se}}(\hat{D})$. Our preference is to achieve a bootstrap upper

two significant digits. That is, we target a large enough result for $\widehat{\text{se}}(\hat{D})$ (or whatever is being computed) is near two significant digits over all bootstrap samples of size should be 100, we want to be assured that generally our boot between about 99 and 101. This does not seem like too much for; yet even this precision may require in excess of 10,000 it is rarely achieved with $B = 1,000$.

Now consider the variation exhibited in Table 5.27 in bootstrap π_1 , \bar{D} , unconditional standard error of \hat{D} , percentile confidence points (95%), and the interval width, $\hat{D}_U - \hat{D}_L$. Only \hat{D} meets our precision criterion for $B = 1,000$. However, we do not get \hat{D}^* : We already have \hat{D} , from the best model and \hat{D} from It is the other quantities in Table 5.27 that we use the bootstrap to compute. We find (empirically or theoretically) that $\hat{\pi}_1$ is generally within 0.42 to 0.48; this does not meet our precision criterion, similarly, none of $\widehat{\text{se}}(\hat{D})$, the confidence interval bounds, or $\hat{D}_U - \hat{D}_L$ meet our (modest) precision criterion when $B = 1,000$. Based on the 10 sets of samples in Table 5.27 we estimate that the percent coefficients of variation on the bootstrap estimates are $\text{cv}(\hat{\pi}_1) = 0.005$, $\text{cv}(\widehat{\text{se}}(\hat{D})) = 0.007$, $\text{cv}(\hat{D}_L) = 0.004$ and $\text{cv}(\hat{D}_U - \hat{D}_L) = 0.01$. Thus here $B = 10,000$ is not too far from producing bootstrap-computed quantities reliable to (almost) two significant digits. When using the bootstrap, think in terms of $B = 10,000$.

5.4.7 Bootstrap Versus Akaike Weights: A Lesson

The estimated unconditional standard error of \hat{D} is 1.23 based on the formula and use of the Akaike weights. However, based on the intensive bootstrap method we obtained 1.48 for the estimated unconditional standard error of \hat{D} . The bootstrap method is telling us that there is more uncertainty in our density estimator than our analytical (i.e., Akaike) method accounts for. We perceived a need to resolve this issue. Unfortunately we took the wrong approach: We assumed that the bootstrap results were wrong and tried to find out why. It is not wrong, but we mention it here before giving the correct resolution of this matter.

The correct analytical variance of \hat{D} , given a model, is $\text{var}(\hat{D}) = \text{var}(n/L) + \text{var}(\hat{P})$. The $\text{var}(\hat{P})$ component is conditional on n (196) and kangaroo locations may not be independent within the spatial units used here as the basis for the bootstrapping are the

South Wales, Australia (from Southwell 1997), π_1 is the selection probability g_1 ; standard errors and (percentile) confidence intervals for \hat{D} are unrounded; include model selection uncertainty. Results for “All” are based on

Set	$\hat{\pi}_1$	\hat{D}^*	$\widehat{\text{SE}}(\hat{D})$	95% Conf. Int.	
1	0.478	10.39	1.47	7.77	13.64
2	0.412	10.41	1.52	7.71	13.84
3	0.473	10.43	1.42	7.98	13.56
4	0.418	10.37	1.49	7.84	13.73
5	0.442	10.40	1.49	8.00	13.80
6	0.410	10.37	1.48	8.03	13.82
7	0.461	10.40	1.50	7.73	13.90
8	0.447	10.40	1.53	7.92	13.82
9	0.448	10.39	1.48	7.84	13.75
10	0.540	10.30	1.44	7.92	13.47
All	0.453	10.39	1.48	7.88	13.78

ments (78 of them). The length of these segments varies. In generating a bootstrap sample the value of g_b^* is not constant over bootstrap samples b . Instead, g_b^* varies considerably. \hat{D} varies substantially over bootstrap samples. Might these variations be incorporated into the bootstrap samples result in an inflated estimate? We investigated this issue very intensively for this example. It turned out that the bootstrap estimate of $\text{se}(\hat{D})$ was acceptable here. The theoretically computed unconditional $\text{se}(\hat{D})$ did not account for model selection in \hat{D} (even though it accounts for all model selection uncertainty).

The resolution of the matter also turned out to be simple. We decided to consider the need for a variance inflation factor \hat{c} . The estimate above was based on theoretical formulae under ML estimation. However, this variance is underestimated if important assumptions are violated (assumption of independence of detections within a line segment, there may be spatio-temporal variation in true detection probability with distance x ; there may be errors in recording detection distances, etc. are). All these problems lead to more variance than theoretical estimate. We can adjust the theoretical $\widehat{\text{se}}(\hat{D})$ to allow for these sources of variation (in a way analogous to what the bootstrap does). The simplest adjustment is $\sqrt{\hat{c}} \cdot \widehat{\text{se}}(\hat{D})$ as our theory-based unconditional standard error.

When all models considered are subsets of one global model, the QAIC, and variance inflation, comes from the goodness of fit of the model: $\hat{c} = \chi^2/\text{df}$. However, here we have four models, b

of freedom, and \hat{c} for these four models, as well as the Akaike weights on use of AIC:

Model	χ^2	df	\hat{c}	w_i
1	25.11	17	1.48	0.499
2	23.73	16	1.48	0.220
3	24.66	15	1.64	0.072
4	23.40	16	1.46	0.209

The weighted average of \hat{c} , weighted by w_i , is 1.49. Here \hat{D} is the \hat{c} from the selected model or this weighted average. It may not be here; hopefully, this would be the usual situation in distance sampling we use here $\hat{c} = 1.48$, $\text{df} = 17$ from model g_1 .

We should, however, have been using QAIC rather than AIC. Akaike weights might then change (along then with other things). In Table 5.25 we obtain $-2\log(\mathcal{L})$ for each model and thus $(-2\log(\mathcal{L})/\hat{c}) + 2K$ and the associated weights w_i :

Model	QAIC	Δ_i	w_i
1	1,384.71	0.00	0.511
2	1,386.47	1.76	0.211
3	1,388.61	3.90	0.073
4	1,386.54	1.83	0.205

The differences between the Akaike weights based on AIC and QAIC are here trivial (this is because of a large sample size here). The Akaike weights with each \hat{D} from Table 5.25 gives a model-averaged \hat{D} ; the original result was 10.19. For an unconditional standard error of \hat{D} , QAIC-derived w_i we get 1.23 (the same as with AIC-based weights). We stay with the originally computed $\hat{D} = 10.19$. In this example, the use of using QAIC is to make us realize that we need to use a correction factor with our theoretical standard errors.

The quick way to adjust the theoretical unconditional standard error is to compute $\sqrt{\hat{c}} \cdot \text{se}(\hat{D}) = \sqrt{1.48} \cdot 1.23 = 1.22 \cdot 1.23 = 1.50$. The result for the unconditional standard error of \hat{D} was 1.44. The use of $\sqrt{\hat{c}} \cdot \text{se}(\hat{D})$ is not the correct formula (we have used it for an epistemological value). Rather, one should adjust each \hat{D}_i to be $\hat{c} \cdot \widehat{\text{var}}(\hat{D}_i | g_i)$ and then apply (4.9), which here becomes

$$\widehat{\text{se}}(\hat{D}) = \sqrt{\sum_{i=1}^4 w_i \hat{c} \cdot \widehat{\text{var}}(\hat{D}_i | g_i) + (\hat{D}_i - \hat{D})^2}$$

The two approaches will give almost identical results. The $(\hat{D}_i - \hat{D})^2$ are small relative to $\hat{c} \cdot \widehat{\text{var}}(\hat{D}_i | g_i)$, as they are

by the bootstrap (this may be a coincidence). The bootstrap method to obtain the unconditional standard error estimator will, if done correctly, automatically include the uncertainty in that standard error. Estimation based on the formulae, for models that do not automatically estimate variation, will not automatically include overdispersion variation, what theory assumes. Thus in these cases we must always include an empirical variance inflation factor \hat{c} in our calculations.

5.5 Summary

Model-based data analysis is very important, as illustrated in this book, and as demonstrated by the much improved estimates (precision, less bias) for the chain binomial survival data in 5.2. For example, rather than try to separately estimate survival, one should produce smoothed estimates of these parameters from suitable parametric models. For such observational data, other examples here—GPA example, gas mileage data, etc.—would rarely, if ever, know a priori the single best model to use.

However, in all such cases the investigator can and should select a small set of suitable candidate models for data analysis, creating a meaningful, reduced number of predictor variables. QAIC_c-based model selection can be very effective at producing the models based on Akaike weights. If it makes sense to do so (if the models mean something as alternative scientific explanations), one can use the expected K-L best model to draw conclusions. In mind that the selection of that model as best is itself uncertain, the resulting standard errors of estimated parameters can and should include selection uncertainty.

If the models are only a means to the end of “smoothing” the data, case for prediction, then we recommend computing model parameter (prediction) estimators and their unconditional sampling errors based on the Akaike weights. Monte Carlo methods showed that they worked well for the chain binomial models; unconditional coverage is generally close to the nominal 95%, while conditional coverage conditioned on the selected best model may achieve only about 80%. Monte Carlo studies in this chapter also show that there is selection uncertainty but that the Akaike weights are effective at managing this uncertainty. The sampling distribution of Δ_p was examined for the GPA example; we found that generally, for a small number of candidate models, the best model corresponds to at least the 95th percentile and more often to the 99th percentile.

For reliable results from simulation we recommend at least 1000 Monte Carlo samples at each set of conditions used to generate the data. For the bootstrap also: For the results to be stable to two significant figures, one must often use at least 10,000 bootstrap samples. Too many simulations if these simulation methods do not use enough replications.

We do not recommend the dimension-consistent criteria for model selection in the biological sciences or medicine where there is no set of well thought out candidate models. Such criteria are often over-conservative, information, are based on poor assumptions, and perform poorly when the sample size is quite large. We do not recommend using any of these methods for testing for model selection.

The choice of models to examine is important. The examples have demonstrated that a class of logistic models performs better than a model class that assumed constant survival rate across the range of GPA. The GPA example demonstrates that in variable-selection problems, model considerations can lead to much better models than unthoughtful model selection.

In Section 5.3 we note that model selection bias occurs in the presence of correlated predictors. Regression coefficient estimators $\hat{\beta}_i$ are biased away from zero if a variable x_i is included in the model only when that variable is correlated with a variable x_j (i.e., when $\hat{\beta}_i$ is sufficiently different from 0). The less correlated x_i is with x_j , the more biasing effect model selection has on $\hat{\beta}_i$. Estimation of $\hat{\sigma}_{y|x}^2$ is negatively biased by model selection. The use of other methods, provides some protection against both model selection bias and model selection.

The best way to minimize model selection bias is to restrict the model class to models fit to the data by thoughtful a priori model selection.

Usually, selection of a best model is needed if scientific inference is the goal. However, often it is better to think in terms of model selection using the full set of models, rather than selecting a single model for basing inferences on that single model. This is especially true for variable selection as practiced in regression, because the response variable is highly variable. Model averaging is then particularly useful. The relative importance of a variable as the sum of the Akaike weights of the models in which that variable appears and examination of the regression parameters.

Erroneous results have stemmed from the frequent misapplication of simulation in judging various model selection approaches. In the simulation, the generating model has had a few parameters (very often < 8) and no or few tapering effects, and the objective has been to select the model method most often chooses the generating model. This contradicts the spirit of BIC and is a sterile exercise as regards real data. Hence there is no reason why results would apply to real data.

amples we have examined, such intervals have excellent
There are surely cases where this simple approach does not
we have not found any during our investigations.

6.1 Introduction

Much of this chapter is new material not in the first edition. Material moved from other chapters because we judged it to be more helpful on a first reading of introductory ideas. In either case, the material here did not fit well in a logical, linear progression of ideas about K-L-based model selection and multimodel inference. The firmatory orientation with a relatively small set of models is an ordering to the sections of this chapter; they can be read in any order.

We consider $R < 100$, or perhaps even $R < 200$, as a rule of thumb, because many classical variable selection analyses, or at least those that consider thousands, tens of thousands, or even millions of models ($R \gg n$). We consider the analysis as exploratory rather than confirmatory, since the number of models exceeds the sample size, which is a situation that real thought has been expended on the issue of meaningfulness. The detailed properties of model selection, and subsequent model comparison, when the number of models considered is huge are not well studied, and require a gargantuan amount of computing required. To illustrate this issue, Section 6.2 looks in some detail at a published all-subsets regression example with sample size 252 (n) and 13 predictors. This includes the no-effects null model, hence $R = 8,191$ models.

Another subject in this chapter is an overview of selected model selection approaches, followed by a more detailed contrasting of BIC and AIC. All current model selection criteria fall into two classes, those that

hope to cast some light on this matter. Another issue is random coefficient (effects) models. This is a rapidly developing area of importance and promise. A partially related issue is determining sample size for a data set. The issue of “the” sample size is often confused with the fact that there is not a single sample size for complex data structures. The issues of how to correspond to noninteger effective sample size.

Also delved into here is goodness-of-fit for count data models; essentially, this is about estimating c for overdispersion. The general handling of overdispersion is considered when the overdispersion parameter can be estimated and used. The importance of formulas for Bayesian model averaging, for the importance of a small-sample version of AIC is discussed and compared (as we have defined it). Another subject considered is comparison of models when the assumed probability distributions are different; there is a single “error” distribution (e.g., normal or multinomial) and model structural aspects vary.

6.2 An Example with 13 Predictor Variables and 8,191 Models

6.2.1 *Body Fat Data*

In this example multiple regression is used to predict percent body fat from predictors that are easily measured. The data are from 130 males, ages 21 to 81. A key reference is Johnson (1996), *Journal of the American Statistical Association* (http://www.amstat.org/publications/jse/toc.html), available on the web in conjunction with Johnson (1996). The data were generously supplied by Dr. A. Garth Fisher, Health Research Center, Brigham Young University, Provo, Utah. He gave permission to freely distribute the data and use them for research purposes.” Reference to the data is also made in Penrose et al. (1999), which have also been used in Hoeting et al. (1999), a seminal paper on model averaging.

We take the response variable as $y = 1/D$; D is measured in g/cm³ (observed sample minimum and maximum are 0.9950 and 1.0500). Weight, lower body density means more body fat because fat is less dense than muscle and bone. The reciprocal of body density is related to percent body fat; however, there is no agreement among experts on the parameters of that calibration, which is why we use $1/D$ as our response variable. Measuring body density requires a time-consuming underwater weighing method. For each

Sampling aspects for the study were not stated, we suspect that this sample was obtained and processed in a short time in the 1980s.

We consider aspects of five possible approaches, ordered as desirable, in our opinion:

- 1) fit the full (i.e., global) model only;
- 2) select one model by standard stepwise selection from among all possible simple regression models, then ignore selection;
- 3) select the best model using AIC_c and consider selection;
- 4) do full multimodel inference, such as model-averaged predictions for all 8,191 models,
- 5) first reduce in number and refine the predictors based on domain logic, to a set of meaningful derived variables, then do multimodel inference (fewer models).

We then explore using the nonparametric and parametric approaches to compare aspects of selection under AIC and BIC. Our motivation for this example is to demonstrate how much model selection uncertainty there is when the model set is huge, and how this uncertainty is reduced when a smaller set of models is used.

6.2.2 *The Global Model*

Table 6.1 shows basic results of fitting the global regression model $\beta_0 + \sum \beta_i x_i + \epsilon$, $\epsilon \sim \text{normal}(0, \sigma^2)$. Hoeting et al. (1999) model checking showed this to be an acceptable model fit. There are also that correlations among the predictors are strong, but not entirely positive, and range from -0.245 (age & height) to 0.98 (knee & hips). The design matrix is of full rank.

The absolute value of the usual t-test statistic (Wald test) for the selection coefficient is $1/|cv|$. Hence, in Table 6.1 any parameter with a P value would be considered “significant” at the $P = 0.05$ level. The results unambiguously suggests dropping knee (x_8), chest (x_5), and age (x_1) from the global model. The issue of other predictors that one might want to keep by the strong correlations among the predictors. However, there is a strong suggestion that not all 13 predictors need to be in the model. It is motivated to apply formal model selection.

6.2.3 *Classical Stepwise Selection*

We used SAS PROC REG, at its defaults, for stepwise selection. Only a few steps a model was selected, thus giving, perhaps,

Variable	$\hat{\beta}_i$	$\widehat{\text{se}}(\hat{\beta}_i g)$	$ t $
INTERCEPT	0.873844	0.04594	0
X1 age	0.000109	0.00007	0
X2 weight	-0.000215	0.00013	0
X3 height	-0.000163	0.00037	2
X4 neck	-0.000971	0.00049	0
X5 chest	-0.000106	0.00021	2
X6 abdomen	0.002036	0.00019	0
X7 hips	-0.000432	0.00030	0
X8 thigh	0.000525	0.00030	0
X9 knee	0.000024	0.00051	21
X10 ankle	0.000571	0.00046	0
X11 biceps	0.000492	0.00036	0
X12 forearm	0.000923	0.00041	0
X13 wrist	-0.003649	0.00110	0

TABLE 6.2. Selected predictors, hence selected best model, for $n/K = 0.7/20$. The methods applied to the body fat data.

Selection	Indices of predictor variables selected									
Stepwise	2		4	6						
Forward	1	2		4	6	7	8		10	
Backward			3			7		9	10	
Mallows C_p	1	2		4	6		8			
AIC	1	2		4	6	7	8			
AIC _c		2		4	6					
BIC		2			6					

pression of confidence in the selected model. Common belief is that stepwise rather than forward or backward, selection is the best of the three selection methods. We give all three results in Table 6.2, plus Mallows's C_p , AIC, AIC_c, and BIC.

There is substantial variation in the best model by method. Stepwise, stepwise versus forward and backward methods. Even though AIC and AIC_c differ by four predictors. However, this may be because we know that we should use AIC_c in deference to AIC, since $n/K = 252/13 \ll 40$. This example shows that using AIC makes a difference even with $n = 252$ and global K (number of parameters plus σ^2). Mallows's C_p does not select the same model. Some literature erroneously claims that these are identical models; in general they give similar, but not identical results. As expected, BIC is more conservative than AIC.

with 8,191 models, there will be considerable model selection uncertainty. It is unconscionable not to evaluate this uncertainty and use unconditional inference about importance of predictors. It is justified in saying that the selected predictors are important if they are not important (this is the same false dichotomy as in hypothesis testing). Even if we can agree on a selection rule, model selection uncertainty, and we consider this next.

6.2.4 Model Selection Uncertainty for AIC_c and BIC

We assume that most readers are now aware of the Bayesian criterion (BIC, Schwarz 1978, Hoeting et al. 1999): $BIC = -2 \log \pi(\hat{\theta}|data, g) + \log(n) \cdot K$ (whereas $AIC = -2 \log \mathcal{L}(\hat{\theta}|data, g) + 2 \cdot K$). For model j , $\Delta BIC_j = BIC_j - BIC_{min}$ (BIC_{min} is the minimum over all R models). Also, the same structural formula that gives ΔAIC from ΔAIC , is used with ΔBIC to give the (posterior) probability of g_1, \dots, g_R . More information about BIC appears in Section 6.3.

For this data analysis we quantify model selection uncertainty by ΔAIC_c and ΔBIC , but mostly by Akaike weights and posterior probabilities (for BIC). Table 6.3 shows Akaike weights, w_j , for all 8,191 models ordered by largest to smallest w_j . The best model ($r^2 = 0.733$) has a weight of only 0.010738; the other model weights (the 8,191 weights sum to 1). The next best model is ranked 52nd by AIC_c . The weight of evidence is essentially identical for both the best and 52nd model (the ratio is 1.0025). Table 6.3 shows a few other evidence ratios. The model ranked 52nd by AIC_c was ranked 1st by BIC. Plausibility sets include those ranked first to 176th, or even to 642nd.

Another way to determine a confidence set on models is to use the sum of these *ordered* Akaike weights is some value (this method is not the best one in general, but it is useful). All we care about here is how large such confidence sets are. Do models have some plausibility (as opposed to models with near certainty)? The number of models in such sets under the AIC_c is shown below:

Σw_j	# of models	ER
0.900	649	$w_1/w_{649} = 34$
0.950	876	$w_1/w_{876} = 68$
0.990	1449	$w_1/w_{1449} = 413$
0.999	2266	$w_1/w_{2266} = 3579$

are shown.

Model order j	w_j	Δ_j	
1	0.010738	0.00000	
2	0.010711	0.00496	
3	0.010333	0.07690	
4	0.009752	0.19273	
5	0.009459	0.25371	
6	0.009104	0.33008	ER: $w_1/w_6 = 1$
.	
36	0.003924	2.01	
52	0.003349	2.33	ER: $w_1/w_{52} =$
176	0.001450	4.00	
642	0.000321	7.02	ER: $w_1/w_{642} =$
1103	0.000071	10.00	ER: $w_1/w_{1103} =$
8191	1.3E-72	322.01	

There is no computational impediment here to using all model-averaged results, but if we were to use Occam's and Raftery 1994), we would want $\sum w_j$ at least 0.95, and. Thus, model-averaged inferences require here using only of the possible models. This is not at all like using just the estimated as) best model and erroneously thinking that this basis for reliable inferences in a repeated sampling (or Bayesian) (the Akaike weights approximate repeated sampling-based frequencies).

One type of model-averaged inference is the variable weight $w_+(i)$, which is the sum of the Akaike weights for models in which predictor i occurs. Table 6.4 shows these weights and the variables included in the six top-ranked models (Table 6.3 shows w_j for these top six models).

Results for BIC, analogous to those for AIC_c in Table 6.5. The model ranked 12th by BIC is the model ranked defining operational (frequentist) property of BIC is that to infinity, the posterior probability of a single model goes both the sampling context and model set to be fixed, in this example the BIC best model has associated posterior 0.14; hence there is again substantial model selection uncertainty window for $\sum \Pr(\text{model } j) = 0.999$ includes 1,611 models whether approached in a K-L or Bayesian context there is uncertainty here, and inferences, after selection, should reflect

It is beyond our intended use of this example to use model-averaged prediction under AIC (see Hoeting et al

$w_+(i)$	variable	i	AIC _c top model			
			1	2	3	4
0.495	age	1	0	1	1	1
0.933	weight	2	1	1	1	1
0.314	height	3	0	0	0	0
0.652	neck	4	1	1	1	1
0.283	chest	5	0	0	0	0
1.000	abdomen	6	1	1	1	1
0.445	hips	7	0	1	0	0
0.588	thigh	8	0	1	1	1
0.293	knee	9	0	0	0	0
0.448	ankle	10	0	0	0	0
0.600	biceps	11	1	0	1	0
0.828	forearm	12	1	1	1	1
0.976	wrist	13	1	1	1	1
$K =$			8	10	10	9

TABLE 6.5. Some posterior model probabilities for BIC, based on a most to least probable, and the cumulative probabilities.

Model		
order j	Pr(model j)	Σ Pr
1	0.13930	0.13930
2	0.08980	0.22911
3	0.05681	0.28591
4	0.03829	0.32420
5	0.03488	0.35908
6	0.03118	0.39027
.
12	0.01484	0.50689
158	0.00060	0.90010
292	0.00023	0.95001
757	0.00003	0.99001
1611	0.00000	0.99900
8191	9.4E-69	1.00000

model-averaged prediction with these data). Moreover, we simply accepting these 13 predictors and using either to doing stepwise selection in the first place. Rather, we encourage approach (noted as #5 in Section 6.2.1): First reduce the number of models, based on theory and/or logic to a set of few variables related to y , then do multimodel inference.

tially, no analysis at all was done. Rather, one of us (KPI) matter intermittently over several weeks and decided to ables below. The actual data were not studied before this done; i.e., none of the above model fitting was done until below were decided on. A knowledgeable health-trained even better at generating derived variables and suitable not consider improved model forms beyond linear regres-

Weight and height jointly ought to be very important for but not as separate predictors in linear regression; they together. Allometric relationships are common in biology body might have a nearly constant ratio of some function of In many animals mass tends to be proportional to the cube species that proportionality is often very stable. Hence, $z_1 = \log(\text{weight})/\log(\text{height})$ would be very stable for the body and thus variation in this derived variable would be with variation in body fat (hence, $1/\text{density}$).

Additional considerations led to five more derived variables ply information about different dimensions of the predictor set of six, ordered as considered most to least important, a tionale for each, and the predicted sign (+ or -) of each are given below:

$z_1 = \frac{\log(\text{weight})}{\log(\text{height})}$	based on ideas of allometry (+)
$z_2 = \frac{\text{abdomen}}{\text{chest}}$	beer gut factor (+)
$z_3 = \frac{(\text{knee}*\text{wrist}*\text{ankle})^{\frac{1}{3}}}{\text{height}}$	heavyset or light-boned (-)
$z_4 = \left[\frac{\text{biceps}*\text{thigh}*\text{forearm}}{\text{knee}*\text{wrist}*\text{ankle}} \right]^{\frac{1}{3}}$	fleshiness index (+)
$z_5 = \text{age}$	standardized by mean and stand
$z_6 = \text{age}^2$	based on standardized age (+).

Part of the thinking here is that because the response variable is percent body fat, only ratios of body measurements should be used. z_1 the most important predictor seemed as if it ought to be but only relative to some other body size metric. Based on shapes of men (and experience with his own measurements) z_2 . Next there are issues of genetic variation in being slight z_3 : For a given weight, more bony is less fat), and there is (hence z_4). Finally, it seemed reasonable that age might be given z_1 to z_4 . Percent body fat would tend to increase with twenties, but asymptotically, so not exactly linear on age. (derived) predictors, and therefore 63 possible models (w

(via feedback) from the data analysis whether our reason-
not. This feedback aspect of data analysis is critical in the

Table 6.6 gives some basic results about model selection
 AIC_c , for the fat data with these six predictors. Notewor-
 AIC_c best model ($r^2 = 0.659$) is also the best model und-
and BIC. The number of models in the confidence set wi-
to different values is shown below:

$\sum w_j$	# of models	ER
0.900	3	$w_1/w_3 = 2.8$
0.950	4	$w_1/w_4 = 5.9$
0.990	5	$w_1/w_5 = 63.2$
0.999	7	$w_1/w_7 = 178.3$

Thus, we need to consider only a few models here, not hun-
Table 6.7 shows the top six models and the variable-impor-

TABLE 6.6. Akaike weights w_j and ΔAIC_c values for some of the
most to least supported by the data, based on the z_i predictors.

Model		
order j	w_j	Δ_j
1	0.48867	0.000
2	0.23627	1.453
3	0.17745	2.026
4	0.08316	3.542
5	0.00773	8.294
6	0.00287	10.278
.
63	2.0E-58	264.258

TABLE 6.7. Variable relative-importance weights and predictors
models (1 if included, 0 otherwise) based on the z_i and K for that

z_i	description	$w_+(i)$	AIC_c top models				
			1	2	3	4	5
z1	wt/ht	1.000	1	1	1	1	1
z2	gut	1.000	1	1	1	1	1
z3	bony	0.323	0	1	0	1	0
z4	fleshy	0.986	1	1	1	1	0
z5	age	1.000	1	1	1	1	1
z6	age*age	0.264	0	0	1	1	0
$K =$			6	7	7	8	5

	$\hat{\beta}_i$	$\widehat{\text{se}}(\hat{\beta}_i g)$	$ \text{cv} $
z_1	0.18693	0.03714	0.199
z_2	0.14404	0.01717	0.119
z_3	0.04520	0.05828	1.290
z_4	0.00554	0.00168	0.303
z_5	0.00310	0.00070	0.227
z_6	0.00011	0.00053	4.785

We did not intuit the correct sign for predictor z_3 , but it is not a problem. The other signs we predicted correctly, but we can discount z_6 as it is relatively unimportant here ($w_+(6) = 0.264$). The relative importance of the z_i is 2, 1, 5, 4, 3, 6, with z_3 and z_6 having the highest importance.

Overall we think that the a priori considerations here are more important and influential than just the purely statistical model selection. It is always possible to do exploratory analysis after any model is selected (not vice-versa). In fact, we did some final exploratory analysis of having only one or two simple predictors and a high model complexity. A deal of probing we ended up with a very competitive 1-pronged model on $z_e = \text{abdomen/height}$ ($r^2 = 0.682$), which gave the results in Table 6.2.6.

parameter	estimate	$\widehat{\text{se}}(\hat{\beta} g)$	$ \text{cv} $
β_0	0.8259	0.00528	0.0064
β_e	0.0924	0.00398	0.0432

This type of *post hoc* hypothesis-generating analysis is a common one, and it is reported for what it is: strictly exploratory, hypothesis-generating, not confirmatory.

6.2.6 Bootstrap Evaluation of Model Uncertainty

We now return to the original 13 predictors problem and use this extended example: to illustrate the extent of model uncertainty. There is a huge number of models ($R \gg n$). We explore the bootstrap model selection probabilities (π) for these selection methods. We want to know how well the estimated selection probabilities agree with the true weights (or for BIC, the posterior model probabilities; however, this is not a fair question).

For each of several model selection methods we created bootstrap samples, all of size 252, from the data. In addition to the relative frequencies we want to know how many different models are possible, ever get selected (this number depends weakly on the method). For AIC_c and BIC we also looked at the sampling distribution of ΔBIC .

These 1,200 models were listed in rank order, most to least selected. Table 6.8 shows a few records from this list. Model 0101010000111 denotes the model with predictors 2, 4, 6, 10, 11, 12 selected as best, by method, with the actual data at hand. Table 6.8.

From Table 6.8 we see that even the most commonly selected model (rank 1, $\hat{\pi} = 0.0215$) has a very low selection frequency. The model selected by the stepwise method from the actual data (rank 23) has a selection frequency of 62. This bootstrap assessment corroborates that there is considerable model selection uncertainty. Corresponding bootstrap-based assessments of the BIC are in Tables 6.9 and 6.10. Applying AIC_c selection to 10,000 bootstrap samples, 1,233 distinct models were selected with $\hat{\pi} = 0.0215$ for the actual data. This $\hat{\pi}$ compares well to the Akaike weights. Applying BIC selection to 10,000 bootstrap samples, the most commonly selected model was rank 18 for the BIC best model. From the data the posterior model probability for the BIC best model was 0.09 versus 0.14, the comparison of the posterior probabilities, which are conditional on the data at hand, is comparable to predata random variable frequencies.

For AIC_c model selection the bootstrap assessment of the model selection matches well to the Akaike weights. However, when we look at the bootstrap-based estimate of the sampling distribution of the ΔAIC_c , it is quite different from our earlier assessments of this distribution. The minimum was 9.1; maximum was 49.3. This maximum depends well on the data at hand.

TABLE 6.8. Model ranks and selection frequencies for a few of the models selected by the stepwise method applied to 10,000 bootstrap samples with the actual data; the models selected by different methods with the actual data at hand.

Rank	Model	Frequency	
1	1101011100011	215	AIC
2	0101010000111	184	AIC _c & stepwise
3	0101010001111	141	
4	1101011101011	141	
...	
16	1100010100011	89	
17	0100010000011	87	BIC
18	1101011100111	87	
...	
21	0101011100111	70	
22	1101010100111	63	C _p
23	0100010100011	62	

selected by AIC, AIC_c, and BIC with the actual data are also indicated.

Rank	Model	Frequency	
1	1101011100011	218	AIC
2	1001011100011	151	
3	1011011100011	149	
4	1101011101011	148	
5	1101010100011	145	
6	0101010000111	130	AIC _c
7	1101010101111	125	
8	0101011100011	112	
.	
53	1011110000111	36	
54	0100010000011	34	BIC
55	1001011101011	34	

TABLE 6.10. Model ranks and selection frequencies for a few of the models selected by BIC applied to 10,000 bootstrap samples from the bootstrap distribution. The model ranks and selection frequencies selected by AIC_c and AIC with the actual data are also indicated.

Rank	Model	Frequency	
1	0100010000011	891	BIC
2	0100010000101	689	
3	0100010001101	470	
4	0010010000001	388	
5	0100010000001	359	
.	
23	0100010100001	90	
24	0101010000111	87	AIC _c
25	1000010000001	87	
.	
142	1101011100011	10	AIC

activated looking at the same sampling distribution information. The guidelines for interpreting Δ are the same for K-L cross entropy. The average $\Delta\text{BIC} = 7.3$, maximum was 49.2. Bootstrap-based percentiles:

Percentile	ΔAIC_c	ΔBIC
0.50	8.0	6.2
0.90	17.5	15.3
0.95	20.8	18.6
0.99	27.5	25.3

These sampling results are not consistent with what we would expect if the number of models is small, such as $R < 100$ (and certainly

6.2.7 Monte Carlo Simulations

We assumed that the measurements y and $\underline{x} = (x_1, \dots, x_{13})$ be suitably modeled as multivariate normal with a variance-covariance matrix taken to be the observed variance-covariance matrix (this is the parametric bootstrap). This full variance-covariance matrix is as below; \underline{c} is 13×1 and Σ_x is 13×13 :

$$\Sigma = \begin{bmatrix} \sigma_y^2 & \underline{c}' \\ \underline{c} & \Sigma_x \end{bmatrix}.$$

The global model is now also the generating model under which the set of true regression parameters is given by $\underline{\beta}' = \underline{c}'(\Sigma_x)^{-1}$. For limited purposes, to set $\beta_0 = 0$ and generate the data in which each observation \underline{x} is generated from the marginal $\text{MVN}(\underline{0}, \Sigma_x)$ and $y = E(y|\underline{x}) + \epsilon$, where ϵ is a normal random variable with mean 0 and variance $\sigma_{y|\underline{x}}^2 = \sigma_y^2 - \underline{c}'(\Sigma_x)^{-1}\underline{c}$. More details, and particularly the bootstrap approach are given in Sections 4.3.6, 5.3.2, and 5.3.4.

The approximate theoretical standard error for each $\hat{\beta}_i$ (see Section 5.3.2). Hence, we computed the “effect size” $\lambda_i = \beta_i / \text{se}(\hat{\beta}_i | \text{global } g)$; this is essentially the mean for $\hat{\beta}_i$ from a normal($\lambda_i, 1$) random variable. These λ_i values are

i	λ_i	i	λ_i	i	λ_i
1	1.687	6	11.123	10	1.048
2	-1.727	7	-1.484	11	1.048
3	-0.453	8	1.783	12	2.048
4	-2.048	9	0.048	13	-3.048
5	-0.509				

There clearly are tapering effects, and the only trivial effect is for the circumference. The actual average ordering of predictor importances may not match the ordering by $|\lambda_i|$ because of the nature of the predictors.

We generated 10,000 independent samples of size 252 and used BIC model selection. Our interest is in regard to, first, the frequency of models selected: Do those relative frequencies match those from the (nonparametric) bootstrap? Second, the frequency of based sampling distributions of ΔAIC_c and ΔBIC match those from the bootstrap. The answers are yes; there was no substantial difference between the bootstrap and Monte Carlo approaches. Some summary statistics and frequencies are given in Tables 6.11 (AIC_c) and 6.12 (for

of the body fat data, the models selected by AIC_c and BIC with indicated.

Rank	Model	Frequency	
1	1101010100011	197	
2	1101011100011	177	
3	1101011101011	165	
4	1101010101011	152	
5	0101010000111	149	AIC_c
.	
9	1101010100111	132	
10	1101010101111	119	
.	
21	1100011101011	72	
22	0100010000011	70	BIC
23	0100010000111	70	

TABLE 6.12. Model ranks and selection frequencies for a few of selected by BIC applied to 10,000 Monte Carlo samples that mimic of the body fat data; the models selected by AIC_c and BIC with indicated.

Rank	Model	Frequency	
1	0100010000011	1063	BIC
2	0100010000101	852	
3	0100010000001	371	
4	0100010001011	349	
5	010100000100	331	
.	
9	0101010000010	221	
10	0100010100011	212	
.	
19	0100010001001	112	
20	0101010000111	111	AIC_c
21	1001011100011	109	

The sampling distribution percentiles for the 10,000 from the Monte Carlo samples are below:

Percentile	ΔAIC_c	ΔBIC
0.50	5.6	5.5
0.90	12.8	14.9
0.95	15.3	18.6
0.99	21.2	26.6

are generally a little smaller than those from the bootstrap inference we only have a sample of size 1). However, both results show larger percentiles of Δ_p than what we have seen of models is very much smaller than 8,191. We believe that here generally apply when R is so large.

The bootstrap simulation relative frequencies of model to the Akaike weights. However, we noticed that the same Δ_p was stretched to the right. We wondered whether this reflected the bootstrap in this case. Therefore, we then did the Monte Carlo to verify the bootstrap; both approaches gave about the same had to reconsider the distribution of Δ_p because the guidelines interpreting Δ as regards inferential evidence strength about questioned, at least to the extent those guidelines were based on sampling distribution ideas. Our inferential guidelines are as those for ΔBIC (Raftery 1996a), and therefore the same ΔBIC is here also out of line with those guidelines.

The resolution of this concern is that we need to realize that basis for inference about model selection uncertainty under the theoretic approach is the model likelihood $\mathcal{L}(g_i|\text{data})$ and not it (evidence ratios and Akaike weights). This is analogous to based on posterior model probabilities for BIC in a Bayesian order both of these approaches inference is conditional on the likelihood, rather than being justified by ideas of sampling. we are justified in retaining our guidelines as being useful. be interpreted strictly in a sampling distributional framework.

6.2.8 Summary Messages

The first general point illustrated by this example is that selection uncertainty should be expected when the number is large, such as under many instances of all subsets (i.e., 2^R). This example has a good sample size ($n = 252$) relative to the number of predictor variables (13) for such applications, and still many models. With so many models we find here that the selected model has a very small Akaike weight (0.010738), and is essentially tied with second- and third-place models (Table 6.3). Moreover, a confidence set easily includes over 100 models. Any all-subsets application with R far exceeding n can be expected to have such extreme uncertainty wherein even the best model has a very small weight.

When all the models have very low weights, such as here, inferential credibility for any single model regarding what the best predictor variables. It is foolish to think that the variables in the model are “the” important ones and the excluded variables

of variable importance are needed (Brieman 2001). The variable i , $w_+(i)$ (Table 6.4), provide a model-averaged measure of the importance of each predictor variable. They are relative, not the baseline value that corresponds to no predictive value (not at $w_+(i) = 0$, but at some value > 0 (randomization) to estimate this baseline value, see Section 6.9.8).

Even if prediction is the goal, it is foolish to think that a model has any special credibility when its Akaike weight is high. Rather, model-averaged prediction should be used (this is the standard in the literature, see e.g., Brieman 1996, 2001, regarding “bagged forests”). For a vector of predictors each fitted model yields a prediction \hat{y}_j and the model-averaged prediction is $\hat{\hat{y}} = \sum w_j \hat{y}_j$. For a given predictor i , $w_+(i)$ implies the best measure of the absolute importance for that predictor (see Section 5.3.6).

If there is so little inferential weight for the best model, why has model selection been considered to be so useful? Because the best model (the one that gives good in-sample prediction, relative to what is possible, measured by the coefficient of determination, r^2). However, the same r^2 is achieved by many competitor models. The best model here has $r^2 = 0.733$. However, in the confidence set of 876 models, the best model by the sum of the ordered (large to small) w_j being $= 0.9$. The maximum r^2 are 0.718 and 0.742. Any of these 876 models has essentially the same average in-sample predictability, but each is a different subset of predictor variables. This phenomenon of many models as judged by r^2 values is acute when there are many predictor variables that are strongly intercorrelated, which is the usual case in variable selection.

A second point we make here is that stepwise model selection should not be used. Almost any thoughtful model selection will find a model with an r^2 above 0.7. Even stepwise selection, ad hoc though it is, can lead to a model with decent r^2 , relative to what is possible. The same model as AIC. So why not use stepwise selection? (1) there is no theoretical basis for stepwise selection, as regards any model selection criterion; (2) there is no simple way to compute model (inference) weights; (3) of stepwise selection and, as practiced, no such model inference is provided by stepwise selection (unless one resorts to the use of bootstrapping, which never seems to be done); this is a major failing. (3) stepwise selection gives the subjective appearance of much less model uncertainty than actually exists. Only a small number of all possible models are fit, and not all of them. For stepwise selection lists only a few (perhaps < 10) models that were fit to the data. As a result, the user is misled into thinking that model selection uncertainty exists. (4) as practiced, stepwise

that the selected model was the one and only a priori model. The results of stepwise selection depend nontrivially on the number of models to enter and drop predictors; there is no theory for this. See McQuarrie and Tsai 1998, 427–429 about stepwise selection.

A third point follows from the results of the bootstrap evaluation of model selection uncertainty and the sampling distribution of Δ AIC values. This is one of the few cases where we have results for the number of models, R , over several hundred, as opposed to a several dozen or fewer. The guidelines we have given for large Δ did not hold up with $R = 8,191$; the same guidelines for BIC differences and they also did not hold up here. Models with Δ being big here, hence discounting the model with the biggest Δ , like $\Delta = 20$ is “big.” However, the model selection results from these simulations were very consistent with Akaike weights. The data analysis of all 8,191 models. Thus, using these weights (and the evidence ratios) as the basis for inference about model selection by the simulations.

We are quite convinced now, from all of our research on this matter, that the w_j are valid and useful inferential statistics. Conversely, rigorous inference should not be based on the sampling distribution of the Δ AIC, even though rough guidelines on this matter can be given if R is small. It is a principle that sampling variation across models is not the same as inferential uncertainty (as reflected in the sampling distribution of models, or parameters in models, given the single data set). The types of “variation” are often similar, but when they differ, inference should be based on the likelihood.

A fourth point is the advantage of reducing a priori the number of models to consider, especially by reducing the number of predictors. This can best be done by thoughtful creation of meaningful models, or dropping meaningless predictors (or ones whose usefulness is doubtful). This idea was illustrated in Section 6.2.5. Instead of 8,191 models, we ended up with only 63 models; of these 63 only 4 had substantial evidence. The best model had $w = 0.489$, with $r^2 = 0.682$. The results show greater interpretability and logic underlying these model selection results. Out-of-sample predictions better than those from the bootstrap approach (alas, we have no other data with which to test this).

A fifth point is that after the a priori analyses one is free to use more judgment-based analyses and model selection, as long as one is aware about the inferences one makes: a priori versus range-finding or exploratory to reckless data dredging. For example, it is clear from Table 6.7) that variables z_3 and z_6 are useless predictors (or at least, on the “objective” methodology and our understanding of the data). So one might choose to drop them and just use the AIC_c.

that health professionals can use to quickly predict perceived health status. You can do the sort of uninhibited exploration we did here. In our analyses, of really simple models to end up with a linear model that used the one derived predictor, $z_e = \text{abdomen}/\text{height}$ ($r^2 = 0.6$). Many models were considered in arriving at this one, so it should be tested without testing it by application to a new set of data. But we will consider this model as being the result of reckless data dredging. We restricted ourselves to single predictor models.

6.3 Overview of Model Selection Criteria

There is a variety of model selection methods. However, from the point of view of statistical performance of a method, and intended consequences, there are only two distinct classes of methods: These have been called *efficient* and *consistent*. We will characterize these two classes in the next section, introducing other model selection criteria.

Under the frequentist paradigm for model selection one has three main approaches: (I) optimization of some selection criterion, (II) hypotheses, and (III) ad hoc methods. One has a further classification: (1) criteria based on some form of mean squared error (e.g., Mallows 1973) or mean squared prediction error (e.g., Press 1989), (2) criteria that are estimates of K-L information or distance (e.g., the special cases AIC, AIC_c , and $QAIC_c$), and (3) criteria that are estimators of K , the dimension of the “true model” (e.g., Burnham 1995). We will discuss (2) and (3) in the following material.

6.3.1 Criteria That Are Estimates of K-L Information

AIC , AIC_c , and $QAIC_c$ are estimates of the relative K-L distance between the true model $f(x)$ and the approximating model $g(x)$. These criteria were developed in the context of the concept that truth is very complex and that no “true model” exists (that it was immaterial to the argument). Thus, one could only approximate truth with a model, say $g(x)$. Given a good set of candidate models, one could estimate which approximating model was best (among those considered, given the data and their sample size). Linhart and Gajdosi (1980) speak of “approximating families” of models. Hurvich and Tsai (1988) argue that these criteria select the best finite-dimensional approximation to the true model for large samples when truth is infinite-dimensional. The basic idea seems reasonable in the biological sciences.

When sample sizes are quite large, there are other criteria that offer advantages in model selection and inference (e.g., Akaike 1973). These criteria specifically allow for “misspecification” of the true model.

literature (but see Konishi and Kitagawa, 1996).

Takeuchi (1976) provides a general derivation from AIC. An intermediate result indicated that a selection when the candidate models were not particularly close a He derived TIC (Takeuchi's information criterion) for mo a more general bias adjustment term to allow $-2 \log_e(\mathcal{L})$ an asymptotically unbiased estimate of relative K-L,

$$\text{TIC} = -2 \log(\mathcal{L}) + 2 \cdot \text{tr}(J(\theta)I(\theta)^{-1})$$

The matrices $J(\theta)$ and $I(\theta)$ involve first and second mix of the log-likelihood function, and "tr" denotes the matrix is an approximation to TIC, where $\text{tr}(J(\theta)I(\theta)^{-1}) \approx K$. T excellent when the approximating model is quite "good" a when the approximating model is poor. One might consid and worry less about the adequacy of the models in the se consideration involves two issues that are problematic. Fi worry about the quality of the set of approximating mode this is not something to shortcut. Second, using the expan term in TIC involves estimation of the elements of the ma (details provided in Chapter 7). Shibata (1989) notes tha these two matrices can cause instability of the results of n that the matrices are of dimension $K \times K$). If overdispers data, then the log-likelihood could be divided by an estimat factor, given QTIC. In most practical situations, AIC and approximations to relative K-L information.

Linhart and Zucchini (1986) proposed a further gener (1993) proposed a network information criterion (NIC) training samples in neural network models. Konishi and K gest even more general criteria for model selection and pro into AIC and TIC and their derivation. Shibata (1989) deve criterion, based on the theory of penalized likelihoods. F called RIC for "regularized information criterion." We w methods, since they would take us too far afield from our they do not have the direct link with information theory a relative K-L distance. However, we note that almost no to evaluate the utility of these extensions in applied probl of these criteria must be reserved for problems where the large and good estimates of the elements of the matrices (bias adjustment term are available.

Mallows's C_p (Mallows 1973, 1995) statistic is well selection, but limited to LS regression problems with norm C_p lacks any direct link to K-L information. Atilgan (19 tionship between AIC and Mallows's C_p , shows that un

linear regression assumptions hold, the two criteria seem to be the same model and rank the contending models in the same way. This is not equivalent. We have not found a small-sample version of C_p that be useful when the sample size is small compared to the number of variables (like AIC_c) (see Fujikoshi and Satoh 1997). Rissanen (1994) provide a robust version of C_p (also see Sommer 1994). Of course, adjusted R^2 has been used in classical multiple regression analysis; however, it has very poor performance (see e.g., Akaike 1998).

6.3.2 Criteria That Are Consistent for K

This section deals with a class of criteria used in model selection that are “consistent” or “dimension-consistent” and with how these criteria have been developed, based on the assumptions that an exactly specified model is one of the candidate models being considered, that the true selection goal is to select the *true* model. Implicit is the assumption that the dimension is of fairly low dimension (i.e., $K = 1-5$ or so) and that the true generating (true) model, is fixed as sample size increases. The criteria were derived to provide a consistent estimator of the order or dimension of the “true model,” and the probability of selecting this “true model” increases as sample size increases. Bozdogan (1987) provides a nice review of “dimension-consistent” criteria. The best known of the “dimension-consistent” criteria was derived by Schwarz (1978) in a Bayesian context. BIC for Bayesian information criterion (or occasionally Bayesian information criterion); it is simply

$$\text{BIC} = -2 \log(\mathcal{L}) + K \cdot \log(n).$$

BIC arises from a Bayesian viewpoint with equal prior probabilities for each model and very vague priors on the parameters, given the assumption that the purpose of the BIC-selected model was often simple prediction. Without scientific understanding of the process or system under study, the estimator of relative K-L.

Rissanen (1989) proposed a criterion that he called the minimum description length (MDL), based on coding theory, another branch of information theory (see also Yu 1996, Bryant and Cordero-Braña 2000). While the details of its justification are difficult to follow without a strong background in information theory, his result is equivalent to BIC. Hannan and Quincampoix (1979) criterion (HQ) for model selection whereby the penalty term

$$c \cdot \log(\log(n)),$$

practice. Bozdoğan (1987) proposed a criterion he called “consistent” and F denoting the use of the Fisher informa-

$$\text{CAICF} = -2 \log(\mathcal{L}) + K \{\log(n) + 2\} + \log$$

where $\log |I(\hat{\theta})|$ is the natural logarithm of the determinant of the Fisher information matrix. He has recently advanced a selection criterion based on a notion of complexity (ICOMP, Bozdoğan 1997). CAICF nor ICOMP is invariant to 1-to-1 transformations and this feature would seem to limit their application. AIC and TIC are invariant to 1-to-1 transformations.

We question (deny, actually) the concept of a simple model in biological sciences (see the Preface) and would surely not accept even if a “true model” existed, it might be included in the set of candidate models! If an investigator knew that a true model existed in the set of candidate models, would she not know which model to select? A little philosophical justification for these criteria in the biological and medical sciences, although they have seen frequent application, by a few people seem to be aware of the differences in the basis for these dimension-consistent criteria relative to criteria based on K-L information. The dimension-consistent criteria are different and more objective than those addressed by criteria that are estimated from the data.

People have often (mis) used Monte Carlo methods to evaluate model selection criteria, and this has been the source of confusion in some of the literature (Rosenblum 1994). In Monte Carlo studies, one *knows* the true model and often considers it to be “truth.” The generating model is known and it is included in the set of candidate models. In the analysis of simulated data, attention is (mistakenly) focused on what criterion best selects this true model (e.g., Bozdoğan 1987, Fujikoshi and Sato 1990, Chen 1997). Under this objective, we would suggest the use of dimension-consistent criteria in this artificial situation, especially if the true model was quite low (e.g., $K = 3-5$), or the residual variance was small, or the sample size was quite large. However, this situation is far from that confronted in the analysis of empirical data in the biological sciences. Monte Carlo studies to evaluate model selection criteria in the analysis of real data must employ generating models with small effect sizes and substantial complexity. Such evaluations are not the selection of a best approximating model and ranking of the models. The notion that the true (in this case, the generating) model is known must be discarded.

Research into the dimension-consistent criteria has often been limited to a model with only a few large effects. More realistic models with many small or tapering effects have been avoided. In addition, the basis

addition of new geographic field sites or laboratories, the number of years, and the inclusion of new animals with genetic variation. Thus, as substantial increases in sample size are achieved, the model also increases. The data-generating model does not change as $n \rightarrow \infty$. We have found that the dimension-consistent criteria perform poorly in open population capture–recapture models even when K is small, but the parameters reflect a range of effect sizes (Lebreton 1998).

Notwithstanding our objections above, the sample sizes in the examples above show the benefits of dimension-consistent estimation of model parameters, very, very large by any usual standard. In the examples where the true model has substantial residual variances) we have seen the need for sample sizes in the thousands or much more before the consistent criteria select the “true model” with a high probability. In cases where the true model was very large, say 100,000, one might merely examine the data to decide on the parametrization, with little regard for the prior (given the assumption that the true model is being sought, the number of candidates). It should be emphasized that these dimension-consistent criteria are not linked directly to K-L information and are “informative” only in the weakest sense. Instead, their motivation veered to consistency of the order (K) of the supposed “true model” by employing K-L terms (but see Section 2.12.2).

When sample size is less than very large for realistic situations, the data, these dimension-consistent criteria tend to select under the true model, the attendant large bias, overestimated precision, and associated problems of inference. Umbach and Wilcox (1996:1341) present the results of simulations conducted under the BIC-type assumptions. For sample size 100,000, AIC performed better than BIC in terms of the selection of the “correct” set. The two criteria were tied at sample size 100,000. However, even at that large sample size, BIC selected the “correct” set 79% of the cases; this is still far from selecting the correct model. While these criteria might be useful in some of the problems of engineering, we suspect that they have relatively little utility in the natural and social sciences or medicine. Findley (1985) notes that consistency can be an undesirable property in the context of selecting

6.3.3 *Contrasts*

As Reschenhofer (1996) notes, regarding criteria that are based on K-L information vs. criteria that are dimension consistent, “the criteria employed in the same situations, which is in contrast to the criteria that have been designed to answer different questions” (also see Pö

AIC_c , $QAIC$, $QAIC_c$, and HQ) are reasonable for the analysis. The dimension-consistent criteria (e.g., BIC, MDL, HQ, C) might find use in some physical sciences where a simple true model exists and where sample size is quite large (perhaps thousands or more). Still, we question whether this true model would be among the candidate models. Even in cases where a simple true model is contained in the set of candidates, AIC might frequently have better properties than the dimension-consistent criteria.

Still other, somewhat similar criteria have been derived (Akaike 1994a, b, and Stoica et al. 1986 for recent reviews). A large number of methods have appeared, including the lasso (Tibshirani 1996), the bootstrap (Breiman 1992), the nonnegative garrote (Breiman 1995), the quasi-deviance (Qian et al. 1996), various Bayesian methods (Gelman and Chen 1997) including the use of Gibbs sampling (Gelman 1993). Some of these approaches seem somewhat ad hoc, and are difficult to understand, interpret, or compute. Often the methods are not applicable, for example, several are applicable only to regression-type data. We do not pursue these methods here, since they take us too far from the main theme.

In summary, we recommend the class of information-theoretic criteria as estimates of relative K-L information such as AIC, AIC_c, QAIC, QAIC_c, and TIC for the selection of a parsimonious approximating model for sample sizes that occur in practice. If count data are overdispersed, then QAIC and QAIC_c are useful. If large sample sizes are available, TIC might offer an improvement over AIC or AIC_c. However, further investigations suggest that the simpler criteria perform as well as the more complex ones we examined (Chapter 7).

6.3.4 Consistent Selection in Practice: Quasi-tru

The original motivation for a consistent model selection criterion is the idea that the true (i.e., data-generating) model is in the set of candidates, or is, or may be, nested within some overly general models. As sample size goes to infinity we want to select that true model with probability 1. This idea is also based on the condition that one can increase sample size while keeping the data generating context fixed: New data may enter as sample size n increases. This sort of sampling can be done in Monte Carlo computer simulation, but not in the real world. More formally, the frequentist motivation for BIC is idealized.

Assume that we have a nested sequence of models, g_1 to g_R . Model g_t is neither the first nor last model. The addition of predictors in models g_i , $t < i \leq R$ are actually not needed. The regression based on predictors x_1 to x_R where $y = \beta_0 + \beta_1 x_1 + \dots + \beta_R x_R$ and x_{t+1} to x_R have zero correlation with y and with all

parametrized. Therefore, the unique, lowest dimension model is g_t . This is the model that consistent criteria must select as n gets large. BIC will do this, and the inferred posterior probability BIC will also go to 1 as n gets large (this holds for any criterion but it suffices to restrict ourselves to BIC).

The inferential model “weights” from BIC selection has the same meaning as the Akaike weights, but may be interpreted as probabilities given the data, the model set, and the prior model probabilities. Define BIC differences as $\Delta\text{BIC}_i = \text{BIC}_i - \text{BIC}_{\min}$, where BIC_{\min} is the minimum BIC value over all models and it is the index of the minimum. By context we mean that the index \min may differ for AIC and BIC. We do not complicate the notation to distinguish these two posterior probabilities.

Under BIC the posterior model probabilities are given by

$$\Pr\{g_i\} = \frac{\exp(-\frac{1}{2}\Delta\text{BIC}_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta\text{BIC}_r)}.$$

If there is a true model, g_t , in the set then $\Pr\{g_t\}$ goes to 1 as n goes to infinity; and of course $\Pr\{g_i\}$ goes to 0 for all other models in the set (some model sizes may be required). For model g_t the Kullback-Leibler distance $I(f, g_t) = 0$. However, when model, g_t is nested in any other model g_j , then $I(f, g_t) = 0$. Hence, from the standpoint of its selection, g_t must mean the smallest dimension representation of this true model. For model g_j with smallest K that has $I(f, g) = 0$. For the nested model g_t , $\Pr\{g_t\}$ goes to zero for all $i > t$ as well as for all $i < t$, while $\Pr\{g_t\}$ goes to 1. This is the large-sample behavior of BIC that we are to look for if there is a true model in our set, and it is this pattern of posterior probabilities that is to be taken as evidence for the fitted model, \hat{g}_t , as being true. The fallacy here that we need to make very clear.

Whereas this is the asymptotic behavior of BIC and $\Pr\{g_i\}$, if a model is in the set and is nested in (unnecessarily) more models, then as the sample size is quite large, it is also what will happen if $I(f, g_t) = 0$, $I(f, g_j) = 0$, $I(f, g_i) = 0$, $I(f, g_t) = I(f, g_j)$. However, asymptotically BIC depends on relative distances $I(f, g_t) - I(f, g_j)$, not absolute. On the other hand, in K-L distances are estimable. Thus, even if we have the true model, we cannot infer that model g_{\min} selected by BIC is the true model. $\Pr\{g_{\min}\}$ is 1, or nearly 1. This is a type of nonidentifiability of models rather than parameters in models. Model g_{\min} may be very close to the true model, it may not even have the correct form or the correct predicted values.

The best that BIC can do asymptotically is identify the true model in general versions of the apparent true model in which g_{\min} is the true model.

asymptotically by both BIC and AIC. Thus, the concern is theoretically based on the idea that for a subset of the effects exist: All models in that subset have the identical minimum of the K-L information loss over all models. Then BIC is a criterion that selects from this subset of models the smallest dimension K .

We need a name for this property of a selected model to be “true” (as explained above) even though it is not the true model, it a quasi-true model. In practice, consistent model selection is an inference that a quasi-true model has been found if the $\Pr\{g_{min}\}$ is virtually 1 and that model is nested in more general models in the model set. We do not need this concept of a quasi-true model if a model is selected; we would not be inclined to think that it is true. The theoretic approach also does not need this concept, especially in practice we expect tapering effects in our model set rather than K-L information loss values.

To make these concepts less abstract we give a simple example. Let x_1, \dots, x_6 be independent normal(0, 1) random variables. Given the x_i let $y = 100 + 15x_1 + 10x_2 + 5x_3 + 3z + \epsilon$, where ϵ is a normal(0, 1) random variable. Let the model set consider the following regression models for response variable y :

model	predictors
1	x_1
2	x_1, x_2
3	x_1, x_2, x_3
4	x_1, x_2, x_3, x_4
5	x_1, x_2, x_3, x_4, x_5
6	$x_1, x_2, x_3, x_4, x_5, x_6$

Although y depends on z , y is uncorrelated with x_4, x_5 and x_6 , $\beta_4 = \beta_5 = \beta_6 = 0$. In this model set, model 3 is the true model. The inferential properties of BIC, as regards model 3, are good. Model 3 was the true model.

Another set of six models was also considered: The 6 models with x_4 replaced by z . Denote these as models 1 z to 6 z . Models 1 z to models 3 z . Now the true model (4 z) is in this set. Model 4 z is nested in models 5 z and 6 z . Table 6.13 presents Akaike model probabilities (used) and BIC model probabilities from one random sample of size n for each power of 10, $n = 10$ to 1 million.

In Table 6.13 consider results in the first model set for model 3. BIC at $n = 100,000$ (also perhaps at $n = 10,000$ and on) gives $\Pr(g_3) = 1$. It would be careless, and wrong, to infer from this that model 3 is true. But this pattern of inferential statistics

Model	Inference by sample size									
	10		100		1,000		10,000		100,000	
	set	w	Pr	w	Pr	w	Pr	w	Pr	w
1		0	0	0	0	0	0	0	0	0
2		0	0	0	0	0	0	0	0	0
3		100	49	66	89	27	87	57	99	6
4		0	35	23	10	44	12	26	1	2
5		0	12	8	1	16	1	12	0	9
6		0	4	3	0	12	0	5	0	4
1z		0	0	0	0	0	0	0	0	0
2z		0	0	0	0	0	0	0	0	0
3z		77	1	0	0	0	0	0	0	0
4z		23	51	57	83	55	96	41	97	2
5z		0	23	32	15	22	4	42	3	4
6z		0	25	11	2	23	0	17	0	3

occur if model 3 is true. In applying model selection with $n \geq 100$ but using the second model set the inferences change dramatically. For $n \geq 100$ model 3z (\equiv model 3) is ruled out; it gets an inference of 0. This is a very different inference than for the first model set. The inference is that statistically we can infer only that a best model (by BIC) has been selected, never that it is the true model. Yet the initial motivation for the underlying BIC was that it would select the true model, or at least, given a large enough sample size (it seemed to be implicit in the discussion) of course the true model was in the set—where else would it be? In Section 6.4.2 what criterion BIC is optimizing and how to interpret the \tilde{Pr} and posterior probabilities for BIC should be interpreted.

A few more comments on the example used here. We did 100 Monte Carlo simulations, and results were as shown in Table 6.4.2. Such results are both extensive and mostly irrelevant to the point of this section. We do note below the average value of $\Pr\{g_3 | \text{models } 1z \text{ to } 6z\}$. The number of Monte Carlo simulations is shown:

# M.C. samples	n	models 1 to 6 $\tilde{Pr}\{g_3\}$	models 1z to 6z $\tilde{Pr}\{g_3\}$
10,000	10	0.2664	0.30
1,000	100	0.8067	0.81
1,000	1,000	0.9294	0.92
100	10,000	0.9694	0.96
10	100,000	0.9947	0.98
1	1,000,000	0.9989	0.99

10,000 (or better yet 100,000).

The summary messages of this subsection: **Truth and the true model are not statistically identifiable from data.** BIC selection procedure justifies only an inference that we have the quasi-true model, and this strained concept requires that there be a subset of models with identical K-L distances (an unlikely event, it seems). Conviction is only for very large sample sizes. So we think that in practice BIC cannot be used. Frequentists want it to do in the unrealistic, idealized context where the true model is not consistent (an asymptotic property), which is the justification given for recommending against AIC. The Bayesian perspective is more general. Finally, to argue, after selection, that you have the true model (an oxymoron) you must argue a priori that the true model is in the model set; true models are not statistically identifiable.

6.4 Contrasting AIC and BIC

6.4.1 A Heuristic Derivation of BIC

The derivation of BIC holds both the model set and the true model fixed as sample size goes to infinity. It is a heuristic derivation. If the model set contains the true (generating) model, then BIC selects the true model with probability 1 to that generating model as $n \rightarrow \infty$ (the probability of that model goes to 1), even if the generating model is in some too-general set of models. The literature has not been clear on this. In his derivation Schwarz (1978) interprets the prior probabilities as being the probability that model g_j is the true model (hence the prior probabilities are to be interpreted this same way). However, Neath (1999) make it clear that the derivation of BIC does not require this assumption about the true model being in the set of models. The difference between AIC and BIC is the $\log(n)$ in BIC (and not in AIC) which is needed for idealized asymptotic consistency.

So a question is, why does the $\log(n)$ arise in deriving BIC? To understand the answer to this question, thinking it over in the context of the issue of the role of the “true model.” It was evident that the different interpretations about prior probabilities are irrelevant in deriving BIC. As used, BIC assumes equal prior probability for all models, which is easily adapted to allow any model priors. Because the derivation of BIC does not depend on any aspect of the priors on the models, its derivation and justification has nothing about how we should interpret model prior and posterior probabilities. The Bayesian literature we have seen simply refers to “the prior probability of g_j ,” without clarifying what these probabilities mean. We

The critical quantity to be approximated is the marginal likelihood for the data:

$$\int \left[\prod_{i=1}^n g(x_i | \theta) \right] \pi(\theta) d\theta.$$

(Section 6.4.4 gives basic formulas for the Bayesian approach to model selection using BIC in context.) The parameter θ has dimension K , as does the function $g(\cdot)$. As a function of θ , the product in $g(\cdot)$ under the integral is the likelihood. As we can write it symbolically as

$$\int [\mathcal{L}(\theta | x, g)] \pi(\theta) d\theta,$$

where x represents the data. Under general regularity conditions, as sample size increases the likelihood function “near” the MLE, $\hat{\theta}$, and the probability distribution of $\hat{\theta}$, can be well approximated by

$$\mathcal{L}(\theta | x, g) = \mathcal{L}(\hat{\theta} | x, g) e^{-\frac{1}{2}(\theta - \hat{\theta})' V(\hat{\theta})^{-1}(\theta - \hat{\theta})}.$$

Here, $V(\hat{\theta})$ is the (estimated) $K \times K$ variance–covariance matrix of the MLE. This form of the likelihood is related to the fact that the distribution of the MLE becomes multivariate normal as sample size increases, $\hat{\theta}$ converging to a fixed value θ_0 (see Section 7.1). As regularity conditions there is no requirement that g be the true model. It suffices that $I(\hat{\theta})J(\hat{\theta})^{-1}I(\hat{\theta})$ (same I, J as used in TIC, Section 7.3) is a positive definite model, $I \equiv J$ and $V^{-1} = I$. Nevertheless, for a random sample, $V(\hat{\theta})^{-1} = nV_1(\hat{\theta})^{-1}$ where the matrix $V_1(\cdot)$ is independent of n . $V_1(\hat{\theta})^{-1}$ converges to $V_1(\theta_0)^{-1}$.

Now we consider the needed integral, which is approximated by

$$\mathcal{L}(\hat{\theta} | x, g) \int e^{-\frac{1}{2}(\theta - \hat{\theta})' V(\hat{\theta})^{-1}(\theta - \hat{\theta})} \pi(\theta) d\theta.$$

As n goes to infinity the approximation becomes exact, the likelihood concentrates near $\hat{\theta}$ (which is converging to θ_0) and the prior is negligible (over the space where $\hat{\theta}$ has any substantial probability of being). We can treat $\pi(\theta)$ as a constant. Alternatively, in the spirit of the Laplace approximation we can just directly use the improper prior $d\theta$. The needed integral is then related to the underlying multivariate normal distribution and can be evaluated because we know the needed normalizing constant:

$$\int (2\pi)^{-K/2} \|V(\hat{\theta})^{-1}\|^{1/2} e^{-\frac{1}{2}(\theta - \hat{\theta})' V(\hat{\theta})^{-1}(\theta - \hat{\theta})} d\theta = 1.$$

$$\int \left[\prod_{i=1}^n g(x_i|\theta) \right] \pi(\theta) d\theta \approx \mathcal{L}(\hat{\theta}|x, g) \left[(2\pi)^{K/2} \|n V_1(\hat{\theta})\|^{-1/2} \right]$$

$$= \mathcal{L}(\hat{\theta}|x, g) \left[(2\pi)^{K/2} \|n V_1(\hat{\theta})\|^{-1/2} \right]$$

and by a property of the determinant, $\|n V_1(\hat{\theta})^{-1}\| \equiv n^K \|V_1(\hat{\theta})\|$, we have the approximation

$$\int \left[\prod_{i=1}^n g(x_i|\theta) \right] \pi(\theta) d\theta \approx \mathcal{L}(\hat{\theta}|x, g) \left[(2\pi)^{K/2} n^{-K/2} \|V_1(\hat{\theta})\|^{-1/2} \right]$$

Taking -2 times the log of the right hand side above, we obtain the BIC criterion:

$$-2 \log(\mathcal{L}(\hat{\theta}|x, g)) + K \log(n) - K \log(2\pi) - \log \|V_1(\hat{\theta})\|$$

The literature drops the last two terms of the expression above because, asymptotically, they are dominated by the term $-2 \log(\mathcal{L}(\hat{\theta}|x, g))$ as well as by the log-likelihood term (which is of order n).

We now see that the $\log(n)$ term arises because of the quadratic feature of marginalization over θ (i.e., integrating out θ) in the mathematical requirement in the derivation of BIC that the model set does not need to contain the true model. However, there is nothing in the foundation or derivation of BIC that addresses the model size tradeoff, and hence addresses parsimony as a feature of BIC. This is not a strike against BIC because this tradeoff is a natural part of the basis to know what sort of parsimony the BIC model selection criterion. Simulation studies of this question have been done, but the results comparing AIC and BIC performance depend on the nature of the true model (such as having many tapering effects or not), on whether the model set contains the generating model, on the sample sizes considered, and on the objective: select the true model or select the K-L best approximation. One can simulate situations where either BIC or AIC is preferred. Thus, it is the unknown context and intent (i.e., true model and the intended use that is critical for deciding which method is “correct”).

6.4.2 A K-L-Based Conceptual Comparison of AIC and BIC

The motivation for this section was to clarify what objective K-L discrepancy meant, in the Bayesian-oriented literature about BIC, by the phrase “the model,” or similar such vague phrases. We precede the discussion of objective K-L discrepancy, as opposed to allowing “model selection” to be subjective, hence meaningless (to us). Here are two examples of objective K-L discrepancy.

denotes your probability that the j th model is the most commonly the literature simply refers to the probability, with no clarity about what this means as regards a model being “appropriate” by some unspecified criterion.

From its operating characteristics we know that BIC’s rule selects the true model with probability 1, that is that it asymptotically will select, with probability 1, the true model is in the set. However, such convergence in a sense to a single model does not, and cannot, logically mean that the true model is the best (Section 6.3.4). In fact, as sample size $n \rightarrow \infty$, the model selected is consistent for the quasi-true model in the model set. We denote this quasi-true model below.

For a set of R models the Kullback-Leibler “distance” from the truth is denoted $I(f, g_r)$ (Section 2.1.3). If $g_r \equiv g_r(x|\theta_r)$ denote a parametric family of models with $\theta \in \Theta$, Θ being a parameter space, then g_r is the family member for the unique $\theta_0 \in \Theta$ closest to truth in K-L distance (see Section 7.2). For our purposes we also assume the models are indexed worst (g_1) to best, i.e. $I(f, g_2) \geq \dots \geq I(f, g_R)$. Let Q be the tail-end subset of models by $\{g_r, r \geq t, 1 \leq t \leq R | I(f, g_{t-1}) > I(f, g_t) = \dots = I(f, g_R)\}$ because $t = R$ is allowed, in which case the K-L best model is unique. For the case when Q contains more than one model we assume the models g_t to g_R are ordered such that $K_t < K_{t+1}$ (in principle $K_t = K_{t+1}$ could occur).

The set Q contains models that are all equally good approximations to truth f . However, we can further distinguish them by their space dimension, and we must prefer the smallest one. If $t < R$, and $K_t < K_{t+1}$ holds, then model g_t is the unique best of the R models. As a matter of inference from data, BIC model selection is consistent for this quasi-true model, which is not absolute truth f . (In principle, there might not be a unique quasi-true model.)

Both AIC and BIC model selection actually depend on the difference $I(f, g_i) - I(f, g_j)$, not on absolute K-L values. Only the absolute value is not estimable. For a random sample we can write $I(f, g_i) - I(f, g_j) = I_1(f, g_i) - I_1(f, g_j)$ being for $n = 1$ is a constant as regards sample size. Because the MLE $\hat{\theta}$ is consistent, $I(f, g_j) = n(I_1(f, g_i) - I_1(f, g_j))$. Because the basic convergence properties of AIC and BIC for large n are based on approximations such as

$$\text{AIC}_i - \text{AIC}_j \approx 2n[I_1(f, g_i) - I_1(f, g_j)] + (K_i - K_j)$$

$$\text{BIC}_i - \text{BIC}_j \approx 2n[I_1(f, g_i) - I_1(f, g_j)] + (K_i - K_j)$$

model.

In the case of tapering effects, so $t = R$,

$$2n(I_1(f, g_i) - I_1(f, g_R)) > 0, \quad i < R$$

Hence, as $n \rightarrow \infty$ all these differences diverge to ∞ at a rate proportional to n . Also, the magnitude of these differences dominates the “penalty” term, so that at best only grow at a rate proportional to $\log(n)$. Therefore, the BIC-selected model will converge to model g_R with certainty.

The case of a nontrivial quasi-true model (i.e., $t < R$) corresponds to model g_t nested in models g_i , $i > t$. In this case here. The relevant differences are

$$\text{AIC}_i - \text{AIC}_t \approx 2n[I_1(f, g_i) - I_1(f, g_t)] + (K_i - K_t)2$$

$$\text{AIC}_i - \text{AIC}_t \approx -\chi_i^2 + (K_i - K_t)2,$$

$$\text{BIC}_i - \text{BIC}_t \approx 2n[I_1(f, g_i) - I_1(f, g_t)] + (K_i - K_t)\log(n)$$

$$\text{BIC}_i - \text{BIC}_t \approx -\chi_i^2 + (K_i - K_t)\log(n),$$

Here, χ_i^2 is a central chi-square random variable on $K_i - K_t$. For all $i < t$ the difference $\text{AIC}_i - \text{AIC}_t$ and $\text{BIC}_i - \text{BIC}_t$ diverge to ∞ as $n \rightarrow \infty$, with probability 1, hence model g_t is always selected over g_i for $i < t$. For all $i > t$ the differences $\text{BIC}_i - \text{BIC}_t$ and $\text{AIC}_i - \text{AIC}_t$ diverge to ∞ , with probability 1, as long as $K_i > K_t$ because $\log(n)$ diverges to ∞ . Hence, if there is a nontrivial quasi-true model (the frequentist sampling theory) of the model being the one selected with probability 1 for a big enough sample size.

By contrast, for $i > t$ we only have $E(\text{AIC}_i - \text{AIC}_t) > 0$ but is independent of n . While these expected differences are positive, the actual AIC values are random variables with variance that grows with n . That AIC does not select model g_t with certainty in this hypothesis is clear for $t < R$. However, even if this were the situation, Shibata (1980) shows that there is a sense in which use of AIC leads to consistent estimates and predictions as $n \rightarrow \infty$.

It is clear to us, as argued for above, that the Bayesian “posterior probability” g_i used in conjunction with BIC can, and must, be interpreted as the posterior probability that model g_i is the best model in the set of R models. Mathematically it is

$$\Pr\{g_i | \text{data}\} = \frac{\exp(-\frac{1}{2}\Delta\text{BIC}_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta\text{BIC}_r)}$$

This posterior probability assumes equal prior probabilities for all models and is conditional on that set of R models. Conceivably, this must be interpreted as the probability that model g_i is the best model in the set of R models. This model will generally be unique.

in inference from data is to have credibility.

Finally, given the necessity of this interpretation on the part of model g_i , then logically one must interpret the model g_i the same way, as the prior probability (it can be degree of belief) of the quasi-true model in the set of R models. If one wants to believe that model g_i might be the true model, then one must believe that the true model is already in the model set (see e.g., Wasserman, 2004). Of possible one-to-one transformations of models the focus is on model dimension, K_i , which is unique, hence the dimensionality.

In summary, K-L distance is fundamental to understanding both AIC and BIC model selection. Both selection criteria are used and applied, without assuming the true model is in the model set. The defining characteristic of BIC (i.e., what is it trying to achieve) is asymptotically in relation to the concept of a quasi-true model. BIC seeks to select only a best model at a given sample size; “BIC” is an expected estimated K-L criterion which serves to recognize the trade-off in model selection. For AIC, “best” varies with n . However, i.e., the quasi-true model, does not depend on n . However, BIC-selected model approaches its target “best” model for a given n of the model ordering imposed here by the $I(f, g_r)$. However, AIC and BIC performance depends on the performance criteria (the model or best model), the assumptions they make (usually based on simulation studies) about the underlying K-L values, $I(f, g_r)$, and the sample sizes considered. Failure to properly recognize these factors and issues has led to much confusion in the model selection literature about AIC versus BIC.

6.4.3 Performance Comparison

A generally accepted measure of model and model selection performance is predictive mean square error (MSE). We evaluated the performance of AIC and BIC using Monte Carlo simulated data that mimics the data in Section 6.2. Thus, the simulated variable y has the properties of the fat data example (such as tapering effects) and is generated using 13 predictors (the \underline{x}) as described in Section 6.2.7. The comparison involved generating a random sample of size n from \underline{x} , generating $E(y)$ (\equiv here to $E(y|\underline{x})$) for each given \underline{x} , and then $y = E(y) + \epsilon$ for the first 252 observations. Model selection for each sample size $n = 252$; \underline{x}_{253} , and $E(y_{253})$ were set aside.

Initially, classical model selection was done; hence we selected the best model selected under AIC_c and under BIC for each sample size. Using these selected best models and \underline{x}_{253} , $E(y)$ ($\equiv E(y|\underline{x}_{253})$)

the 10,000 values of $(\hat{E}(y_{AIC_c}) - E(y))^2$ and $(\hat{E}(y_{BIC}) - E(y))^2$ are listed below, respectively.

The model selection procedure with the smallest MSE was chosen. The results of this simulation were a MSE of 5.6849 for AIC_c and 7.6590×10^{-6} for BIC. Thus the ratio of predictive MSEs was estimated as 0.74. Both MSEs have a coefficient of variation of 0.0908. The statistical comparison is based on the paired nature of the comparisons. The mean of $diff = (\hat{E}(y_{BIC}) - E(y))^2 - (\hat{E}(y_{AIC_c}) - E(y))^2$ at the 10,000 values of $diff$ was 1.9741×10^{-6} with a standard error of 0.0908. Hence, a 95% confidence interval on $E(diff)$ of 1.8×10^{-6} was obtained. The results about the sample distribution of the 10,000 values of $diff$ and $(\hat{E}(y_{AIC_c}) - E(y))^2$ are listed below:

performance measure	sample percentiles $\times 10^6$				
	5	10	50	90	95
$(\hat{E}(y_{AIC_c}) - E(y))^2$	0.019	0.089	2.5	14.9	22.7
$(\hat{E}(y_{BIC}) - E(y))^2$	0.029	0.110	3.3	20.5	29.3

This shows that the predictions based on the AIC_c model are stochastically closer to the true $E(y)$ values than are the predictions based on the BIC-selected model (as opposed to the result for MSE which favored the BIC-produced outliers).

In 5.7% of the simulated sample $diff = 0$ occurred; all other cases $diff < 0$. In 5.7% of cases AIC_c and BIC selected the same model. In the other 94.3% of cases, AIC_c selected a different model. Thus, overall by our procedure, AIC_c -model selection performed as well as or better than BIC in 94.3% of samples.

Simple linear regression of $\hat{E}(y)$ on $E(y)$ was also done for each model. The intercept (β_0) and slope (β_1) parameter by each model was estimated. The results are listed below, with standard errors in parentheses.

procedure	$\hat{\beta}_0 \times 10^6$	$\hat{\beta}_1$
AIC_c	-8.83(23.8)	0.9936(0.0016)
BIC	12.80(27.6)	0.9818(0.0019)

These results suggest that prediction based on the AIC_c model is much closer to being unbiased than prediction based on the BIC model. It is interesting that the smaller prediction MSE of the BIC model is at the expense of greater prediction bias.

A second set of 10,000 simulations was done wherein $\hat{E}(y)$ was based on model averaging the $\hat{E}_r(y)$, $r = 1, \dots, R = 8,191$ different models, one for each fitted model. Under such a model-averaged procedure, the model weights vary by selection (i.e., weight generation) and the mean of $\hat{E}_r(y)$. The results of this simulation example of model averaging are listed below:

MSEs have a coefficient of variation of 1.9 %. Other aspects of prediction-based model averaging are about the same as the best-model strategy.

A final noteworthy comparison is that the MSE values are smaller under model averaging, as shown below, for MSE

	model	best	
method	averaged	model	ratio
AIC _c	4.8534	5.6849	0.85
BIC	5.8819	7.6590	0.77

For this simulation scenario (which mimics real data), prediction beats the traditional best model approach, and

Producing the above simulation results for the best-model approach required 9 hours of CPU time on a dedicated 1.9 GHz speed computer with 4GB RAM and 80G hard drive (and about 14 hours clock time). The averaged example took 12 hours of CPU time and 19 hours of wall time on the same computer. Thus, whereas, extensive simulation studies are needed, it will be very computer intensive.

We also examined a few aspects of AIC_c versus BIC model selection in the simulation scenario of the example in Section 3.4. The results are shown on a real capture–recapture experiment, and the simulated data. The reader would benefit from scanning parts of Section 3.4.6, before reading the next two paragraphs. In the simulation, Figure 3.4.6 shows the results of AIC model selection on a set of 14 models. The results on 50,000 Monte Carlo data sets. BIC model selection results are shown in Figure 3.4.7. The key point we make here is that the results are very similar between the models selected by AIC and BIC.

BIC selection in the face of tapering treatment effects is of “only” 2,500 (but still quite less than ∞), and a general model with 34 parameters that was not (quite) in the set of candidate models. BIC selection poorly, as theory would suggest. BIC selection most frequently selected model $g_{2\phi}$ (26.8%), followed by model $g_{1\phi}$ (24.5%), and model g_{3p} (21.5%). Selection frequencies fell rapidly for models $g_{3\phi}$, g_{2p} , g_{1p} , $g_{0\phi}$, 9.8, 3.1, and 1.3, respectively). The BIC selected model was underfit and would have poor confidence interval coverage. If the candidate models were to be expanded to include the general model, then as sample size increased, BIC should select model g_{34} more often. The initial number of nestlings required for BIC to select model g_{34} ($K = 34$) with probability approaching one in the complex example is approximately 108,000 birds (instead of 2,500 in the example).

Clearly, it would often be impossible to find and band 108,000 nestlings on one small island for a particular year; it might be quite u

several years. However, in so doing, other factors become conceptualization of truth must include obvious factors such as in addition to slightly less obvious factors such as technological resighting probabilities and islands with differing vegetation resighting probabilities. The “year” effect is not so much a year, but a host of covariates (most unmeasured) that affect resighting probabilities in complex, nonlinear ways across years. There is individual heterogeneity that is substantial (e.g., with growth rate, dispersal distance). Thus, the concept of truth is very complex. To think that such reality exists as an exact model is useful; to think that such a true model is included in the set of models seems absurd.

The primary foundations of the BIC criteria do not apply to sciences and medicine and the other “noisy” sciences. Real sample size is increased by orders of magnitude in biology. The target “true model” sought by BIC increases in size with sample size. The simple fact is a violation of the assumptions that form the basis of this is allowed under the AIC-type criteria.

Many other published works have compared AIC or AICc to BIC selection. Some, but not all, of these studies note the real-world limitations. For example Hjorth (1994:46) says “. . . the asymptotic theory of measures is theoretically correct, but in the model selection context one should be skeptical against asymptotic results when we are analyzing 50 or even 500 observations. The asymptotics may require ten more data before a reasonable accuracy is achieved, if possible.” McQuarrie and Tsai (1998) report on extensive comparisons of model selection methods. A parsimonious summary of the book is on pages 410–411 of that book.

6.4.4 *Exact Bayesian Model Selection Formulas*

The analytical formulae for Bayesian model selection are presented here to help put BIC in context. Implementation of these formulas requires computer-intensive methods such as Markov chain Monte Carlo.

For structural model g_i the likelihood of θ_i (a $K_i \times 1$ vector) is denoted as $g_i(x|\theta_i) (\equiv \mathcal{L}(\theta|x, g))$ as used in Section 6.4.1). The prior for θ_i is denoted $\pi_i(\theta_i)$. A key quantity needed is the marginal likelihood

$$g_i(x, \pi_i) = \int g_i(x|\theta_i)\pi_i(\theta_i) d\theta_i.$$

This is taken (by us) as the likelihood of model structure g_i given the prior on θ_i . In essence this quantity is what was approximated in Section 6.4.1 to obtain the BIC criterion. The posterior distribution

$$m_i(x, \pi_i) = \frac{g_i(x, \pi_i)}{g_i(x, \pi_i)}.$$

Let the prior probability for model structure g_i be denoted p_i . This prior probability should be interpreted as the prior probability that model g_i is the true model in the set of R models. The posterior probability of model g_i given the data x is the quasi-true model in the set of R models is given by

$$m_i(x, g_i, \pi_i) = \frac{g_i(x, \pi_i)p_i}{\sum_{r=1}^R g_r(x, \pi_r)p_r}.$$

From data analysis alone one cannot conclude that a model g_i is the true generating model even if it gets $m_i = 1$. There is an identification problem here. All we can know from the data is that a model is the true model, i.e., for a large enough sample size we can infer that there is a model in the model set that has a smaller K-L discrepancy, even if the true model is nested in some other models. In order to infer that a model is the true model we need to believe (i.e., *know*) a priori that the true model is in the model set (we just do not know a priori which model it is). This seems reasonable, it is justifiable to believe that there is always a single best model in the model set, if we have a suitable criterion for what is a best model.

The posterior odds ratio is informative:

$$\frac{m_i(x, g_i, \pi_i)}{m_j(x, g_j, \pi_j)} = \left[\frac{g_i(x, \pi_i)}{g_j(x, \pi_j)} \right] \left[\frac{p_i}{p_j} \right].$$

The ratio $g_i(x, \pi_i)/g_j(x, \pi_j)$ is called the Bayes factor; it is the information-theoretic evidence ratio. Often, BIC is motivated as a simple approximation to the Bayes factor, inasmuch as the Bayes factor can be difficult to compute exactly. The prior odds ratio is not informative; the posterior odds ratio solely by the Bayes factor.

6.4.5 Akaike Weights as Bayesian Posterior Model Probabilities

For a large sample size a good approximation to the Bayesian posterior model probability can be based on BIC, provided one is willing to assume that the model probabilities (useful background Sections here are 6.4.1 and 6.4.4). However, the expression BIC can be used more generally to compute model priors. Let p_i be the prior probability placed on model g_i .

¹A note to the reader: At the time we were checking the second edition of this book (late March 2002) we found relationships of such importance that we decided to include them. The results are in this section.

$$\Pr\{g_i|data\} = \frac{\exp(-\frac{1}{2}\Delta\text{BIC}_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta\text{BIC}_r)}$$

To get Akaike weights we use the model prior

$$p_i = B \cdot \exp\left(\frac{1}{2}\Delta\text{BIC}_i\right) \cdot \exp\left(-\frac{1}{2}\Delta\text{AIC}_i\right)$$

B is a normalizing constant (p_i simplifies, as will be shown below)

$$\exp\left(-\frac{1}{2}\Delta\text{BIC}_i\right) \cdot \exp\left(\frac{1}{2}\Delta\text{BIC}_i\right) \cdot \exp\left(-\frac{1}{2}\Delta\text{AIC}_i\right) = \exp\left(-\frac{1}{2}\Delta\text{AIC}_i\right)$$

hence, with this prior probability distribution on models

$$\Pr\{g_i|data\} = \frac{\exp\left(-\frac{1}{2}\Delta\text{BIC}_i\right)p_i}{\sum_{r=1}^R \exp\left(-\frac{1}{2}\Delta\text{BIC}_r\right)p_r} = \frac{\exp\left(-\frac{1}{2}\Delta\text{AIC}_i\right)}{\sum_{r=1}^R \exp\left(-\frac{1}{2}\Delta\text{AIC}_r\right)}$$

which is the Akaike weight for model g_i .

Moreover, this model prior is actually simple and not difficult as shown by simplifying how it is expressed (there is an alternative formula for the normalizing constant):

$$\begin{aligned} p_i &= B \cdot \exp\left(\frac{1}{2}\Delta\text{BIC}_i\right) \cdot \exp\left(-\frac{1}{2}\Delta\text{AIC}_i\right) \\ &= B \cdot \exp\left(\frac{1}{2}[\Delta\text{BIC}_i - \Delta\text{AIC}_i]\right) \\ &= C \cdot \exp\left(\frac{1}{2}[\text{BIC}_i - \text{AIC}_i]\right) \\ &= C \cdot \exp\left(\frac{1}{2}K_i \log(n) - K_i\right) \end{aligned}$$

and

$$C = \frac{1}{\sum_{r=1}^R \exp\left(\frac{1}{2}K_r \log(n) - K_r\right)}.$$

The result easily generalizes to AIC_c and formally to Bayesian model selection. A Bayesian model selection would not handle overdispersion by simply using AIC_c . Formally, the Akaike weights from AIC_c are Bayesian model probabilities for the model prior

$$\begin{aligned} p_i &= C_c \cdot \exp\left[\frac{1}{2}K_i \log(n) - \frac{nK_i}{n - K_i - 1}\right] \\ C_c &= \frac{1}{\sum_{r=1}^R \exp\left[\frac{1}{2}K_r \log(n) - \frac{nK_r}{n - K_r - 1}\right]}. \end{aligned}$$

We will call this the K-L model prior. Because BIC is for large n , the validity of the result (i.e., AIC_c as Bayesian) can be questioned. It does not apply for large n and small or large K_i .

Do not be confused by the $\log(n)$ appearing in both the model prior. BIC actually arises in the context of obtaining an approximation to the Bayes factor, a quantity that is un-

Radner, the $\log(n)$ in BIC arises owing to the re-normalization of the prior. If we assume that the prior is a uniform distribution over the parameter space, then the BIC formula is used with any prior p_i that is not a function of the sample size. If the prior is a function of the sample size, then the posterior will converge (as $n \rightarrow \infty$) to the quasi-true model.

When the K-L model prior is used with BIC so that we get the result, the interpretation of the model probabilities is not the same as for AIC. Whereas the quasi-true model g_i is the “target” of BIC (for a given sample size), the target model of AIC is the model that minimizes the estimated K-L information loss. From Section 7.3 (7.18) that minimizes

$$E_{\underline{y}} \left[I(f, g(\cdot | \hat{\theta}(\underline{y}))) \right] = \text{constant} - E_{\underline{y}} E_{\underline{x}} \left[\log[g(\underline{x} | \hat{\theta}(\underline{y}))] \right]$$

For simplicity we call this target model the K-L best model. The sample size n is fixed, and it is essentially a fitted model as opposed to $g_i(x | \theta_o)$ for BIC, which is the model of smallest dimension. The minimum $I(f, g_i)$. Thus as a Bayesian result we must interpret w_i and a model prior p_i as the probability that model g_i is the true model.

The target K-L best model has a variance-bias trade-off. The model prior is a function of the sample size n and the number of parameters K to be estimated. Still, we did not expect that the model prior in a Bayesian context had to depend on n and K : That prior should be independent of these values (which are known prior to the data). It is easy to numerically explore the K-L model prior so we will do so in the next numerical examples.

There are now two ways to compare AIC and BIC. One is to use the frequentist framework of looking at sampling measures of performance: predictive mean square error and confidence interval coverage. The other is to consider AIC as Bayesian and think about and compare the model prior to the K-L model prior in conjunction with known properties of the models for selection are different, and the interpretation of the model prior is different, for AIC versus BIC. Such comparisons must be based on what we are assuming about the information in the data. The model prior is an estimation and the models as approximations to some common generating distribution (“truth”). It is useful to think in terms of the standard error $|\theta|/\text{se}(\hat{\theta})$. We would assume few or no effects are truly zero. We would assume meaningful, informative data and thoughtfully selected predictors. We assume tapering effects: some may be big (values like 2 or 3) and some are only 2, 1 or 0.5, or less. We assume we can only estimate a few parameters reliably; m might be 20 or as small as 10. These ideas are based on the use of savvy model priors, with properties like the K-L prior. The model prior and K . The K-L prior is a particularly important savvy model prior.

unity. However, the interpretation of what the probability is different for AIC versus BIC. Also, to use the Bayes factor provided by BIC in the context required by AIC implies the prior probability on model g_i be an increasing function of function of K (i.e., a savvy prior). The implicit BIC prior is in the information-theoretic context.

6.5 Goodness-of-Fit and Overdispersion

Overdispersion of count data, relative to a theoretical model, is with to obtain valid inferences. A more sophisticated approach advocated here is to incorporate one or more variance parameters into the parametric model, hence into the likelihood. As noted, this is a desirable approach that can be quite flexible and is not the simple omnibus approach that QAIC and QAIC_c. Lindsey's ideas here would take us too far afield from our objective. Instead, we present a generalization of QAIC to a general overdispersion parameter. This, however, assumes that the data, by some factor, or factors, that allow partitioning the data into subsets of overdispersion, c , can vary by data subset. Before discussing issues, we consider a strategy for obtaining \hat{c} when there is no model.

6.5.1 Overdispersion \hat{c} and Goodness-of-Fit: A Caution

When there is a global model we can usually compute from the data \hat{c} . The logic is that the global model is theoretically the best because all other models are special cases of the global model. If cases cannot, on average, fit the data better than the global model, then the global model is adequate. If the global model is inadequate, then the goodness-of-fit statistic will exceed its degrees of freedom (as shown by a χ^2 test). This will be true even if the global model is structurally adequate. To find \hat{c} meaningfully > 1 we do not know whether this is a problem of overdispersion or inadequate model structure. In the end this distinction is operationally if we cannot generalize the global model. If the global model is not adequate, this must be accounted for, and if it is, then we prefer to say that the lack of fit problem is due to overdispersion. We prefer to say that the lack of fit problem is due to overdispersion then the global model is structurally adequate.

However, the use of \hat{c} (when $c > 1$ is clearly indicated) for inference for count data (such as inflated standard errors) means that the inferences are based on empirical residuals.

tomatically reflect both lack of model structural adequacy and stochastic variation about “truth.” Thus using \hat{c} we are confident in our inferences because importance of a structural data feature (if any) is judged against the totality of unexplained stochastic variation, not lack of fit, just as with models for continuous data.

With no global model, there is the question of how to choose \hat{c} . For example, in the dose-response example of Section 2.1, we have distinct a priori models, all on an equal footing. None is known to be better than all have the same number of parameters. We must compare \hat{c} (and \hat{c}) for each model, because there is no theoretical basis for choosing the best-fitting model. If the best-fitting model has a \hat{c} value use $c = 1$. If even the best-fitting model is a poor fit (i.e., $\hat{c} < 1$) use \hat{c} from that model. The logic here is that overdispersion will show up in each goodness-of-fit test, but so will inadequacy of fit. If at least one model is adequate, then on average, its goodness-of-fit will only overdispersion, and that would, on average, be the best model. If more than one model fits the data (no overdispersion, $\hat{c} = 1$), the chance get $\hat{c} < 1$. Fortunately, this is not a problem, because $c = 1$.

The general strategy to obtain \hat{c} is as follows. Partition the data into s subsets, where each subset (size R_i , $i = 1, \dots, s$) is used to fit a global model. Thus, one gets s ($\leq R$) subglobal models, each of which is a subset of each other. Compute \hat{c}_i for each subglobal model. The data are used in fitting each subglobal model. Even if one of the subglobal models structurally fits the data, if there is problem with the data will stochastically inflate all of the goodness-of-fit tests. Otherwise, if a subglobal model structurally fits the data, its \hat{c} value estimates actual overdispersion. Use the smallest \hat{c} as the estimate: $c = 1$ if for that goodness-of-fit test, $\hat{c} = 1$ (a guideline). Otherwise, use the computed \hat{c} .

We do not see a simple alternative to this strategy, although there are potential biases arising from taking the minimum of a set of \hat{c} values. The smaller is s , the better, as this minimizes potential selection bias. Such bias is not a big concern here, for the reason that the \hat{c} values in each calculation of goodness-of-fit. As a result, the \hat{c} values are pairwise correlated and this reduces selection bias (which is not the case for selection over independent random variables). Of course, if there be small degrees of freedom for a goodness-of-fit statistic, then for other subglobal models have much larger degrees of freedom, then the small degrees of freedom for goodness-of-fit leads to less reliable \hat{c} . However, if the model allows ample, or similar degrees of freedom for \hat{c} and \hat{c} (or 4) this strategy should work. Assessing goodness-of-fit (in particular), requires some judgment. If nothing else, we should use $c = 1$, i.e., when we judge that the model is adequate.

restrictive for many instances of count data. Such data can have overdispersion due to being collected across factors such as gender, age, and so forth. The degree of overdispersion may vary across factors. For example, it is easy to imagine that in a survival study, the dispersion parameter for males versus females that a common \hat{c} should not be appropriate.

Assume that the data are naturally partitioned into V subsets, hence symbolically $data = \cup_{v=1}^V data_v$. Then the log-likelihood is a simple sum of separate log-likelihoods

$$\log \mathcal{L}(\theta|data, g) = \sum_{v=1}^V \log \mathcal{L}(\theta|data_v, g)$$

We let θ represent the full vector of structural parameters, which may be in common over different parts of the likelihood. The log-likelihood for $data_v$ has likelihood component denoted by $\mathcal{L}_v(\theta|data_v, g)$ and the dispersion parameter for $data_v$ be c_v . The appropriate log-likelihood is actually

$$\sum_{v=1}^V \frac{\log \mathcal{L}_v(\theta|data_v, g)}{c_v}.$$

The overdispersion parameters must be estimated under the same approaches apply here as are discussed in Section 2.1. The log-likelihood for each subset of the data. The overall goodness-of-fit is the sum of the separate chi-square statistics on their summed log-likelihoods. However, if even by this overall result we judge the global fit is poor, we should check whether any separate component, $data_v$, clearly shows overdispersion. We can use some $c_v = 1$ if warranted, while other $\hat{c}_v > 1$ are used.

There are some caveats. In theory, the V data subsets should be independent, hence without correlations across subsets. However, one cause of overdispersion is correlation structure in the data, which may exist both within and between subsets. Hence, it would be best to assume that subsets are independent. Still, even if there are weak correlations, data subsets, if the c_v are quite variable, this generalized approach may be preferred to the simple case of only a single \hat{c} . We assume the data partitioning is a priori to data analysis. In a sense, for any model, even though it is based on factors that may affect the data construction.

Determination of the data partition, hence also the variance, should be based on subject-matter knowledge. Strive to have V subsets of data, partition on factors such as sex (i.e., males vs. females) or age (vs. adults), especially if the sexes and ages have different biological processes. Be sure to use the total number of parameters as $K = p + V$, where p is the number of structural parameters.

The QAIC formula is not changed by this data partitioning.

$$\text{QAIC} = \sum_{v=1}^V \frac{-2 \log \mathcal{L}_v(\hat{\theta}|data_v, g)}{\hat{c}_v} + 2$$

However, its small sample version, QAIC_c is problematic involving reciprocals in n_v should be used, but these tend to be a poor partitioning of K over v . We considered the issue of a simpler, but more complex, formula versus the value of just using the simpler formula, hence minimizing the number of formulae one has to remember. Now we suggest using just

$$\text{QAIC}_c = \sum_{v=1}^V \frac{\log \mathcal{L}_v(\hat{\theta}|data_v, g)}{\hat{c}_v} + 2K + \frac{2K}{n - K}$$

When $V = 1$ we can ignore the estimate \hat{c} while finding the MLE, is, we can work directly with $\log \mathcal{L}(\theta|data, g)$. Moreover, we can use the nominal empirical variance-covariance matrix $\hat{\Sigma}$ (as the inverse of the Hessian) directly from $\log \mathcal{L}(\theta|data, g)$, and then the appropriate variance-covariance matrix of $\hat{\theta}$ is taken as $\hat{c}\hat{\Sigma}$. Also, the degrees of freedom for the variance estimates is the df of \hat{c} . However, for $V > 1$ we must use the quasi log-likelihood. To find the MLE, we must

$$\sum_{v=1}^V \frac{\log \mathcal{L}_v(\theta|data_v, g)}{\hat{c}_v}$$

over θ . The likelihood equations to solve are

$$\sum_{v=1}^V \left[\frac{1}{\hat{c}_v} \right] \left[\frac{\partial \log \mathcal{L}_v(\theta|data_v, g)}{\partial \theta_i} \right] = 0, \quad i = 1, \dots, K$$

The elements of the $K \times K$ appropriate Hessian are

$$\sum_{v=1}^V \left[\frac{1}{\hat{c}_v} \right] \left[\frac{\partial^2 \log \mathcal{L}_v(\theta|data_v, g)}{\partial \theta_i \partial \theta_j} \right] = 0, \quad i, j = 1, \dots, K$$

Each \hat{c}_v has associated degrees of freedom $\text{df}_v = n_v - K_v$, the number of structural parameters in the global model for v . In this global model there should be no parameters in common with other v . Hence, another complication is determining the appropriate degrees of freedom for the variance estimates, hence the variance of any component of $\hat{\theta}$. As an exact method, we simply add as $\text{df} = \sum \text{df}_v$. However, to keep it simple we use the summed df for the applicable degrees of freedom.

A single \hat{c} should often suffice, and should not be used in more complicated approaches: Favor parsimony even in model selection.

be ignored (e.g., 1.5 versus 1.7) given that the v_p are on sufficient df_v , the levels of differences to be concerned versus 2, or 1.4 versus 2.8. Bear in mind that exact model overdispersion is not as important as having at least a bare $V = 1$) for overdispersion, as by quasi log-likelihood.

6.5.3 *Model Goodness-of-Fit After Selection*

Often the set of models under consideration contains a most global model), in which case we recommend assessing the global model to the data (preferably before commencing with the global model fits, as by some standard goodness-of-fit selected model will fit the data. We think that this is true selection (but we are not sure). If the global model does one might decide that the lack of fit is not of concern and or QAIC_c. (In fact, if the global model does not fit, but K-L-based model selection, you must use QAIC or QAIC_c.)

There is a philosophy under which one would want to are using BIC, even when the context is such that AIC is clear that BIC selects more parsimonious models than AIC to model fit; if the AIC selected model fits, will the BIC the same data also fit? This is a question we have never seen literature, and we do not know the answer.

The paper by Leroux (1992) motivated our interest in oux (1992) reports the observed versus expected count for automobile accident data ($n = 9,461$). A pure Poisson and two mixture models are fit to the data. BIC selects the two mixture model ($K = 3$), while AIC selects the three-component ($K = 5$). Model selection tends to lead to overly optimistic model fit. Hence, model selection may result in optimistic model fit for the selected model. However, a goodness-of-fit to the global model will not be biased, because no selection occurred. The usual chi-square goodness-of-fit procedure applied selected in Leroux (1992) (this entails some pooling of sparse $\chi^2 = 1.11$ (1 df, $P = 0.2921$) for the AIC-selected model, $P = 0.0031$) for the BIC-selected model.

Whereas the goodness-of-fit $P = 0.0031$ is small, the size here, and perhaps therefore it is acceptable to use for that statistically is not a good fit to the data. We think that acceptable, but it must be argued for on a case-by-case basis. Statisticians have consistently cautioned about drawing inference that does not fit the data. We should not ignore the issue of selection procedures systematically select models that do

data at the nominal α levels used in goodness-of-fit tests, to understand this general issue of the fit, and assessing data after model selection.

6.6 AIC and Random Coefficient Models

6.6.1 Basic Concepts and Marginal Likelihood

Parameters are sometimes considered as “random effects” or as random coefficients; see Longford 1993). In the simple case, all parameters are of the same type (e.g., survival rates). We consider the K elements of $\underline{\theta}$ as random variables. Thus, we conceptualize $\theta_1, \dots, \theta_K$ as independent random variables with variance σ^2 . Now the inference problem could be entirely in terms of fixed population-level parameters μ and σ^2 . However, the direct inference interest is for $\underline{\theta}$ as if the elements of $\underline{\theta}$ were the direct inference interest. In using the likelihood $\mathcal{L}(\underline{\theta})$ we are assuming how $\theta_1, \dots, \theta_K$ may have arisen from some process or some population. The likelihood $\mathcal{L}(\underline{\theta})$ is appropriate for when $\underline{\theta}$ are “fixed effects.” This is a valid approach if we interpret the θ_K as deterministic. It is then possible to fit this global model with likelihood methods and also fit simpler models based on deterministic $\underline{\theta}$, such as $\theta_i \equiv \mu$, where the likelihood is $\mathcal{L}(\mu)$.

However, we may also want to consider an intermediate model where only the two parameters μ and σ^2 , where we regard θ as random with mean μ and variance σ . Thus while we directly have $\mathcal{L}(\mu)$ (hence models g_K and g_1), we also want the likelihood for the two-parameter model g_2 . The parameter σ^2 in model g_2 allows for the possible stochastic nature of the θ_i . Model g_K allows arbitrary θ_i , but this freedom costs us $K - 2$ extra parameters compared to g_2 . The unexplained variation in these K parameters is substantial with them being considered as exchangeable random variables. We select model g_2 rather than model g_K . Model g_2 with only μ and σ^2 , parsimoniously allows for variation in the θ_i (since σ^2 does not do parsimoniously and model g_1 does not do at all).

Whereas the likelihood $\mathcal{L}(\mu)$ is a special case of either g_K or g_1 , the conceptually intermediate model g_2 is not mathematically a model between models g_K and g_1 in the simple sense of a minimistically constrained version of the global model g_K . Model g_2 cannot be fit by standard likelihood methods based only on the likelihood $\mathcal{L}(\underline{\theta})$ and deterministic constraints. These random-effects models are different from other models

Our focus here is on how we can compute a valid MLE for models, such as this case of having only the parameters μ . Our only obvious starting point is model g_K and its likelihood approach to stochastic parameters is to postulate a distribution, hence a model, for the random variable θ , say $h(\theta)$, the needed (proper) likelihood $\mathcal{L}(\mu, \sigma^2)$ based on the ma

$$g(\underline{x} | \mu, \sigma^2) = \int g(\underline{x} | \underline{\theta}) h(\underline{\theta} | \mu, \sigma^2) d\underline{\theta}$$

thus $\underline{\theta}$ has been integrated out. Considering $g(\underline{x} | \mu, \sigma^2)$ parameters given the data, we have $\mathcal{L}(\mu, \sigma^2) = g(\underline{x} | \mu, \sigma^2)$. Model g_2 is computed based on $g(\underline{x} | \mu, \sigma^2)$, which we call the integral in (6.1).

A more informative way to think about random-coefficient models is sometimes a few parameters, defined by deterministic coefficients, explain all the variation in the much larger set $\theta_1, \dots, \theta_K$, (like 2 or 3). We might have $K = 10, 15$, or 20 (or more). If all θ_i are of the same type of parameter (perhaps for 20 years or sites), then θ_i is some “explainable” (i.e., consistent, simple, and understandable) low-level pattern, such as a linear trend, to the variation in θ_i . However, to be consistent with our philosophy of models, the actual values of the 20 parameters will not perfectly fit a linear trend (a two-parameter linear trend). There will be unignorable residual variation in the parameters; the modeling issue is to explain this residual variation we can detect with the data. If the θ_i behave like *iid* random variables, then random-coefficient models are effective tools for data analysis when there are large numbers of types of parameters. The explainable variation is fit by simple structural models, and the unexplainable (not smooth) variation is swept into σ^2 .

Conceptually, we still think of a parsimonious structural model on $\underline{\theta}$, but one allowing homoskedastic residuals. For example, one reason to try the constrained model structure $\underline{\theta} = X\underline{\psi}$ (where X is the matrix of regressors in matrix X), but we think that this model would be too restrictive even if we could apply it to those exact θ_i . If the model is too restrictive (relative to the size of components of $\underline{\psi}$), we could safely ignore the interpretation of $\underline{\theta}$ and define any new structural model based on $\underline{\psi}$ and directly get the likelihood for the new parameters $\underline{\psi}$.

$$\mathcal{L}(\underline{\psi}) = \mathcal{L}(\underline{\theta} | \underline{\theta} = X\underline{\psi});$$

$\underline{\psi}$ would have only, say, 1 to 4 components.

It is often no more reasonable to assume that $\underline{\theta} = X\underline{\psi}$ is a reasonable servable $\underline{\theta}$ than it is to assume that $Y = \underline{Z}'\underline{\theta}$ is exact for the

and model $h(\underline{\theta} | \underline{\psi})$ imposed on $\underline{\theta}$ in terms of fixed parameters. We conceptualize $\underline{\theta}$ as a normal random variable with $\underline{\theta} = X\underline{\delta}$ and the variance–covariance matrix of $\underline{\delta}$ as $\sigma^2 I$. This serves to define $h(\underline{\theta} | \underline{\psi}, \sigma^2)$ and we compute the actual parsimonious reduced model for the

$$g(\underline{x} | \underline{\psi}, \sigma^2) = \int g(\underline{x} | \underline{\theta})h(\underline{\theta} | \underline{\psi}, \sigma^2)d\underline{\theta}$$

It is more likely that we will be faced with mixed models where some random parameters in one or more of our models are shared. To be explicit we extend our notation by partitioning the generic $\underline{\theta}$ into two parts: $\underline{\theta} = (\underline{\alpha}', \underline{\beta}')'$ with $\underline{\alpha}$ fixed and $\underline{\beta} = (\beta_1, \dots, \beta_K)'$. We of course may choose to consider $\underline{\beta}$ as fixed for some models. For others we will have a model imposed on $\underline{\beta}$ in terms of a distribution $h(\underline{\beta} | \underline{\psi}, \sigma^2)$ variable, hence $h(\underline{\beta} | \underline{\psi}, \sigma^2)$. The needed marginal distribution

$$g(\underline{x} | \underline{\alpha}, \underline{\psi}, \sigma^2) = \int g(\underline{x} | \underline{\alpha}, \underline{\beta})h(\underline{\beta} | \underline{\psi}, \sigma^2)d\underline{\beta}$$

Further generalizations are possible, but for our purposes we consider the simple case, hence (6.1) and issues of likelihoods for fixed-effects models,

$$\mathcal{L}(\underline{\psi}) = \mathcal{L}(\underline{\theta} | \underline{\theta} = X\underline{\psi}),$$

versus random-effects reduced models, hence

$$\mathcal{L}(\underline{\psi}, \sigma^2) = g(\underline{x} | \underline{\psi}, \sigma^2),$$

where computing (6.3) requires the multidimensional integral in (6.2) by (6.1). The MLE $\hat{\underline{\psi}}$ from these two models (6.2 vs. 6.3) will be the same value even if $\hat{\sigma}^2 > 0$ occurs ($\hat{\sigma}^2$ might be zero). The likelihood functions will differ. Therefore, AIC values for the two models will be different, and it is important to consider the random-effects model as well as its more restrictive fixed-effects version (wherein $\underline{\beta}$ is fixed) in sets of models fit to data. Also, the variation represented by σ^2 is of interest in its own right. This is the case in Section 3.5.6 where the coefficient model is indirectly fit to real sage grouse survival data. The process variation, σ^2 , from a set of annual survival rates is of interest.

One classical basis for the concept and use of random effects is in experiments where levels of some factor are selected randomly from a population, for example animals from a herd, farms from a region, cultivars from a population of cultivars (see, e.g., Cox and Wermuth, 1994). Then a well-defined inference to a real population. However, in many models need not have this “random selection” feature, especially when K is involved. The set of K annual survival probabilities, S_1, \dots, S_K

level smoothing model imposed on the parameters, such as or even just $S_i = \alpha + \epsilon_i$. If the imposed smoothing (or explainable variation in the survival parameters over time) residuals $\epsilon_i, \dots, \epsilon_K$ to have the properties of exchangeability. Hence, we can treat them as independent, identically distributed and $\text{var}(\epsilon) = \sigma^2$. It is these residuals, relative to some level that are the random effects. Therefore, there is no concern that the years (hence the S_i) be in any way selected at random from a population of years.

In collapsing the problem to $g(\underline{x} | \underline{\psi}, \sigma^2)$ we are restricting $\underline{\psi}$ and σ^2 , hence ignoring the original individual $\theta_1, \dots, \theta_K$, “have our cake and eat it too,” that is, get estimators of $\underline{\theta}$ and σ^2 . This can be accomplished using shrinkage estimators. Shrinkage estimators $\tilde{\theta}$ arise in both Bayesian and frequentist theories (see, e.g., Morris 1975, Morris 1983, Longford 1993, Casella 1995, and Carlin and Louis 1996). Shrinkage estimators, $\tilde{\theta}$, for random-coefficient models such as $E(\delta) = \underline{0}$, $E(\delta\delta') = \sigma^2 I$ can be based on the MLE of $\hat{\theta}$ in such a way that the residuals from direct simple linear regression reproduce the estimate of σ^2 computed in obtaining $\tilde{\theta}$. This is as saying that we may be able to find a suitable proxy for the likelihood of the fitted model in (6.3), $\mathcal{L}(\hat{\underline{\psi}}, \hat{\sigma}^2)$, by use of $\tilde{\theta}$ evaluated at such a shrinkage estimator, hence via $\mathcal{L}(\tilde{\theta})$.

It would be a considerable advantage if such random effects could be fit without ever computing the integral in (6.1). This is a practical approach to getting a nearly correct AIC value for $g(\underline{x} | \underline{\psi}, \sigma^2)$ yet based on $\log(\mathcal{L}(\tilde{\theta}))$. As it is, the shrinkage estimator is a natural way to fit what amounts to model $g(\underline{x} | \underline{\psi}, \sigma^2)$, thus avoiding $\underline{\theta}$ subject to the stochastic “constraint” inherent in the model, without making distributional assumptions.

If we do have to compute the integrals as in (6.1), it is possible using the recent developments from Bayesian methods; see, e.g., Smith (1990), Zeger and Karim (1991), and Carlin and Louis (1996).

6.6.2 A Shrinkage Approach to AIC and Random Effects

Making inferences about all the random and fixed effects in a random effects model can be accomplished by Bayesian methods (Morris 1995), empirical Bayes methods (Carlin and Louis 1996), or shrinkage methods (Tibshirani 1996, Royle and Link 2002). However, has formal model selection been considered in conjunction

section is to give an overview of one way AIC can be extended to mixed models (Burnham in review; Burnham and White 2002). The extension to simple random effects based only on the likelihood $\mathcal{L}(\underline{\alpha}, \underline{\theta})$ is not straightforward because we need to consider the θ_i as random effects. In this likelihood they are technically to be considered as fixed effects. It is not proposed that this is the ultimate methodology we should use for model selection with random effects. But it is a practical one factor is a random effect and sample size is not small.

Fixed effects inference based on the likelihood $\mathcal{L}(\underline{\alpha}, \underline{\theta})$ focuses on the MLE $\hat{\underline{\theta}}$, which is considered conditional on $\underline{\theta}$, a r -dimensional vector. This MLE has conditional sampling variance-covariance matrix W (it may depend on $\underline{\theta}$) which applies to $\epsilon = \hat{\underline{\theta}} - \underline{\theta}$, and hence $\mathcal{L}(\underline{\alpha}, \underline{\theta})$, with a random-effects model with $\text{E}(\underline{\delta}) = \underline{0}$, $\text{VC}(\underline{\delta}) = \sigma^2 I$; $\underline{\beta}$ is an r -dimensional vector (r is the number of random effects) that an unconditional structural model applicable to the (other) fixed effects. The MLE is

$$\hat{\underline{\theta}} = X\underline{\beta} + \underline{\delta} + \underline{\epsilon}, \quad \text{VC}(\underline{\delta} + \underline{\epsilon}) = D = \sigma^2 I$$

VC means variance-covariance matrix. Parameters $\underline{\beta}$, as well as the random effects. In practice, $\hat{E}_{\theta}(W) = \hat{W}$ from standard likelihood theory.

From generalized least square theory, for σ^2 known, the best linear unbiased estimator of $\underline{\beta}$ is

$$\hat{\underline{\beta}} = (X'D^{-1}X)^{-1}X'D^{-1}\hat{\underline{\theta}}.$$

Assuming normality of $\hat{\underline{\theta}}$ (approximate normality suffices), the conditional sum of squares $(\hat{\underline{\theta}} - X\hat{\underline{\beta}})'D^{-1}(\hat{\underline{\theta}} - X\hat{\underline{\beta}})$ has a central chi-squared distribution on $k - r$ degrees of freedom. Therefore, a method of moments estimator of σ^2 is obtained by solving the equation

$$k - r = (\hat{\underline{\theta}} - X\hat{\underline{\beta}})'D^{-1}(\hat{\underline{\theta}} - X\hat{\underline{\beta}}).$$

Under random effects for inference about $\underline{\theta}$ we use shrinkage estimators, not the MLE (see e.g., Burnham in review). Shrinkage is a form of smoothing. Computing the shrinkage estimator requires the

$$H = \sigma D^{-1/2} = \sigma(\sigma^2 I + \hat{E}_{\theta}(W))^{-1/2} = \left(I + \frac{1}{\sigma^2} \right)^{-1/2}$$

evaluated at $\hat{\sigma}$. Then $\tilde{\underline{\theta}} = H(\hat{\underline{\theta}} - X\hat{\underline{\beta}}) + X\hat{\underline{\beta}}$. Move into the form $H + (I - H)AD^{-1}$, where $A = X(X'D^{-1}X)^{-1}X'$. The matrix G such that $\tilde{\underline{\theta}} = G\hat{\underline{\theta}}$. This shrinkage estimator is the sum of squares of the shrunk residuals (i.e., $\tilde{\underline{\theta}} - X\hat{\underline{\beta}}$), divided by

model from $\underline{\theta}$ using only the fixed effects model. Let ℓ be the dimension of $\underline{\alpha}$. Then as a fixed effects model $\text{AIC} = -2 \log \mathcal{L}(\hat{\underline{\alpha}}, \hat{\underline{\theta}}) + 2K$. The random effects log-likelihood

$$\log \mathcal{L}(\tilde{\underline{\alpha}}, \tilde{\underline{\theta}}) \equiv \log \mathcal{L}(\hat{\underline{\alpha}}(\tilde{\underline{\theta}}), \tilde{\underline{\theta}}) = \max_{\underline{\alpha}} [\log \mathcal{L}(\underline{\alpha}, \tilde{\underline{\theta}})]$$

Reoptimizing over α at fixed $\tilde{\underline{\theta}}$ is necessary. The dimension of the parameter space to associate with this random effects model is K_{re} ,

$$K_{re} = \text{tr}(G) + \ell.$$

We note that $1 \leq \text{tr}(G) \leq k$. This corresponds to the fixed effects model for $\theta_1, \dots, \theta_k$ is intermediate between a random effects model where $\theta_1 = \dots = \theta_k$ is unstructured and a no-effects model wherein $\theta_1 = \dots = \theta_k$.

AIC for the random effects model on $\underline{\theta}$ is $-2 \log \mathcal{L}(\tilde{\underline{\alpha}}, \tilde{\underline{\theta}})$. The small sample version is

$$\text{QAIC}_c = \frac{-2 \log \mathcal{L}(\tilde{\underline{\alpha}}, \tilde{\underline{\theta}})}{\hat{c}} + 2K_{re} + 2 \frac{K_{re}(K_{re} + 1)}{n + K_{re} + 1}$$

If it is justified to take $\hat{c} = 1$, then the above becomes AIC. The application of this methodology is given in Burnham and White (2002). The inference formulae not given here.

Examples are given in Burnham (in review). In particular, a recovery example for which $k = 41$ (units are years); θ_i is a probability. The fixed effects model means fitting 40 means, but just a mean (μ) of the θ_i merely to “capture” average abundance. This might be better represented by a single parameter σ^2 , which is on its own right. The full model also requires 42-band recovery parameters. Results for three models (the r_i are unrestricted) are below. The model is justified because sample size was 42,015 birds banded.

Model	K	ΔQAIC	Akaike weight	survival
$\theta_{\mu, \sigma}$	73.26	0.00	0.9984	random time
θ_t	83	12.87	0.0016	fixed-time
θ	43	100.11	0.0000	time-con

Summary points we want to make from this section are: First, the methodology generalized to random-effects models, ultimately probably given here. Second, the correct measure of parameter dimension for random effects is the trace of the associated projection (smoothing of computing K_{re} is not unique). Extensions of AIC to other models (not simple, fixed effects) models often require computing a matrix.

stood for models that correspond to fixed-effects likelihoods. In Section 6.6.2 it is applicable to, and has been developed for, mixed-effects models. Extensions of AIC-type model selection are with respect to a best approximating model, not a true model. Inference to other contexts are mostly state-of-the-art. It has been successfully, Our point here is that AIC does not deal with mixed-effects models. It has a wider spectrum of applications, so we do not mention here, but without details.

Data smoothing by semi- or non-parametric methods, is an important area in statistics. Hurvitch et al. (1998) provide versions of Akaike's information criterion for parameter selection on nonparametric regression (see also Shi and Tsai 2001). For data vector \underline{y} Hurvitch et al. (1998) note that the fitted values correspond to $\hat{\underline{y}} = H\underline{y}$ for a smoothing matrix H that depends on \underline{y} as a function of a smoothing constant; this constant is not a single traditional parameter. Moreover, for several AIC-type selectors "Each of these selectors depends on H through $\text{tr}(H)$ and can be interpreted as the effective number of parameters used in the model. . . ." (Hurvitch et al. 1998:273). Indeed, the role of the number of parameters is taken here by $\text{tr}(H)$ and then, generally, by $\text{tr}(H)/2$ because there is also one variance parameter, σ^2 . There are extensions for generalized smoothing and generalized cross-validation for generalized additive models (Hastie and Tibshirani 1990:49, 151). In this case the trace of a smoother matrix takes the role of the number of structural parameters.

Other general statistical-modeling methodologies for model selection has been extended include generalized cross-validation (Pan 2001a,b). The method of Pan also serves to extend K-L model selection. Robust regression is sometimes approached with quantile deviations (also called L1 regression, a type of quantile regression). Shi and Tsai (1990a) give the small-sample AIC for L1 regression and call it as L1cAIC. The formula for L1cAIC is very different from AIC. However, one of their conclusions was that both AIC_c and L1cAIC are good. (Shi and Tsai 1990a:263) "... provide good model selections in a linear regression model with double exponential errors. Our recommendation that AIC_c is useful in general. Recent work on model selection for robust and nonparametric regression (Shi and Tsai (1998, 1999), Hurvich and Tsai (1998), and Shi and Tsai (1999)).

A general approach to K-L model selection when the model is unknown remains elusive. However, the recent deviance information criterion (DIC) approach within a Bayesian framework may provide a way (Spiegelhalter et al. 2002). As with AIC the intent of DIC is to

MCVE approach to model fitting. Heuristically, the solution to minimize the posterior expected value of (relative) Kullback-Leibler information loss (see also Thabane and Haq 1999, regarding a general likelihood approach to fitting and selecting among random effects (without integrating out the random effects) within the framework of h-likelihood (Lee and Nelder 1993)). Research along this line seems especially worthwhile. Some results in Section 6.6.2 combined with ideas from the DIC approach that a likelihood solution may be possible (i.e., AIC for general models).

6.7 Selection When Probability Distribution is Unknown by Model

6.7.1 Keep All the Parts

Most model selection focuses on questions about modelable plainable variation in data, within the context of a single probability distribution. For example, if y is the response variable and x_1 to x_6 are predictors, we may assume that the structural aspect of the data is some functional form for $E(y|\text{predictors})$ based on the predictors transformed with impunity. The (so-called) error distribution is the residual $\epsilon = y - E(y|\text{predictors})$, i.e., unexplained variation. Leibler-based model selection allows the error distribution to be normal. For example, we can compare models wherein ϵ has a normal distribution to models wherein ϵ has a lognormal distribution. However, when so doing: No component part of either probability distribution can safely be dropped in forming the likelihoods. Also, the comparison must be based on y for some models and $\log(y)$ for other models (see 2.11.3). We can contrast models for y as normal versus lognormal. This must be done in a correct way.

Denote the model structural aspects by $\mu_i = E(y_i|\text{predictors})$. If y_i depend functionally on some smaller number of parameters, the normal probability distribution, variance homogeneity, and the likelihood, as $\prod g_1(y_i|\mu_i, \sigma)$, is

$$\mathcal{L}_1 = \left[\frac{1}{\sqrt{2\pi}} \right]^n \left[\frac{1}{\sigma} \right]^n \exp \left[-\frac{1}{2} \sum_{i=1}^n \frac{(y_i - \mu_i)^2}{\sigma^2} \right]$$

As long as every model considered is concerned just with the response given the assumed normality of y , then we can drop the σ from the likelihood involving 2π . This is a general principle: If the likelihood

and comparisons based on relative expected K-L distance (and siblings) remain valid.

As another example assume that you think that the probability for y should be a special case of the gamma: $g_2(y|\mu) =$ Then,

$$\mathcal{L}_2 = \left[\prod_{i=1}^n y_i \right] \left[\prod_{i=1}^n \frac{1}{\mu_i} \right]^2 \exp \left[- \sum_{i=1}^n \frac{y_i}{\mu_i} \right]$$

and as long as this is the global likelihood, the leading term is dropped. However, to compare g_1 to g_2 no parts of either

Keep All the Parts to Compare Different Distributions

To compare two models g_1 and g_2 that are based on different distributions (those models may have the same or different μ_i) we must keep in \mathcal{L}_1 and \mathcal{L}_2 all component parts of the underlying probability distributions.

If in both (6.4) and (6.5) we dropped the leading term, the resultant AIC_1 to AIC_2 is invalid: it gives meaningless results. This is because we would have confounded real model “effect” arising just because $\log \left[\frac{1}{\sqrt{2\pi}} \right]^n$ and $\log \left[\prod_{i=1}^n y_i \right]$ are different.

There is a link here to comparing models on data transformed to log space. This also cannot be done directly. For example, rather than comparing normal distribution models for y versus $\log(y)$, we must compare log-normal wherein y can have either a normal or log-normal distribution. To do this is use software for generalized linear models, which we recommend. Under GLM (McCullagh and Nelder 1989) the response is transformed, rather it is the parametric part of the model that is transformed and linked with different assumed “error” distributions.

6.7.2 A Normal Versus Log-Normal Example

To compare a model based on the normal distribution (g_1) to a log-normal distribution (g_2) we have to use the log-normal

$$g_2(y|\theta, \sigma) = \frac{1}{y\sigma\sqrt{2\pi}} \exp \left[-\frac{1}{2} \frac{[\log(y) - \theta]^2}{\sigma^2} \right]$$

(Again, a GLM approach would be doing this for you.) To know the correct probability density form for y being log-normal, in general, is now

$$\mathcal{L}_2 = \left[\prod_{i=1}^n \frac{1}{y_i} \right] \left[\frac{1}{\sqrt{2\pi}} \right]^n \left[\frac{1}{\sigma} \right]^n \exp \left[-\frac{1}{2} \frac{\sum_{i=1}^n [\log(y_i) - \theta]^2}{\sigma^2} \right]$$

be compared based on likelihoods (6.7) and (6.8) for normal distributions, respectively. It will generally make more sense to compare the same structural models; hence the comparison with a plausible probability model for the data. (Note, the primary focus of data analysis is the question of a suitable model; then the issue of what to use for the error distribution part is a minor matter, as long as the distribution chosen is not a terrible one.)

We provide results of a small Monte Carlo study to select between the log-normal model, based on likelihoods in (6.4) and (6.6), and the normal model wherein the μ_i are constant, and hence $\mu_i \equiv \mu$. The results are from either the normal or log-normal model, with the scale parameter fixed at 1 (results are scale invariant). Because both alternative models are being compared, issues of using AIC or AIC_c (or even BIC) are moot. Table 6.14 shows, in terms of selection relative frequency of the normal distribution, the expected (i.e., average) value of the Akaike weight for the normal model. Results in Table 6.14 are accurate to the two decimal places.

This likelihood-based discrimination between the normal and log-normal distributions performs well (Table 6.14). The two distributions are well separated at small values of $E(y)$ (which actually should be interpreted as $E(y/\sigma)$). However, these distributions are increasingly overlapping as $E(y)$ increases. Correspondingly, the two distributions can be distinguished with given sufficient sample size, if they are meaningfully different.

TABLE 6.14. Results from 10,000 Monte Carlo trials to select between the normal and log-normal models, variance(y) = 1 for both distributions; π is the probability of selecting the normal model; $E(w)$ is the average of the Akaike weight for the normal model, given as a percentage.

		E(y), generating model: normal					
n	5		10		25		50
	π	E(w)	π	E(w)	π	E(w)	
10	58	54	53	51	52	50	50
50	77	71	65	57	56	51	50
100	88	83	72	63	59	53	50
500	100	100	92	88	71	62	60

		E(y), generating model: log-normal					
n	5		10		25		50
	π	E(w)	π	E(w)	π	E(w)	
10	38	44	44	49	48	50	49
50	21	29	34	43	43	49	49
100	12	17	27	36	41	47	49
500	0	1	9	13	29	38	39

identical distributions). The π reflect sampling variation, whereas the $E(w)$ reflect inferential uncertainty. Sampling variation and inferential uncertainty are different in general, and numerically not identical in Table 6.14 in some cases. Notice that the average inferential uncertainty is about 100% sampling variation. That is, $E(w)$ is always intermediate between π and $100\pi\%$ in Table 6.14.

It is also worth noting again that these results are equivalent to those arising from AIC or BIC (because $K = 2$ for both models). We are not entitled (pragmatically) to interpret w_i , for large sample sizes, as the probability that model g_i is the K-L best model, but not that it is the true model. We further illustrate this last point about inferring truth, with a simulation from the negative exponential distribution $g_3 = \exp(-y/\lambda)$ with $\lambda = 1$, between the normal and log-normal models. The parameter λ is unknown, but the parameter so results are invariant to its value. Hence, we can compare model selection results for different sample sizes. As with the normal model, the results are for selection of the normal model as the best approximation to the true model.

n	π	$E(w)$
5	24	32
10	17	22
20	10	12
50	2	3
100	0	0

For n greater than about 100 the result is selection of the normal model as the best approximation to the true model, essentially with no inferential uncertainty. Of course, this is only true if we generate the data: Selecting a model with (inferential) certainty that the model is truth.

6.7.3 Comparing Across Several Distributions: Lindsey and Jones (1998)

Lindsey and Jones (1998) gives an example based on observations of white blood cell count per cubic millimeter of blood. This is a type of leukocyte count used to assess the immune system. The data are from 20 patients in remission of leukemia (considered “treatment” here) and 20 patients as controls (their Table I, shown here in Table 6.15). The interest is in the difference in the mean of this white cell count between the two groups of patients. The standard deviations are $\bar{y}_T = 823$, $\bar{y}_C = 522$, $s_T = 566$, and $s_C = 200$. The use of these as the basis of parametric inference is problematic and the case is in favor of distribution-free inference, except that the standard deviation across patients for this type of count is of interest in its own right.

Standard practice has been to assume a normal distribution for the white cell count, test, or preferably, a point estimate and a confidence interval.

Treatment				Control		
171	397	795	1212	116	375	440
257	431	902	1283	151	375	503
288	435	958	1378	192	377	675
295	554	1004	1621	208	410	688
396	568	1104	2415	315	426	700

TABLE 6.16. AIC results for the T_4 blood cell count data (Table 6.15). Jones (1998, Table II), “difference” means a treatment effect, hence more explanation.

Model	No difference			K
	K	AIC	Δ	
normal	2	608.8	22.8	3
log-normal	2	590.1	4.1	3
gamma	2	591.3	5.3	3
inverse Gaussian	2	590.0	4.0	3
Poisson	1	11652.0	11066.0	2
negative binomial	2	589.2	3.2	3

the t -distribution. However, one might postulate several probability distributions as the basis of the model and then use AIC to compute a value for the suitability of each model. As long as we compute AIC based on the complete probability distribution (i.e., not just the observed data), the approach is valid, as noted in Lindsey and Jones (1998). Table 6.15 by Lindsey and Jones (their Table II) for models he considered. The results in Table 6.16. The parameters of these models are either $E(y) = \mu$, or are a dispersion parameter such as for the gamma distribution. Within a distribution there is a pair of models: either $\mu_T = \mu_C$ or $\mu_T \neq \mu_C$. We will show aspects of calculation of three models in Table 6.16.

For the Poisson model when $\mu_T = \mu_C = \mu$ the probability mass function is given by

$$g(y|\mu) = \frac{e^{-\mu} \mu^y}{y!}.$$

Hence the likelihood is to be taken here as (T is the sum of the y_i 's)

$$\mathcal{L}(\mu) = \prod_{i=1}^{40} \frac{e^{-\mu} \mu^{y_i}}{y_i!} = \left[\prod_{i=1}^{40} \frac{1}{y_i!} \right] [e^{-40\mu} \mu^T]$$

$$= \exp(-195150.88) [e^{-40\mu} \mu^{26905}]$$

The huge term involving factorials could be ignored (due to our only model and inference were just about μ (given that we would care only about the ratios $\mathcal{L}(\mu)/\mathcal{L}(\hat{\mu})$). But we need parts to make a comparison of this model (i.e., probability distribution) with a different probability distribution. The MLE is $\hat{\mu} = 26905$, $\log \mathcal{L}(\hat{\mu}) = -5824.9478$ and $\text{AIC} = 11651.896$, rounded to 6.16.

Consider the normal model for the no-difference case. The log-likelihood is given by (6.4). The MLEs are $\hat{\mu} = 672.625$ and $\hat{\sigma}^2 = (464.57)^2$; $\log \mathcal{L}$ at the MLEs is given by

$$\log \mathcal{L} = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\hat{\sigma}^2) - \frac{n}{2} = -3$$

Hence, here $\text{AIC} = -2 \log \mathcal{L} + 4 = 608.8$.

For the normal model with a difference by treatment and dispersion parameter σ^2 , the likelihood is

$$\mathcal{L} = \left[\frac{1}{\sqrt{2\pi}} \right]^{40} \left[\frac{1}{\sigma} \right]^{40} \exp \left[-\frac{1}{2} \frac{\sum_{i=1}^{20} (y_{T_i} - \mu_T)^2 + \sum_{i=1}^{20} (y_{C_i} - \mu_C)^2}{\sigma^2} \right]$$

The MLEs are $\hat{\mu}_T = \bar{y}_T = 823$, $\hat{\mu}_C = \bar{y}_C = 522$ and $\hat{\sigma}^2 = 193151.45 = (439.49)^2$ (SS denotes sum of squares). The log-likelihood at the MLEs, evaluated at the MLEs, reduces to the same as the no-difference model case, but $\hat{\sigma}^2$ is different:

$$\log \mathcal{L} = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log(\hat{\sigma}^2) - \frac{n}{2} = -3$$

$\text{AIC} = 606.364$.

To consider the question of the best approximating model, we make comparisons that confound structural and stochastic elements. So in Table 6.16 we must compare within model pairs (columns). The correct comparison is automatic when we use the same structural model pairs. For example, to compare the normal and log-normal distributions the two evidence ratios are

$$\begin{aligned} 11499 &= \frac{e^{-4.1/2}}{e^{-22.8/2}}, & \text{forcing } \mu_T = \mu_C, \\ 7332 &= \frac{e^{-2.6/2}}{e^{-20.4/2}}, & \text{allowing } \mu_T \text{ and } \mu_C \text{ to vary} \end{aligned}$$

Either one of these ratios is overwhelming evidence against the log-normal distribution, relative to the log-normal, as a best approximating

In practice, to address the distribution issue we would use the difference-allowed cases. There it is clear that the binomial, but the log-normal (evidence ratio in favor of the binomial is 3.7), gamma (2.7), and inverse Gaussian (3.0) are plausible. Akaike weights for this subset of four models (ordered as normal, binomial, gamma, inverse Gaussian) are 0.14, 0.18, 0.17, 0.51.

A final observation here is that inference about the “treatment effect” is much affected by choice of model, for models that fit well (the normal is a terrible fit). For the other model pairs, evidence ratios are

model	ER for an effect
normal	3.3
log-normal	2.1
gamma	2.7
inverse Gaussian	2.5
negative binomial	5.0

How would we analyze these data for an effect? First, we would fit the no-effect models. Second, we would reparametrize each model from μ_T and μ_C to $\delta = \mu_T - \mu_C$ and $\gamma = \mu_T + \mu_C$. This leads to six reparametrized models. For each reparametrized model leads to $\hat{\delta}$, $\widehat{\text{se}}(\hat{\delta}|g)$, and Akaike weights. For the six models (those w_i are 0.000, 0.138, 0.186, 0.169, 0.186, 0.186) we would base inference about δ on model averaging (no model selection requires substantially more data).

6.8 Lessons from the Literature and Other

6.8.1 Use AIC_c , Not AIC , with Small Sample Sizes

It is far too common that papers examining AIC , by itself or in comparison, fail to use AIC_c when the latter must be used because the number of parameters, at least for some models considered, is not small relative to the sample size. For example, Chatfield (1996) considered model selection in a time series example with $n = 132$ and $R = 12$ a priori designated predictors. The number of models ranged from 6 to 61. Overall we commend the paper; however, for example the conclusion that AIC performed poorly is misleading. AIC_c performed poorly; but it is well known, documented, and commented on K-L-based model selection that in such an example it performs poorly. AIC_c , not AIC (e.g., Sakamota et al. 1986, Bozdogan 1989, Hurvich et al. 1990).

Consider for $n = 132$ the effect of the bias-correction term for these two models with likelihoods denoted by \mathcal{L}_6 and

K	AIC	AIC _c
6	$-2 \log(\mathcal{L}_6) + 12$	$-2 \log(\mathcal{L}_6) + 12.6$
61	$-2 \log(\mathcal{L}_{61}) + 122$	$-2 \log(\mathcal{L}_{61}) + 230$

The difference here between AIC and AIC_c is huge for $K = 61$ and greatly affect which model is selected.

We present in Table 6.17 the results that Chatfield (1996) presented as regards K-L-based model selection versus BIC. We use Δ_i values, not absolute values of these model selection criteria. To show for comparison the Δ AIC values implicitly used by Chatfield, the results in Table 6.17 are based on the results in Table 1 of Chatfield (1996). The nature of the models need not concern us, so we label them 1 through 12. We keep them in the same order as used in Table 1 of Chatfield (1996). The results in Table 6.17 for the AIC_c criterion do have here the interpretation described in Section 2.6. Those interpretations are not true for the Δ_i derived here from AIC, because for large K relative to n , AIC is a biased estimator of the expected K-L distance.

As noted by Chatfield, in this example AIC and BIC lead to different selected models. However, AIC_c (which must be used here for model selection) gives acceptable results. In fact, the Akaike weights here (see Table 6.18) for only four fitted models have any plausibility in this set of models.

TABLE 6.17. The Δ AIC_c that must be used for K-L model selection, as considered in Table 1 of Chatfield (1996), and corresponding Δ AIC values. Also, the Akaike weights based on AIC_c.

Model	K	Δ AIC _c	Δ AIC	Δ BIC	w_i
1	6	4.5	68.0	0.0	0.000
2	11	0.0	62.0	13.0	0.000
3	21	2.5	58.3	47.0	0.000
4	9	86.1	148.8	92.2	0.000
5	17	98.6	157.4	130.9	0.000
6	41	155.5	181.4	246.0	0.000
7	11	83.7	145.7	96.7	0.000
8	21	94.7	150.5	139.3	0.000
9	13	0.5	61.6	20.0	0.000
10	25	10.1	62.0	65.7	0.000
11	31	15.9	60.3	86.8	0.000
12	61	43.9	0.0	139.8	0.000

definitely eliminates 6 of the 12 models, because relative to the set, they are extremely implausible. Also, most of the values put on models with low K . Based on comparing AIC to BIC, we concluded that BIC was a better criterion than AIC. This is justified in that here one must compare AIC_c to BIC for a fair test of K-L information-theoretic model selection versus BIC.

6.8.2 Use AIC_c , Not AIC, When K Is Large

Leirs et al. (1997) report the analysis of an extensive set of data from Tanzania on the rat *Mastomys natalensis*. The data analysis was to examine factors potentially important to the survival dynamics of the species. The data were collected in 1986 and February 1989 by live trapping on a 1 ha grid of 25 positions (several traps per position). There were three capture-trapping each month (hence 29 primary trapping periods) and a total of 6,728 captures of 2,481 individual animals. We take the sample size to be the latter, i.e., $n = 2,481$. Leirs et al. (1997) considered six a priori models to represent how environmental (rainfall) and density factors might affect survival probabilities (S) and subadults maturing to adults (ψ). Capture probabilities (p) were assumed to be constant. The number of parameters in these models was 6. Data analysis was by ML methods for capture–recapture models (see, e.g., Brownie et al. 1993, 1995) with incorporation of covariates for rainfall and population density. The goodness-of-fit of the global model was quite acceptable.

The global model used by Leirs et al. (1997) allows full temporal variation in all model parameters (hence S and p vary by time, and ψ varies by time). In their Table 1, this is model g_1 with 6 parameters. Their model g_2 is the most restricted model: no temporal variation in parameters (S and p vary by age only, subadult versus adult, and one maturation probability parameter, ψ). Model 3 allows p to vary by time, but S and ψ are not time-varying. Model 4 allows structured time variation in the three classes of parameters. The most complex models are based on population dynamics models. Model 5 is a general capture–recapture model. Model 4 has temporal variation in p as functions only of population density (internal factors regulating capture); Model 5 has functions only of rainfall (i.e., external factors regulating capture); Model 6 has temporal variation in p as functions of both population density and rainfall. Leirs et al. (1997) used AIC_c for model selection. In Table 6.18 we present the results of their analyses, as well as results for AIC_c , which we compare to the AIC weights.

Model	K	AIC results		Δ_i
		Δ_i	w_i	
1 Global	113	0.0	0.99	0
2 No effects	5	540.7	0.00	529
3 No dynamics	49	207.4	0.00	198
4 Density effects	52	205.8	0.00	197
5 Rainfall effects	55	25.9	0.00	17
6 Rainfall and Density	64	9.7	0.01	2

Even though the sample size is large here (2,481), the model with 113 parameters means that AIC_c should be used (the term added to AIC to get AIC_c is $(2 \times 113 \times 114)/2367 = 10.9$, is not trivial. Clearly, using a different interpretation of the relative evidence for model selection using AIC, Leirs et al. (1997) opted to select model g_6 (the model and therefore to infer that there were population dynamics could be substantially explained only by both external (rainfall) (population density) factors. They worried some about the model (Nichols, personal communication); they did not need to. AIC_c should be done here, model g_6 is a tenable model.

6.8.3 When Is AIC_c Suitable: A Gamma Distribution

The K-L approach to model selection is exact, philosophically sound (sample size (i.e., it is not intrinsically asymptotic). However, variation in the face of truth being unknown means some degradation to the target model selection criterion, usually large-sample selection. To get a practical estimator of this criterion, the simplest is to use asymptotic results, which produces TIC; AIC is a practical implementation of TIC. Simulation studies and experience suggest these “large sample” formulas will perform very poorly when n is small. Useful insights to small sample model selection are obtained by deriving exact versions of the target criterion for various assumed-true models (Section 7.4). The results here are expressible as AIC + one or more terms of the form $h(K)/(n^a)$ and functions $h(\cdot)$ and $m(\cdot)$.

The term added to AIC to get AIC_c is just $2K(K + 1)/n$. It is not unique, AIC_c is especially compelling as an omnibus version of AIC because essentially it requires only that the likelihood be proportional to a normal distribution. This will be a good approximation at quite small sample sizes if the sample elements are (i.i.d.)

needed adjustment term to AIC for the gamma distribution is $2K(K+1)/(n-K-1)$. We strongly recommend using a version of AIC when K is not large relative to n ; use AIC_c unless n is known.

Theoretically, no small sample adjustment term is needed for example regression when the residual variation, σ^2 , is estimated. However, performance of model selection will be improved in this case by use of AIC_c . The added term basically prevents overfitting from reaching and exceeding n , which must somehow be avoided.

The theoretical bias of AIC (Section 7.2) is given by

$$2 \left(E_x \log[\mathcal{L}(\hat{\theta}(x)|x)] - E_x E_y \log[\mathcal{L}(\hat{\theta}(x)|y)] \right),$$

For the normal distribution model as truth, and with $g = f$, the bias is

$$\frac{2K(K+1)}{n-K-1}.$$

We evaluated a simpler form of (6.7) for $f = g$ under the normal distribution using Monte Carlo methods to compare those results to (6.8). Table 6.1 shows the results of these Monte Carlo evaluations; however, we first need to discuss the details of what was done.

The gamma model probability density function is

$$g(x|\alpha, \beta) = \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha)\beta^\alpha}, \quad 0 < x, 0 < \alpha, \beta < \infty.$$

Here $E(x) = \alpha\beta$ and $\text{var}(x) = \alpha\beta^2$. Because β is a scale parameter, it is invariant to the value of α . Thus it sufficed to just set $\alpha = 1$ (we still must estimate β from the data).

The gamma distribution is in the exponential family. The simpler form of (6.7) is

$$2\text{tr}[\text{COV}(\hat{\theta}, S)] - 2K,$$

where $\text{COV}(\hat{\theta}, S)$ is here a 2×2 matrix and θ and S are the parameters of the gamma distribution. The canonical form for the gamma distribution is

$$g(x) = \exp \left[x \left[-\frac{1}{\beta} \right] + \log(x)(\alpha) - \log(x) - \log(\Gamma(\alpha)) \right]$$

Hence, $\theta_1 = -1/\beta$ and $\theta_2 = \alpha$ are a 1-1 transformation. The sufficient statistic, S (in 6.9), has the components $\sum x$ and $\sum \log(x)$.

For given α , a random sample of size n was generated from the gamma distribution, and then S was computed and $\hat{\theta}$ found by maximizing the log-likelihood. This was repeated for 100,000 independent trials.

	Sample Size, n			
α	10	20	50	100
0.25	2.78	1.06	0.34	0.16
0.50	2.03	0.80	0.28	0.14
0.75	1.93	0.77	0.28	0.13
1.00	1.83	0.76	0.25	0.13
5.00	1.74	0.66	0.26	0.14
10.00	1.73	0.71	0.27	0.13
25.00	1.68	0.69	0.28	0.14
50.00	1.71	0.69	0.24	0.13
100.00	1.75	0.71	0.27	0.13
From (6.8)	1.71	0.71	0.26	0.12

was found, and finally (6.9) was computed. Results, and Δ_6 for α , are shown in Table 6.19, along with the value of (6.8) for $\alpha = 0.25$. The value based on (6.9) is reliable to (almost) two decimal places. From Table 6.19 that AIC_c is very adequate for use with the $K = 2$ model (at $K = 2$), except for when both α and n are both small. The case $\alpha = 0.25$ and $n = 10$ is quite extreme in that the underlying distribution is very skewed, hence the likelihood is not near to having a normal form at sample size 10. However, even for $n = 10$ the approximation is good at $\alpha = 1$, which corresponds to the negative exponential distribution.

More research is desirable on the issue of small-sample size. The issue on the general suitability of AIC_c . However, this example does not have examined support AIC_c as generally suitable unless the underlying probability distribution (for a single sample) is extremely nonnormal in terms of being strongly skewed.

6.8.4 Inference from a Less Than Best Model

We continue with some ideas, exemplified by the example of Leirs et al. (1997), about inference from other than the K-L best model. In circumstances this is justified, especially if (1) the model $g_{(2)}$ (notation for the second-best AIC model), used for inference, is not the best model, $g_{(1)}$; and (2) the unexplained “effects” in the data are due by the additional parameters added to model $g_{(2)}$ to generate the model, small relative to the explained effects represented by model $g_{(1)}$. This assumes a parametrization of the models as $g(\underline{x} | \underline{\theta})$, with model $g_{(2)}$ arising under the imposed constraint $\underline{\theta}_1 = \underline{\theta}_2$.

We elaborate these ideas further using the Leirs et al. (1997) example. The AIC_c best model (g_1) was not interpretable in its entirety, but the second-best model (g_6) was interpretable, and because $\Delta_6 = 2.3$, that model was a good model for the data. Moreover, that second-best model is not

additive to (and preferably orthogonal to) the 64 parameters of the best fitted model. Therefore, their results provide overwhelming support (for the set of models used) for the joint importance of rainfall and temperature as at least good predictors of the observed population growth and capture probabilities (if not outright support for a causal role for these variables).

In choosing to make inferences based on model $g_{(2)}$ (ignoring model $g_{(1)}$ (their model 1), Leirs et al. (1997) argue that they cannot interpret the meaning of the additional 49 parameters that constitute the difference between their best and second-best model, but this does not in any way invalidate inference from the second-best model. This is because $g_{(2)}$ is nested within $g_{(1)}$. This sort of argument holds for all models are nested.

The only pressing concern here, in ignoring the best model, is that the 49 “effects” defining the difference here between the two models, is the issue of the relative magnitude of the two effects. In an analysis of variance terms this issue is about the partitioning of effects represented by the difference in their fitted models. It is a sum of squares for effects of g_6 versus g_2 plus a sum of squares for effects of g_1 versus g_6 . Analogous to ANOVA, we can use here ANODEV (see, e.g., McCullagh and Nelder 1989, Skellam 1991) to accomplish a useful partition.

In this example, ANODEV proceeds as follows to measure the relative importance of the ignored effects left unexplained in model g_6 versus explained effects in model g_2 . First, some baseline “model” is needed; here that baseline is model g_2 of Leirs et al. (1997). Since $g_2 \subset g_6 \subset g_1$ and corresponding values of K : 5, 64, and 108, we proceed by obtaining the log-likelihood values and computing the difference of total deviance of model g_2 versus g_1 as

$$\begin{aligned} \left[2 \log(\mathcal{L}(\hat{\theta} \mid g_1)) - 2 \log(\mathcal{L}(\hat{\theta} \mid g_2)) \right] &= \left[2 \log(\mathcal{L}(\hat{\theta} \mid g_6)) - 2 \log(\mathcal{L}(\hat{\theta} \mid g_2)) \right] \\ &\quad + \left[2 \log(\mathcal{L}(\hat{\theta} \mid g_1)) - 2 \log(\mathcal{L}(\hat{\theta} \mid g_6)) \right] \end{aligned}$$

The result here is $756.8 = 649.0 + 107.8$. The above three values are also interpretable as likelihood ratio test statistics on 108, 59, and 49 df, respectively.

The above partitions a measure of the magnitude of the difference between the total deviance (on 108 df) represented by fitted model g_1 into a measure of the difference between the deviance explained by model g_6 alone (649.0, on 59 df), plus the additional deviance explained (107.8, on 49 df) explained by the added 49 parameters of model g_1 from model g_6 . Based on this partition we can define

$$K^* = \frac{2 \log(\mathcal{L}(\hat{\theta} | g_1)) - 2 \log(\mathcal{L}(\hat{\theta} | g_2))}{756} = 756.$$

The interpretation is that 86% of the total structural information variation in model g_1 is contained in model g_6 . Thus, information potentially interpretable effects has been lost by making inference on model g_6 (i.e., the second-best AIC_c model), rather than on model g_1 . However, that other 14% of information was left as not judged to be real information, as evidenced by AIC_c selecting the best model, but ignoring it does not invalidate the inference on model g_6 .

Clearly, the addition to model g_6 of all the structure of the additional 49 parameters (to get model g_1) does, for the data, produce a K-L best-fitted model. However, in principle there is some information between models g_6 and g_1 , that adds far fewer than 49 parameters to produce an even smaller AIC_c than model g_1 . Such an additional model can extract additional useful information from the data; it might be a random-effects model, or some interaction effect of rainfall and population density. The situation faced here is, essentially, considered in sections 3.5.6, and 6.9.3, where we point out that if there are two models in the other and differing by a large number of parameters, then anomalies can arise in data analysis based on K-L model selection.

In general, there are situations where choosing to make inference on other than the AIC_c best model can be justified. However, it is not satisfied if the AIC_c best model has many additional parameters to the model one uses for the basis of inference. If we find a better situation, it suggests that we did not think hard enough about the choice of models, because we probably left out at least one good model a posteriori (to the initial data analysis) model building process; just admit, then, which models were a priori and which were chosen by initial data analyses.

6.8.5 *Are Parameters Real?*

Consideration of what is a parameter seems important, but has been focused entirely on parametric models. With only one class of model, regard a parameter as a hypothetical construct. Hence, a parameter is an embodiment of a concept and does not have the reality of a measurable variable. As such, a parameter in a statistical setting is (usually) virtually essential, conceptual abstraction based on the frequency of the expected value of a measurable variable that is not measurable. There also needs to be a large number of actual occurrences of the measurable variable, or at least a well-defined concept of the number of occurrences. Then the concept of an average

one, or a few, simple measurements. There is no instrument to record the exact value of a parameter used in a statistical model. An exception occurs in measurement error models where the parameter is real but becomes the parameter of interest because each recorded value is recognized to be imprecise at a nonignorable level of imprecision.

We go a step further and recognize two classes of parameters in statistical models: (1) parameters that appear in the log-likelihood; that is, they have any associated physical or biological reality; and (2) parameters above that are directly related to expectations of measurable variables. The second class of parameters are tied to measurable events and need not appear in the likelihood (they often do appear).

As an example, consider the analysis of cohort survival data represented by examples in Section 5.2. The age-specific survival probabilities S_r cannot be directly measured (such as the probability of survival can be). However, the concept represented by S_r has clear meaning: survival of an animal over a defined time interval. This can be repeated based on a sample of animals (from a large, random, infinite, population of animals). These survival probabilities are in the second class of parameters above. To provide both a useful and parsimonious estimation of this set of parameters from limited data, ecological science adopts smooth, deterministic parametric models (“models” for short) such as

$$S_r = \frac{1}{1 + \exp[-(\theta_1 + \theta_2 \cdot r + \theta_3 \cdot r^2)]}$$

(as emphasized in this book, we should not pretend that this model holds). The parameters θ_1 , θ_2 , and θ_3 appear in the likelihood function. These parameters are in our first class of parameters above: they do not have any direct physical or biological reality. In this case, S_r is very useful in making parsimonious predictions of the S_r values. The derived parameters based on the interpretable and parsimonious model. Often, interpretability is as important as parsimony. It is not that the two criteria of model usefulness are complementary; they conflict. (Interpretability is a subject-matter criterion, not a statistical one; we have not focused on it here).

The relationship of a parameter to prediction and expectation (and themselves concepts) is straightforward in a simple linear model.

$$E(y | x) = \beta_0 + \beta_1 x.$$

If we can measure the values of y when separately x and $x+1$ (or be able to control x), then

$$\beta_1 = E(y | x + 1) - E(y | x).$$

y can be discovered by direct measurement. As Mayr (1982) argues, hypotheses are often the driving force in science, much more so than data. This is especially true in biology. The concept of parametric models in statistical science is a particularly powerful force.

6.8.6 *Sample Size Is Often Not a Simple Issue*

Students are introduced to statistical data using the concept of a single response variable, y , and possibly some explanatory variable, x . This simple data structure fails to convey the possible complexities of real data, especially when that structure is not describable by a simple model. For example, there may be a sample n_s of subjects, and on each subject n_t measurements are taken at n_t time points. To then claim that the total sample size is $n = n_s \times n_t$ can be misleading; it is only defensible under a certain inference model. However, if subjects are treated as random effects, there is one sample size for subjects and a different sample size for time points within subjects. This latter case creates difficulties in the computation of AIC_c and thus the sample size and thus the computation of AIC_c .

A related example is the Durban storm data (Section 4.4). The sample size as all 2,474 weeks in which a storm event occurred. One might argue that the sample size should be the number of storms, or maybe n should be taken as 52, the number of weeks in a year. This makes some sense because if we knew the weekly storm probability, then $n = 52$ would be correct (we would still use a model with weekly storm probabilities). Results of the Durban storm analysis are slightly different if we use $n = 52$ and hence use $QAIC_c$. There is no definitive basis to know the “correct” sample size.

For some data structures there may be two (or more) defensible sample sizes. In line transect sampling (simple designs) there will be a sample of k transect lines (often k ranges 10 to 30) and a sample of n detected animals on each line, with total length L . For modeling the detection function, the sample size is n , whereas k is the relevant sample size for estimating the variation of encounter rate, n/L . For capture–recapture models, n_s for the number of distinct animals captured once or more, and n_r for the potential number of recapture events. There are logical arguments for either of these as the sample size to use for AIC_c (see Burnham and Anderson, 2002, for choice now). Further thought suggests that the sample size for survival parameters (under fixed-effects modeling) should be n_s , while for the capture probability parameters it should be n_r . In general, one can envision different sample sizes for different subsets of the parameters in models. Certainly, if some parameters are treated as random effects.

issue. Under a Bayesian approach Spiegelhalter et al. (2002) propose a model selection statistic called DIC that seems to behave well with complex sample-size structure. A similar approach is possible under likelihood-based inference.

The issue of sample size can be complex and has implications as AIC_c and $QAIC_c$. We do not pursue solutions here, but leave this issue as a future research area.

6.8.7 *Judgment Has a Role*

We have seen published statements to the effect that formal (including automated or mathematical) model selection methods are used presumably because they do not allow for judgment and thus avoid an inference based on a single model. Apparently, a proper approach would want to select a model by some ill-defined approach and possibly do some form of ill-defined multimodel inference with such thinking, even though we value and respect this thinking in science and data analyses. To paraphrase the first sentence in Stewart-Oaten (1995), statistical analyses are based on informal, not on formal, knowledge and a (hopefully) deep understanding of the limitations of the formal data analysis methods used. The common use of formal (“objective”) methods with the admission that judgment is needed for data analysis and interpretation of results.

Judgment does play an important role in data analysis. It is needed for the decisions on the general approach and the specific analysis that will be used. Also, a great deal of judgment may be required to select a priori set of models to be considered. However, once this set is established, then we maintain that a formal and well-grounded (including automated or mathematically) model selection methodology can be used to select these R models. The result will be at least an objectively selected best model, and preferably, model credibility will be achieved by full multimodel inference. This level of inferential model selection can be achieved by objective, criterion-based methods such as AIC_c . Interval estimation of a parameter given a model surely may be argued that inference should be just a matter of judgment.

Given a model, the statisticians insist on exactly this: a well-founded, objective criterion, or approach, that is precise in the mathematics and that, for given data, gives the same result to all who do the analysis. After various such formal analyses, the conclusions and recommendations may indeed be the result of judgment to the situation. But at least the formal results can be argued to be of objectivity. Just as we insist on a well-founded objectivity.

multimodel inference). There is room for subsequent progress about a final best model or inference but only after the relevant selection methods have been presented and defended as reasonable.

The other aspect of the critique of formal-objective (i.e. selection) that is sometimes heard is that it forces one to select. This is a valid concern; it is also now mute. Formal model selection methods will greatly reduce the need for what are difficult means to cope with the obvious inadequacies of methods for a single best model.

The real issue about the use of judgment is not whether but rather who is qualified to exercise good judgment, and whether they know they are qualified? Good judgment can be taught. However, it takes a lot of training and experience to achieve (as opposed to say medicine), we do not get frequent advice about whether our judgments (decisions) are good ones.

6.9 Tidbits About AIC

The section contains miscellaneous ideas and results from elsewhere, but are worth understanding.

6.9.1 Irrelevance of Between-Sample Variation

Likelihood-based inference, including AIC-based model selection, is concerned with relative evidence about alternatives, conditional on the data. For AIC the technical nature of the inference is about comparing relative K-L information loss, as bias-adjusted log-likelihoods, for different models; all such comparisons must be based on a single sample, neither meaningful nor valid to compare individual AIC values between different samples. However, people have seen the data and noted the sample-to-sample variation in the values of the fitted model and concluded that our guidelines about AIC could not be correct. This “insight” is wrong because sample-to-sample variation for a single model (rather than for Δ) is totally irrelevant to inference about alternative models given the data. The situation is like that to a randomized complete block (RCB) experiment where the treatments are estimated only from within-block differences.

This error is common enough that we illustrate it here. Consider a linear model $y_i = \mathbf{x}_i' \boldsymbol{\beta} + \epsilon_i$, $\epsilon \sim \text{normal}(0, \sigma^2)$ where $\boldsymbol{\beta}$ has unknown parameters. For a random sample under this model, conditional on the data, the likelihood is

$$-2\log(\mathcal{L}) = n \log \left[\frac{2\pi\sigma^2}{n} \right] \frac{1}{(n-p)}$$

where RSS is the residual sum of squares. As a random variable, $\text{RSS}/(n-p)\sigma^2$ is central chi-square on $\nu = n-p$ degrees of freedom.

$$-2\log(\mathcal{L}) = n \log(\chi^2_\nu) + n \log \left[2\pi\sigma^2 \left(\frac{n-p}{n} \right) \right]$$

Thus, ignoring the additive constant, the sample-to-sample variation in the maximized $-2\log$ -likelihood is easy to simulate in this case. The variation depends strongly on sample size. For example, for $n=100$, 500, or 5000, Monte Carlo repetitions of this model are shown below for $-2\log(\mathcal{L})$:

sample %-tile	n	
	100	500
1	412.6	3,021.2
10	429.3	3,056.1
25	438.9	3,075.3
50	449.0	3,096.5
75	458.7	3,117.7
90	467.8	3,136.9
99	482.5	3,168.5

The “absolute” variation over samples for a given model is much larger than our guidelines, such as a $\Delta \geq 10$ being generally strong evidence for the model with the bigger AIC. But this sampling variation in the AIC for a single model is meaningless and misleading. The only relevant evidence about alternative models are differences over models. The variation of such a difference is not related to, and much smaller than, the (irrelevant) sampling variation of a single AIC value. The treatment effect in an RCB experiment can be inferred from the difference in AIC values, not from the variation of a given treatment.

A related misconception arises because AIC values can be very close together. Sometimes authors and analysts have AIC values of 4999 and 5000 for three models under consideration and conclude that the models are a short distance apart and “one model is nearly as good as the other.” This is a poor interpretation and is probably influenced by the small sample size that contributes to the fact that AIC values in this case are very close together. The focus of attention must always be on the difference in AIC values, the Δ_j and the associated Akaike weights, w_j , and the relative support for the models based on these w_j .

relationship between the G-statistic for goodness-of-fit and distance. The G-statistic is usually written as

$$G = 2 \sum_{j=1}^k O_j \log \left(\frac{O_j}{E_j} \right),$$

where O_j is the observed count and E_j is the expectation under the model. Under mild conditions, G is asymptotically distributed under the null hypothesis that the model is an adequate fit to the data. Such G-statistics are additive, whereas the more traditional goodness-of-fit test statistic

$$\text{Pearson} = \sum_{j=1}^k ((O_j - E_j)^2 / E_j)$$

is not. The K-L distance for discrete data is written as

$$I(f, g) = \sum_{i=1}^k p_i \log \left(\frac{p_i}{\pi_i} \right)$$

and is almost identical in form to the G-statistic.

Given a sample of count data n_1, \dots, n_k ($n = \sum n_j$) which correspond to the observed relative frequencies. Denote the true probabilities under the approximating model by $\hat{\pi}_j(\theta)$; then in the discrete case, we have $0 < p_i < 1$, $0 < \pi_i < 1$, and the sum to 1, as do their estimators. Then $I(\hat{f}, \hat{g})$ can be rewritten as

$$\sum_{j=1}^k (n_j/n) \log \left(\frac{n_j/n}{E_j/n} \right).$$

Now K-L distance between these (estimated) distributions is

$$\frac{1}{n} \sum_{j=1}^k (n_j) \log \left(\frac{n_j}{E_j} \right),$$

or

$$\frac{1}{n} \sum_{j=1}^k O_j \log \left(\frac{O_j}{E_j} \right).$$

Thus, the G-statistic and K-L information differ by a constant factor, i.e., in this context, $G = 2n \cdot I(\hat{f}, \hat{g})$. Similar relationships exist for information expectations of likelihood ratio statistics for continuous data (a likelihood ratio test (LRT) for discrete data). Thus, the LRT is directly related to the K-L distance.

The use of the information-theoretic criteria in model selection is different from that of hypothesis testing, and this is an important point to understand. These differences can be illustrated by considering a set of candidate models, each successive model differing by one parameter. Model g_i is the null model with i parameters, and model g_{i+j} is the model with $i + j$ parameters. Model i is nested within model $i + j$; tests (LRT) can be used to compare the null model with alternative models g_{i+j} , where $j \geq 1$. Thus, if model g_i has 12 parameters, g_{i+1} has 13, g_{i+2} has 14, and so on.

This concept of a set of nested models is useful in illustrating the differences between AIC versus LRT for model selection. First, the AIC value for each of the models is exactly the same; thus, model g_i has more support than any other model. Second, in each test, the null hypothesis is model g_i and assume that it is an adequate model. Then, we entertain a set of alternative hypotheses, models g_{i+j} , hypothesized to offer a “significantly” better explanation of the data. Model g_i (the null) is tested individually against the $j \geq 1$ alternative models. The first test statistic (g_i versus g_{i+1}) here is assumed to follow a χ^2 with 1 df, while the second test statistic (g_i versus g_{i+2}) follows a χ^2 distribution with 2 df, and so on. The following relations

$$\begin{aligned} \text{AIC}_i &= -2 \log(\mathcal{L}_i) + 2i, \\ \text{AIC}_{i+j} &= -2 \log(\mathcal{L}_{i+j}) + 2(i + j), \\ \text{LRT} &= -2(\log(\mathcal{L}_i) - \log(\mathcal{L}_{i+j})) \text{ with } j \text{ df} \end{aligned}$$

Then, in general,

$$\text{LRT} = \text{AIC}_i - \text{AIC}_{i+j} + 2j.$$

Now, for illustration of a point about the difference between AIC and LRT for model selection, assume

$$\text{AIC}_i \equiv \text{AIC}_{i+j}.$$

If this boundary condition were to occur (where K-L selection is indifferent to the model), then we would have,

$$\text{LRT} = 2j \text{ on } j \text{ degrees of freedom.}$$

Now, a difference of 1 df between g_i and g_{i+1} corresponds to a P -value of 0.157 (Table 6.20). Similarly, a difference of 4 df ($j = 4$) between g_i and g_{i+4} corresponds to a χ^2 value of 0.092. If the degrees of freedom is less than about 7 (a condition where hypothesis-testing methods support the null model), then hypothesis-testing methods support the null model over the alternative models ($g_{i+1}, g_{i+2}, g_{i+3}, \dots$) (Table 6.20). This is not the case with AIC selection, where in this example all the models are

j	χ^2	P
1	2	0.157
2	4	0.135
3	6	0.112
4	8	0.092
5	10	0.075
6	12	0.062
7	14	0.051
8	16	0.042
9	18	0.035
10	20	0.029
15	30	0.012
20	40	0.005
25	50	0.005
30	60	0.001

Test results change in this scenario when there are more estimable parameters in the alternative model (Table 6.20). If (g_i) is rejected with increasing strength since the alternative model has more parameters. For example, the likelihood ratio test of g_i versus g_{i+10} has 10 df, $\chi^2 = 20$, and $P = 0.029$. More strongly, g_i versus g_{i+30} , which has 30 df, $\chi^2 = 60$, and $P = 0.001$, and the P -value is the same for all the models (the null and the various alternatives). In these cases (i.e., > 8 parameters difference between the null and the alternative model), the testing method indicates increasingly strong support for the alternative model with many parameters and strong rejection of the simple model. (See Sakamoto 1991 and Sakamoto and Akaike 1978:196 for a discussion of this issue).

More extreme differences between the two approaches arise if we let $\text{AIC}_i = \text{AIC}_{i+j} - x$ for x in the range of about 0 to 10. To work with the Δ_{i+j} values; then relative to the selected model g_{i+j} is x . If $x = 4$, the choice of model g_i is competitive with g_{i+j} of nested models, as judged by AIC. For comparison, let $x = 2j - x$. Let $x = 4$ and $j = 20$; then the LRT statistic has 20 df, $\chi^2 = 40$, and $P = 0.0154$. Most would take this P -value as compelling support for model g_{i+j} . Thus, AIC can clearly support the simple model g_i and can clearly support model g_{i+j} with 20 additional parameters. This dilemma is entirely a matter of which the model selection method has a sound theoretical basis: Information criteria based on Kullback-Leibler likelihood ratio testing does not.

and may even believe that AIC “loses power” as the difference between models increases beyond about 7. [Note: The correct answer is that AIC has no utility in the information-theoretic approach because it is not a proper scoring rule.] Akaike (1974) noted, “The use of a fixed level of significance for the comparison of models with various numbers of parameters is not appropriate because it does not take into account the increase of the variability of the test statistic when the number of parameters is increased.” The α -level of significance is a function of sample size and the degrees of freedom if hypothesis testing is used as a basis for model selection (see Akaike 1974). However, the α -level is usually kept fixed, regardless of sample size or degrees of freedom, in the hypothesis testing approach. This practice of using a fixed α -level constant corresponds to asymptotically inconsistent hypothesis testing. For example, if the null hypothesis is true and the significance level is, say, 0.05, then even as the degrees of freedom approach ∞ , the probability of rejecting the null hypothesis, even with near-optimal tests, does not converge to 0. The inconsistency is that statistical procedures in this situation do not converge on truth with probability 1 as $n \rightarrow \infty$.

6.9.4 A Subtle Model Selection Bias Issue

Consider having 10 independent one degree-of-freedom random variables, denoted as x_1 to x_{10} . Let z be the minimum of these variables, $\{x_i, i = 1, \dots, 10\}$. Because of the selection of z as a minimum of 10 independent variables, z is not distributed as central chi-square on 1 df. It is distributed as a central chi-square on 1 df, but it is stochastically smaller (“biased”) compared to a central chi-square on 1 df. The selection process induces what may be considered a selection bias. A few authors have expressed concern for a similar sort of selection bias when using AIC as a random variable; they are both right, yet mostly wrong. The selection process as such. By right, we mean that if one focuses on the true K-L-based criterion, T (Section 7.2), to be estimated, then the selection process is unbiased (good) model we can indeed have $E(\hat{T}) = T$. Thus \hat{T} is a nearly unbiased estimator of T (which is related to relative loss) if only one model is considered. However, if we have multiple models, we compute \hat{T} (i.e., AIC) for each model and then select the smallest \hat{T} , then there is a bias by this selection process as regards \hat{T} for the selected model.

Technically, just assume $E(\hat{T}_j) = T_j$, where the expectation is over all possible samples. If model j produces the minimum \hat{T}_j , then \hat{T}_j is the minimum of all \hat{T}_i and we can define $E(\hat{T}_j | \min = j)$. The selection bias that arises is that we mistakenly worry about occurs because $E(\hat{T}_j | \min = j) < E(\hat{T}_j)$.

However, the situation of AIC is not really analogous to the situation of hypothesis testing. For example for two reasons. First, the same data are used to estimate T and to compute \hat{T} (i.e., each \hat{T}); this induces a strong positive correlation in the estimates of T and \hat{T} .

then the $10 \log$ would all be equal, so selecting their minimum would be a bias; the selected variable would always be just the original random variable.

Second, we do not care about a single AIC; we only care about Δ AIC. Most of any bias induced by selection (i.e., ordering) will drop out of these differences. In particular, let the Δ AIC values) smallest to largest be AIC_1 to AIC_{10} . Then we can select the correct (not biased) for models close to each other in Δ AIC, producing on average small Δ , such as expected for Δ AIC. Here, one model was selected as best and the other as second best. We expect that the selection biases are very similar, and hence cancel out of the difference. Moreover, it is highly relevant that the same model is fitting both (all) models; only the models differ, not the data. Again, this means potential selection bias in an individual AIC is irrelevant. Rather, the possible selection bias noted herein stems from the differences, Δ AIC. Because only these differences matter, the possible biasing of a single AIC by selection of the minimum is an almost irrelevant issue. Possible selection bias of the K-L differences is relevant, but should be quite small for models close to each other in K-L information loss, especially if R is not large.

6.9.5 The Dimensional Unit of AIC

The dimensional unit of AIC is the unit of $-\log(g(x|\theta))$, where g is a probability distribution, not a pdf (which is for a continuous random variable). Thus, g is strictly for a discrete random variable. From a philosophical point of view all random variables have a countable number of possible values increment by some minimal step size, δ . For example, weight is continuous, but we can only measure it to some number of digits (and at best to the nearest discrete atom of mass). A random variable properly be considered a discrete random variable. For convenience, we treat it as continuous, which allows both the use of models based on probability and all the convenience and power of calculus.

The reason for this observation is simply that likelihood is derived from probability, not pdf's (see Lindsey 1999b). The likelihood is $g(x|\theta) \equiv \Pr\{\tilde{X} = x|\theta\}$, where \tilde{X} is the random variable. If \tilde{X} is not correct, we prefer to use the briefer form. However, the issues about the dimensional unit of $-\log(\mathcal{L}(\theta|x)) = -\log(g(x|\theta))$ are irrelevant. In particular, the dimensional unit of x is irrelevant to the likelihood of $\Pr\{\tilde{X} = x|\theta\}$, which is *always* probability. The probability $\Pr\{\tilde{X} = x$ is invariant to the measurement units used for x .

units being negative log-probability. But Claude Shannon, founder of information theory, established that $-\log(\text{probability})$ is a mathematical representation of information (and K-L is a model of f). The intrinsic unit for $-\log(\mathcal{L})$ is information. Units of AIC are always information.

There does not seem to be an established name for a unit of information. Claude Shannon died in February 2001. It has now been proposed that a unit of information be called the “Shannon.” Regardless of the name, here is that the units of the data are irrelevant when we compare models (which then become the units); and likelihood, properly, is a unit of data.

More clarification is needed. A continuous random variable X is a convenient way of dealing with the situation when the range of X is very small relative to the range of \tilde{X} , for example, a range of $\delta = 0.0001$. The range could be 0 to infinity, as long as the pdf of \tilde{X} drops off sufficiently fast for large x . If \tilde{X} is the pdf for a (conceptual) continuous random variable, δ is a small constant approximation $\Pr\{x|\theta\} \approx g(x|\theta)\delta$. For sufficiently small δ , this approximation is excellent and justifies using $g(x|\theta)\delta$ for $\mathcal{L}(\theta|x) = \Pr\{x|\theta\}$. δ drops out of all uses of the likelihood as a measure of relative likelihood of parameter values given the model, or about models given the data.

Thus, for models based on the presumption of continuous data, $-\log(\Pr\{x|\theta\}) = -\log(g(x|\theta)) - \log(\delta)$ is technically correct. It suffices to use $-\log(\mathcal{L}) = -\log(g(x|\theta))$ for purposes of comparing models. If computed, negative log-likelihood has information as its unit. However, $-\log(g(x|\theta))$ does not have this unit of dimensionality, missing $\log(\delta)$. This carries over to AIC. Because of such a constant possibly left out of the likelihood (i.e., dropped for simplicity), there is confusion over the inferential dimensional unit of AIC: it is information. It is justified because the only inferential way we use AIC is in terms of ΔAIC . The dimensional unit of Δ is information, and the additive constants common to, but left out of, each under consideration.

It is easy to lose sight of these deep matters when considering $-\log \mathcal{L}$ for normal distribution-based models. We should

$$-\log(\mathcal{L}(\hat{\theta}|x, g)) = \frac{n}{2} \log \left(\frac{\text{RSS}}{n} \right) + C(g)$$

where RSS is the residual sum of squares for the fitted model, and C depends only on the assumed normality of “errors,” ϵ . If all models considered assume normality and variance σ^2 , $C(g)$ is identical over all R models, and thus it drops out of ΔAIC .

When we bypass all or this important information about x , usually is, we risk confusion about what is the dimensional unit of AIC. For example, for the normal distribution one might expect the dimensional unit is $\log(\text{dimension of } x^2)$. In fact this is not the case because we are now dealing with units in $-\log\text{-probability}$.

6.9.6 AIC and Finite Mixture Models

Finite mixture models (see e.g., McLachlan and Peel 2000) are a class of statistical models, however, they are nonstandard in many ways. In particular, the usual likelihood ratio test statistic is not applicable because about parameter values under the null hypothesis being one parameter space (McLachlan and Peel 2000: 185–186). This has caused concerns about the use of AIC for model selection of mixture models are in the model set. People wonder, should the “2” in AIC be modified, just as the likelihood ratio test for mixtures must be modified. We considered the issue enough to provide a preliminary opinion. At a shallow level, AIC does not need to be modified for mixture models. However, at a deeper level, this is because AIC is not a test. At a deep level, when AIC to mixture models we encounter two issues: (1) the MLE of K is on a parameter space boundary, in light of failure of parameter estimability when the MLE is on a parameter space boundary, and (2) model redundancy occurs, and must be dealt with (redundant models are dropped from the model set). Thus, it is an aspect of how AIC should be modified.

Consider an example of a strictly positive response variable. A two-component mixture negative exponential model be useful (McLachlan 1988). The mixture model is

$$g_3(x) = \pi \left[\frac{1}{\lambda_1} e^{-x/\lambda_1} \right] + (1 - \pi) \left[\frac{1}{\lambda_2} e^{-x/\lambda_2} \right]$$

$0 < \pi < 1$, $0 < \lambda_i$, and $\lambda_1 \neq \lambda_2$. The model of (6.10) can be written

$$g_3(x) = \pi g_1(x) + (1 - \pi)g_2(x),$$

where

$$g_i(x) = \frac{1}{\lambda_i} e^{-x/\lambda_i}, \quad i = 1, 2.$$

The nominal number of parameters in model g_3 is $K = 3$, or g_1 , or g_2 , $K = 1$. The problem with fitting mixture models is that the MLE will be on a boundary if either $\hat{\pi} = 0$ (or 1), and then the model is no longer a mixture. In the first case, π is not estimable. In the second case π is not estimable; $\hat{\pi}$ can take

At a minimum, model selection here would have two choices equivalent to either g_1 or g_2 , i.e., no subscript is needed. The mixture model always has $K = 1$. When the likelihood based on \mathbf{y} is maximized with $\hat{\pi}$, $\hat{\lambda}_1$, and $\hat{\lambda}_2$ not on boundaries, only the mixture model is fitted and compute the usual AIC or AIC_c for the mixture model. If $\hat{\pi}$ is on a boundary, then the fitted mixture model actually collapses to a single model. In that case fitted models g_3 and g are redundant (Section 7.2). Model selection occurs because the only fitted model one has is g .

If model g_3 is not estimable, for the data at hand, it is as if one had an AIC for it as if $K = 3$. Such an AIC would actually be the log-likelihood identical to that for model g but would be for a model with $K = 3$. The AIC for model g would be the AIC for model g . However, model selection or multiple comparisons is meaningless in this situation where the mixture model in \mathbf{y} becomes \hat{g} .

The ideas here generalize if the mixing models are of different types. If the mixture model has more than two components. In the mixture model have $g_3 = \pi g_1 + (1 - \pi)g_2$, with models g_1 and g_2 being different (e.g., negative exponential and half-normal) with K_1 and K_2 components. If model g_3 will collapse to either model g_1 or g_2 if $\hat{\pi}$ is on a boundary, model redundancy occurs, and there really are only models g_1 and g_2 . The mixture model is estimable does one have three models? $K_3 = 1 + K_1 + K_2$ for model g_3 .

If the same type of basic model is used in a three component mixture, the possible complications increase. But the problem is the full mixture model cannot be fit because the MLE of π is on a boundary, then the model set must be adjusted for model redundancy. Also, the true nature of the reduced-component mixture must be recognized and its correct K determined. For example, a three component mixture, when fit to the data, might collapse to a 2-component mixture (it may have already been a model in the model set). Thus, model selection and AIC-model selection entails some traps that must be avoided.

The above does not address AIC and mixture models and model selection. One way to examine theory is to numerically compare the AIC, as properly used for mixture models, to the theoretical AIC for estimating (Section 7.2), e.g., $target = -2E_{\hat{\theta}}E_{\mathbf{y}}[\log(g_3(\mathbf{y}))]$ are with respect to the actual data-generating distribution. If the models were “on track” we did some Monte Carlo evaluations with n is (6.10) and data are generated under either (6.10) or (6.11). n were 25, 50, 100, and 200, with 25,000 or 50,000 Monte Carlo replications sufficed to fix $\lambda_1 = 1$ and vary λ_2 over 1, 5, 10, 15, 20. The results were over 0 to 0.5 by 0.1. The resulting comparisons, as $\delta = |E(AIC) - target|$ quite good for $n = 100$ and 200, and (to us) acceptable at $n = 25$ and 50.

the generating distribution was a mixture of a simple and a mixture (mean and maximum δ are over the set of values used for

generating distribution	n	δ	
		mean	maximum
simple	25	1.63	1.79
	50	1.10	1.31
	100	0.82	0.97
	200	0.55	0.90
mixture	25	0.66	1.32
	50	0.29	0.66
	100	0.24	0.53
	200	0.23	0.57

6.9.7 Unconditional Variance

The formula for estimating the unconditional variance of a parameter estimate is a derived result (Section 4.3.2, denoted (4.9)):

$$\widehat{\text{var}}\left(\hat{\theta}\right) = \left[\sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i)} + \left(\hat{\theta}_i - \hat{\theta}\right)^2 \right]$$

A corresponding formula is given in Section 4.3.2 for the covariance, $\widehat{\text{cov}}(\hat{\theta}, \hat{\tau})$. However, in obtaining (6.11) we assumed a wise conditional (on the models) correlation of estimator parameters across models. Thus, just on that ground there might be weak support for a theoretical exploration of an alternative to (6.11).

A better, but not overwhelming in practice, motivation for the underlying (6.11), as for example using linear models and p predictors x_1 to x_p . Thus, model g_i uses some subset of predictors and $\theta_i = E(y | \underline{x}, g_i)$ is the appropriate linear combination of parameters under model g_i . Because all the models considered have (Section 5.3.6)

$$\hat{\theta} = \tilde{\beta}_0 + x_1 \tilde{\beta}_1 + \cdots + x_p \tilde{\beta}_p = \underline{x}' \tilde{\underline{\beta}}$$

(see also Section 4.2.2 regarding $\tilde{\beta}_i$). Using (6.11) and the underlying formula for unconditional covariances we can obtain the variance–covariance matrix for $\tilde{\underline{\beta}}$, say $\hat{\underline{\Sigma}}$; $\hat{\underline{\Sigma}}$ is very complicated, but $\widehat{\text{var}}\left(\hat{\theta}_i | g_i\right)$, hence (6.11), is simple. An alternative for $\widehat{\text{var}}\left(\hat{\theta}_i | g_i\right)$ is, however, because of the nonlinearities involving the weights

One heuristic approach is to note what is used with DIC: the averaged posterior is a mixture distribution. Similarly, the pseudo-likelihood can be defined in the K-L framework of each model-specific likelihood. Either motivation lead to a possible alternative to (6.11) as

$$\widehat{\widehat{\text{var}}}(\hat{\theta}) = \sum_{i=1}^R w_i \left[\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2 \right]$$

This formula is linear in the weights. Using the Cauchy–Schwarz inequality can show that $\widehat{\widehat{\text{var}}}(\hat{\theta}) \leq \widehat{\widehat{\text{var}}}(\hat{\theta})$, hence (6.12) actually yields a lower bound. We emphasize that (6.12) is not a derived result in the K-L framework. We have not studied this matter further; however, the unconditional variance and covariance for AIC are subjects for more research.

6.9.8 A Baseline for $w_+(i)$

The measure $w_+(i)$ of the relative importance of variable x_i is defined in Section 4.2.2. We envision it as applied mostly when there are many variables (thus, $w_+(i)$ for $i = 1, \dots, p$), and the R model subsets of variables, such as just main-effect terms like x_i (in the R models), or these terms plus interaction-type terms, such as $x_i x_j$ (if p is small). One important point is that the interpretation of $w_+(i)$ within both the context of the set of models and predictors is different from the model set or the set of predictors, and $w_+(i)$ can change if the set of predictors changes. That $w_+(i)$ is not expected to be 0, even if x_i has no context effect, is at all. Rather, in this case $E(w_+(i)) > 0$, regardless of sample size. For this reason these summed weights give only a relative importance.

A randomization method can be used to estimate the $w_+(i)$ if x_i has no predictive value. We denote this unknown value by $w_{0+}(i)$. The data structure is an n by 1 response-variable column and the full design matrix X , which is n by $p + 1$ if the models are full (they usually do). Based on this data structure all R model AIC is obtained; Akaike weights are obtained, and $w_+(i)$ is computed.

To estimate $w_{0+}(i)$ there is one variation on this scenario. For each of the R models, permute the n values of x_i that are in column $i + 1$ of the design matrix, with the other columns unaltered. Then proceed in the usual way to compute what is nominally $w_+(i)$, except what you get is $w_{0+}(i)$. This computation is quite easy, especially if one already has the data. A slight alteration in code changes the bootstrap procedure to generate a sample in terms of the needed randomly permuted values of x_i .

i	variable	$w_+(i)$	$\hat{w}_{0+}(i)$
1	age	0.50	0.31
2	weight	0.93	0.31
3	height	0.31	0.29
4	neck	0.65	0.29
5	chest	0.28	0.29
6	abdomen	1.00	0.30
7	hips	0.45	0.31
8	thigh	0.59	0.31
9	knee	0.29	0.30
10	ankle	0.45	0.31
11	biceps	0.60	0.31
12	forearm	0.83	0.29
13	wrist	0.98	0.31

x_{i1}, \dots, x_{in} . Obviously, the random permutation renders y on average; the permuted x_i and all $x_j, j \neq i$, are also un-

We have tried this methodology; it seems it could be useful. One drawback. A single permutation sample is not enough; $\hat{w}_{0+}(i)$ varies from permutation sample to sample. We suggest doing a large number of permutations. More might be needed. Also the random variable $\hat{w}_{0+}(i)$ can be used to estimate the distribution. Hence, rather than average the resultant sample medians we suggest using the sample median as the single best $\hat{w}_{0+}(i)$.

An example of estimating baseline values for the $w_+(i)$ is given in Table 6.21. One hundred independent permutation samples were used to estimate $w_+(i)$ for each variable i , as described above, and the sample median was used. Because of the correlations among these predictors we thought they would vary distinctly, but they did not do so here. Perhaps the systematic variation in the set combined with the randomization process will always result in $w_{0+}(i)$ being the same for all x_i . From Table 6.21 we judge that age, weight, and in this context of usage the predictors height, chest circumference, and knee circumference have essentially zero importance. To estimate the importance of each variable we define a measure of absolute variable importance based on the difference $w_+(i) - \hat{w}_{0+}(i)$. Similar ideas appear in Breiman (2001). Combining these methods and ideas is possible and worthwhile.

Another informative idea for a baseline here is to look at the full set of $w_+(1), \dots, w_+(p)$ by leaving the X matrix intact and permuting the elements of vector y . Then fit all models, get the $w_+(i)$ values (the reader might think formal null hypothesis test here is not the intention). If each $w_+(i)$ and $\hat{w}_{0+}(i)$ pair are about the same, it indicates no predictability of y by the entire set of x_i . We have seen

e.g., $p = 50$ and $n = 100$.

For the body fat data example we also randomly permuted the y vector to estimate the $w_{0+}(i)$. This is a much faster approach. The median of the randomization sample was used as the estimate. For the body fat data, the results varied from 0.30 to 0.34 (mean of 0.32). For the other data sets, the results varied from 0.31 to 0.34 (mean of 0.32). For either x or y , we note that $w_{0+}(i)$ will depend upon p , and in general, the results are not too far from 0.31 or 0.32.

We conclude that more research on these ideas seems warranted.

6.10 Summary

This chapter is a more in-depth examination of some aspects of model selection; included are some comparisons to other methods. In particular, AIC is contrasted to BIC to better understand both methods. Those results are mostly in Sections 6.3 and 6.4. For their importance we will start this summary with references to them and compare AIC and BIC.

The derivation of BIC (Section 6.4.1) can be done without assuming that the set of models contains the true model. Thus, neither AIC (or Bayesian, in general) model selection methods require the true model to be in the set, nor do they require the true model to be valid, or use that the true data-generating model is in the set. However, for their consideration. Moreover, in the commonly assumed BIC framework (the true data-generating model and a fixed model set) as sample size goes to infinity, that selection converges with probability 1 to a single model, which is the true model, valid to infer that the selected model is truth (see e.g., Section 6.4.2).

Rather, the model selected by BIC converges to the model with the smallest dimension (i.e., the minimum K) in the subset (of size n) of models that all have the identical minimum Kullback–Leibler distance to the true model (6.4.2). Denote this model, which BIC selection converges to, as $g_b(\underline{x}|\underline{\theta}_o)$. The value of $\underline{\theta}$ that minimizes K-L information loss $I(f, g_b(\underline{x}|\underline{\theta}))$ of models. Model $g_b(\underline{x}|\underline{\theta}_o)$ is the model with both the smallest dimension K_i and for which $I(f, g_i(\cdot|\underline{\theta}_o))$ is minimized.

In reality with real data we expect the model that minimizes the K-L information loss to be unique in the model set. But in many, if not most, situations, the set of models used includes the data-generating model, and has it nested in some overly parametrized models. Thus, the set of models is being a subset of the R models that all have the same K-L information loss to the generating (“true”) model. BIC is then consistent for the true model, that generating model, which has the smallest dimension in the set. Thus the dimension-consistent property of BIC is motivated.

analysis the BIC target model cannot validly be inferred to be merely the model with minimum dimension that is nearest to f by K-L distance.

Whereas both the context and the target model for BIC are of sample size, the context for AIC is sample-size specific. The context is unchanging, under the information-theoretic approach the target model is allowed (in fact, assumed) to grow if n increases (by an order of magnitude). It is not realistic to let n go to infinity in the set of candidate models fixed (as in BIC) because such a limit means both more information and more factor levels, for any study. It then follows that more parameters need to be estimated. In this spirit the AIC target model g_a depends on n is the model for which $E_{\hat{\theta}}[I(f, g_i(\cdot|\hat{\theta}))]$ (expectation is over $\hat{\theta}$) is minimized over $i = 1, \dots, R$.

Thus the AIC target model (the K-L best model) is generally different from the BIC target model. Because both g_a and the model g_b are specific it is not logical to compare g_b to the model that g_a approaches as $n \rightarrow \infty$. Furthermore, the model that is g_a at huge sample size is different from the realized model g_a at actual n .

Section 6.4.5 provides a milestone result: It shows that AIC can be used as a Bayesian model selection criterion. The key is to use AIC with a quite different prior probability distribution on the model set. In the Bayesian context BIC assumes a uniform prior distribution on the model set. The K-L model prior that yields AIC is proportional to expected squared error. This can be generalized for QAIC_c. One consequence of this result is that for interpreting the Akaike weights as a posterior probability over the model set. Hence, w_i = the probability that model g_i is the best model.

Another result of Sections 6.3 and 6.4 is a clear understanding of the significance we must associate with prior and posterior probabilities. Saying “ p_i is the probability of model g_i ” we must be referring to the probability that this model is the target model of the model selection process. The target models, g_a and g_b noted above, are different for AIC and BIC (neither target is necessarily truth f). Now a Bayesian understanding of AIC and BIC is to argue for or against their respective priors. In general to understand the implications of those priors. For BIC $p_i = 1/R$. The prior for AIC makes p_i an increasing function of the number of parameters and a decreasing function of the number of estimable parameters. Such a prior corresponds to the idea that we expect the number of parameters we can reliably estimate to depend on n and K_i , and that the amount of information in the data.

The alternative way to understand and compare selected models in frequentist terms of their actual performance and expected

$n = 252$ ($K = 6, 191$), the AIC_c selected model has an AIC_c of 232. A 95% confidence set on the K-L best model includes 191 models. Thus it is totally unacceptable to say the best model has a 95% chance of being correct. In particular, it is absurd to interpret the included variables as important. This example illustrates the point that with high dimensionality, it typically be the case that even the best model will have very little support. The extensive body fat example of Section 6.2 has a similar problem. Section 6.2.8; we direct the reader there rather than repeat the details.

A few theoretical simulations were done based on the data in Section 6.2 to examine predictive mean square error (MSE) for BIC best-model selection strategy and for model averaging. The predictive mean square error (MSE) for AIC_c was 15% lower than BIC for the best model strategy and 17% lower under model averaging. Model averaging was superior to the traditional best-model selection strategy. AIC_c (MSE 15% lower) and BIC (23% lower). Ongoing research shows that model-averaged inferences are generally superior in all situations.

The remainder of Chapter 6 (Sections 6.5 to 6.9) has to do with overdispersion. Section 6.5 shows that overdispersion adjustment can be done with far more than one \hat{c} . The key idea is to partition the data into subsets and estimate overdispersion estimates by these data subsets. Partitioning can be by year, area, treatment, and so forth. When this is done, the number of different \hat{c} values used. Issues of goodness-of-fit with estimation of overdispersion factors. Therefore, a discussion of these issues is given in Section 6.5.1, including suggestions for no global model. We also note the matter of goodness-of-fit for a global model, assuming there is a global model and that it fits the data. This is, does the selected model then also fit the data? This is an ongoing research. Ongoing work (not included here) has shown that (global model fits) BIC can select a model that is in fact a good model but AIC virtually never does so. It is not clear it either, but operating characteristics should be of any concern.

Almost all model selection literature is only really applicable to parameters as fixed effects. However, the range of applications has expanded to random effects. It can also be expanded to other applications, such as generalized estimation equations. Section 6.6 discusses about these matters, especially application of AIC to models with multiple random effects. These extensions are in an early stage of development. What is important is that AIC can be extended beyond the traditional ML approaches of this book. An AIC-like Bayesian procedure can be applied in general to Bayesian hierarchical models.

Another seeming nonstandard application occurs when the data are based on the set based on different probability distributions. For example, we want to compare models for the data based on a gamma distribution.

the standard formula for AIC. That section gives some details and details.

In Section 6.8.3 we use numerical methods to evaluate sample bias of AIC when the data are from a gamma distribution; this result to the AIC bias-correction term used to get AIC_c (by assuming a normal distribution). The result is basically that AIC_c is an adequate small sample version of K-L model selection distribution for the situation studied ($n \geq 10$). This is in line with confirming the general usefulness of AIC_c even when the data are from a normal distribution.

This book focuses on the formal (i.e., objective) aspects of model selection and multimodel inference. Whereas we do not offer much advice on how to apply professional judgments in the course of data analysis, we do discuss the substantial role that judgment plays. We therefore recommend that statisticians try to be clear in their work about what supports the conclusions they draw from data. What is the quantifiable evidence; upon what assumptions do the assessment of evidence rely; are the persons making the judgments qualified to be making these judgments?

In Section 6.9.5 we consider the dimensional unit of AIC. Information is measured in the Shannon sense that negative $\log(\text{probability})$ measures information. Although the data have associated units, those units are lost, in a sense, as soon as one interprets the data by using a likelihood, which properly has units of probability. The units for $-\log(\mathcal{L})$ are information, regardless of the units of the data. It also follows that the dimensional interpretation of AIC is the loss when using model \hat{g}_i to approximate model \hat{g}_{min} (the best model).

Model selection has been applied to finite mixture models. Model selection for a mixture model is nonstandard and the likelihood ratio test must be modified for use with mixture models. This motivates a concern that the formula for AIC might need changing to apply to selection of finite mixture models. Our thought (Section 6.9.6) at this time is that the formula for AIC need not be modified for use with mixture models (heuristics are not a test, such as a LRT is). Instead, the big issue with finite mixture models is model redundancy that arises when a mixture model is fitted to data, i.e., the fitted mixture model actually collapses to a single component mixture model. In finite mixture models it is critical to properly deal with this model redundancy (Section 6.9.6).

We did some theoretical evaluation of whether AIC_c is a reasonable target for a two-component mixture negative exponential model. It did quite well at $n \geq 50$. Research on an improved small sample version of AIC_c is acceptable even for $n \geq 50$. Research on an improved small sample version of AIC_c is acceptable even for $n \geq 50$.

The relative importance of a predictor variable x_i can be given by $w_{+}(i)$. This number is relative, not absolute, because it will be given only when a predictor x_i has no predictive value at all in the given context. The value for $w_{+}(i)$ can be estimated by computer-intensive data partitioning. One can then better judge what predictors or factors are important in the data at hand. The method is simple, but computer intensive. See Section 6.9.8.

There is a lot of material in this chapter; some of it we have not covered here. Many sections probe issues about AIC without fully resolving them, thereby suggest additional research areas. A couple of highlights follow: The body fat example shows that one can expect large selection uncertainty with all subsets selection applied to the data, and illustrates ways to deal with this uncertainty. Those are some types of multimodel inference. A seminal result is that AIC is to the same extent as BIC and shows that the difference is all in the interpretation over the model set (i.e., model probabilities). Finally, how Bayesian model probabilities must be interpreted for model selection interpretations are different.

Statistical Theory and Numerical Results

This chapter contains theory and derivations relevant to K-L information-theory-based model selection. We have tried to make the chapters of this book readable by a general audience, especially students in various fields. Hence, we have reserved this chapter for material we believe should be made available to statisticians and biologists. For many, it will suffice to know that this theory encourages researchers, especially if they have some mathematical training, to read and try to understand the theory given. A deeper understanding provides a much deeper knowledge of model-based model selection in particular, and of some general model selection also.

The material given here is a combination of our distillation of the existing literature and what we feel are clarifications of the existing theory. In the former case we have not drawn from any one source; hence there is no particular reference for these derivations. We have not indicated what results in the literature about the estimation of expected K-L information this is sometimes not clear even to us.

7.1 Useful Preliminaries

The sole purpose of this section is to provide a summary of concepts, and mathematical background needed to produce

in that it establishes much of the notation and conventions of 7.2 and beyond in this chapter.

As a model selection criterion, it is clear what AIC is: for a model with K estimated parameters, $\hat{\theta}$ being the MLEs, computed from the data \underline{x} , under an assumed model. However, we need more detailed notation than just $\mathcal{L}(\hat{\theta})$, and we need to alternate between the likelihood and the pdf in the model. Therefore, without loss of generality we take the notation $\mathcal{L}(\underline{\theta} | \underline{x}) = g(\underline{x} | \underline{\theta})$ by simply then interpreting g as a function instead of using this convention we had constantly switched between $g(\underline{x} | \underline{\theta})$ and $\mathcal{L}(\underline{\theta} | \underline{x})$, that would be more confusing than the single notation $g(\underline{x} | \underline{\theta})$. This dual usage of the notation is noted; the reader must follow the mathematics with an eye being made at any point.

A second dual usage of notation for the random variable \underline{x} denotes the data (as a random variable), and sometimes \underline{x} denotes the space of integration, always with respect to $f(\underline{x})$, under an integration (in n -dimensional space). Because we are dealing with random variables, the expectation is usually denoted in terms of the statistical expectation operator, which is just an integral. At times we must have both a random variable and, separately, data, say \underline{y} . But the notation for data versus random variable is arbitrary and sometimes must be switched back and forth. It becomes impossible always to use \underline{x} for a variable of integration and data; hence, we do not try to do so, and instead we often use \underline{y} for data, even though at other times \underline{x} is an integrand variable and \underline{y} is a random variable; however, the data, no matter how denoted (\underline{x} or \underline{y} , or other), are from truth $f(\cdot)$, not from $g(\cdot | \underline{\theta})$ (when $f \neq g$); this is a point.

AIC has been motivated, justified, and derived in a variety of ways (for example, Akaike 1973, Sawa 1978, Sugiura 1978, Chow 1983, Shibata 1989, Bozdogan 1987), but these derivations are often difficult to follow. Here we give a general derivation in some detail, being rigorous about all required conditions (they are not listed, but we do note where approximations are made. The data have n dimensions and the general result is justified for “large” n . That is, the result is asymptotically as $n \rightarrow \infty$. Also, the integrals and expectations are over an n -dimensional sample space, although that fact is not fully reflected in the notation used.

The most general approach to deriving AIC uses the Taylor expansion to second order. An elementary introduction to the Taylor expansion is given by Peterson (1960) (or any introductory calculus book); a more advanced treatment including results for real-valued multivariable functions,

to $\underline{\theta}$ is given below:

$$h(\underline{\theta}) = h(\underline{\theta}_o) + \left[\frac{\partial h(\underline{\theta}_o)}{\partial \underline{\theta}} \right]' [\underline{\theta} - \underline{\theta}_o] + \frac{1}{2} [\underline{\theta} - \underline{\theta}_o]' \left[\frac{\partial^2 h(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right] [\underline{\theta} - \underline{\theta}_o] + Re$$

($\underline{\theta}$ and $\underline{\theta}_o$ are just two different points in the space over which $h(\cdot)$ is defined. Here, Re represents the exact remainder term for the quadratic expansion; the exact nature of Re is known (see Apostol, 1957). For heuristic approximations for the error that results from ignoring Re (see Appendix A), we can claim that an approximation of order

$$O(\|\underline{\theta} - \underline{\theta}_o\|^3).$$

Here, for any vector argument $\underline{z} - \underline{w}$,

$$\|\underline{z} - \underline{w}\| = \sqrt{\sum_{i=1}^K (z_i - w_i)^2}$$

denotes the Euclidean distance between the two points in K -dimensional space. Thus, the order of the approximation error is the cube of the distance between $\underline{\theta}$ and $\underline{\theta}_o$. This is quite a simplification that makes the point that the error of approximation is quite small if the distance is small.

The notation $O(x)$ denotes an unspecified (but positive) function of the scalar argument x that satisfies the condition that it is approximately equal to cx for small x , where c is a constant. It goes to 0 at least at a linear rate in x as x gets near 0. In the context of a quadratic approximation to $h(\underline{\theta})$ “near” $\underline{\theta}_o$ is arbitrarily close to $\underline{\theta}_o$ for $h(\cdot)$ a suitably smooth and bounded function.

In (7.1) the notation

$$\left[\frac{\partial h(\underline{\theta}_o)}{\partial \underline{\theta}} \right]$$

denotes a $K \times 1$ column vector of the first partial derivatives of $h(\cdot)$ with respect to $\theta_1, \dots, \theta_K$, evaluated at $\underline{\theta} = \underline{\theta}_o$; hence,

$$\left[\frac{\partial h(\underline{\theta}_o)}{\partial \underline{\theta}} \right] = \begin{bmatrix} \frac{\partial h(\underline{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial h(\underline{\theta})}{\partial \theta_K} \end{bmatrix}_{\underline{\theta} = \underline{\theta}_o}.$$

$$\begin{bmatrix} \partial^2 \theta^2 \\ \partial \theta_i \partial \theta_j \end{bmatrix}_{|\underline{\theta}=\underline{\theta}_o}$$

denotes the $K \times K$ matrix of second mixed partial derivatives with respect to $\theta_1, \dots, \theta_K$, evaluated at $\underline{\theta} = \underline{\theta}_o$. This matrix is the Hessian of $h(\underline{\theta})$.

The expansion in (7.1) when terminated at the quadratic approximation to $h(\underline{\theta})$. In this deterministic case, as indicated, the error of approximation is related roughly to the cube of the distance between $\underline{\theta}$ and $\underline{\theta}_o$. For a sufficiently small distance, this approximation is good. For the cases of interest, $h(\cdot)$ will be a log-likelihood based on a probability distribution. One special value of interest is the value needed in these expansions is the large-sample (hence asymptotic) value of the MLE $\hat{\underline{\theta}}$; that is, $E(\hat{\underline{\theta}}) \approx \underline{\theta}_o$ for large n (the exact relation to K-L information will be given below). The approximation error is often of order $1/n$, denoted by $O(1/n)$. This notation means that the approximation in $E(\hat{\underline{\theta}}) \approx \underline{\theta}_o$ is less than or equal to a constant divided by sample size for large sample sizes (the constant might even depend on n).

Stronger statements about large-sample limits are possible. As sample size $n \rightarrow \infty$, $\hat{\underline{\theta}} \rightarrow \underline{\theta}_o$ with probability 1, and the quadratic approximation given by (7.1) is quite good. In this case (7.1) becomes

$$h(\hat{\underline{\theta}}) = h(\underline{\theta}_o) + \left[\frac{\partial h(\underline{\theta}_o)}{\partial \underline{\theta}} \right]' [\hat{\underline{\theta}} - \underline{\theta}_o] + \frac{1}{2} [\hat{\underline{\theta}} - \underline{\theta}_o]' \left[\frac{\partial^2 h(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right] (\hat{\underline{\theta}} - \underline{\theta}_o)$$

Now the error of approximation in (7.3) is stochastic, but is generally on the order of $1/n$ with probability going to 1. The added “ p ” notation of the form $O_p(\cdot)$. The exact size of the error of approximation in expansions like (7.3) is not known (but is negligible for large sample sizes, subject to mild regularity conditions of the same type needed to ensure that the MLE is well behaved (Lehmann 1983)).

In the context of parametric MLE the standard approach is that the data are generated by one specific member of a family of models. The family of models, denoted here by $g(\underline{x} | \underline{\theta})$, is a set of probability distributions indexed by an unknown parameter that may be estimated by any valid method in parameter space Θ . By assumption, truth corresponds to one specific member of the family, of $\underline{\theta}$, which we could for clarity denote by $\underline{\theta}_o$. One would like to know where it comes from; it simply exists as (unknown) truth. Thus even if we know that the known model structure of g is true, there is still a fundamental question of an underlying unknown truth to the problem of inference. We cannot know, metaphysically, where this truth $\underline{\theta}_o$ comes from.

When we acknowledge that g is just a model of truth, and not the truth itself, specified, the issue arises as to what unique parameter in Θ corresponds to the truth.

given the concept of a fixed underlying unknown truth, and of this essential conceptualization of the inference problem that the data arose from some deep truth, denoted without f . Now, one cannot usefully ask where truth f comes from in physical sense that one cannot ask where $\underline{\theta}_o$ comes from under that $g(\underline{x} | \underline{\theta}_o)$ is truth, but we just do not happen to know it.

Given this essential framework of f as truth (rather than as truth) we can, and must, ask whether there is a unique model in the class of models $g(\cdot | \underline{\theta})$ that best describes the data. Here, given a class of models $g(\cdot | \underline{\theta})$, $\underline{\theta} \in \Theta$, is there a unique $\underline{\theta}_o$ that the MLE is this $g(\cdot | \underline{\theta}_o)$ a best model in some sense? In fact, the samples) estimating a unique parameter value that we will call this parameter value that indexes our target model under $\underline{\theta}_o$ (we will say more on this below).

Approached theoretically, ignoring issues of data and estimation, approximating model g in the class of models considered, under the K-L information measure, is simply the model that produces the minimum K-L value over Θ . Hence we look for a unique value of $\underline{\theta}$ that we denote by $\underline{\theta}_o$, that provides the K-L best approximating model. This is the solution to the optimization problem

$$\min_{\underline{\theta} \in \Theta} [I(f, g)] = \int f(\underline{x}) \log \left(\frac{f(\underline{x})}{g(\underline{x} | \underline{\theta}_o)} \right) d\underline{x}$$

Clearly, $g(\underline{x} | \underline{\theta}_o)$ is the best model here, and this serves, in fact, as a target $\underline{\theta}_o$ given f and given the class of models $g(\cdot | \underline{\theta})$. Below, the MLE of $\underline{\theta}$ under model g is estimating $\underline{\theta}_o$.

Given the assumed regularity conditions on the model, the following equations

$$\frac{\partial}{\partial \underline{\theta}} \int f(\underline{x}) \log \left(\frac{f(\underline{x})}{g(\underline{x} | \underline{\theta}_o)} \right) d\underline{x} = 0.$$

Rewriting (7.4) using that $\log(a/b) = \log(a) - \log(b)$, we have

$$\frac{\partial}{\partial \underline{\theta}} \int f(\underline{x}) \log(f(\underline{x})) d\underline{x} - \frac{\partial}{\partial \underline{\theta}} \int f(\underline{x}) \log(g(\underline{x} | \underline{\theta})) d\underline{x} = 0.$$

Because $\underline{\theta}$ is not involved in $f(\cdot)$, the first term of the above equation (ignoring the minus sign) can be written as

$$\int f(\underline{x}) \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]_{|\underline{\theta}=\underline{\theta}_o} d\underline{x} = E_f \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]_{|\underline{\theta}=\underline{\theta}_o} \right]$$

$$E_f \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right] = \underline{0}.$$

The well-known asymptotic consistency property of convergence of means of *iid* random variables allow, (7.5), another interpretation of $\underline{\theta}_o$. If \underline{x} represents an n from pdf $f(\underline{x}) \equiv \prod_{i=1}^n f(x_i)$ and we consider the $g(\underline{x} | \underline{\theta}) \equiv \prod_{i=1}^n g(x_i | \underline{\theta}_o)$, then for every n we have the K (expressed as a mean, without loss of generality)

$$\frac{1}{n} \left[\sum_{i=1}^n \frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \hat{\underline{\theta}})) \right] = \underline{0}.$$

As $n \rightarrow \infty$ two limits are approached with probability 1 (almost sure convergence). The sequence of MLEs $\hat{\underline{\theta}}(n)$ (adding no new MLE as a function of sample size), converges to some limit. This limit has to converge to $\underline{\theta}_o$, because the means on the left-hand side of the likelihood equations converge (as n gets large) to their expected values, which must all equal 0. Under suitable regularity conditions, (7.5) implies that the unique value of $\underline{\theta} = \underline{\theta}_o$. Hence, the sequence $\hat{\underline{\theta}}(n)$ converges almost surely to $\underline{\theta}_o$, which is the K-L minimizer (see, e.g., White, 1984).

Some deep ideas and philosophy are involved in the above discussion. In particular, we have the distinction that unknown truth $f(\underline{x})$ implies that the numerical values on the (often only) conceptual, but the parameter $\underline{\theta}$ of interest to us. Yet $f(\underline{x})$ is not a mathematical function. The model $g(\underline{x} | \underline{\theta})$ is a mathematical function of $\underline{\theta}$, because the parameter space is interpretable, hence useful to consider, and varies over the parameter space. Even if we think that $g(\cdot)$ represents truth, this is still a single point $\underline{\theta}_o$ in the parameter space (in a frequentist perspective). Hence, in this case we would be saying that $f(\underline{x}) \equiv g(\underline{x} | \underline{\theta}_o)$ is a single fixed point; thus even in this context $\underline{\theta}$ is not a function. Therefore, in this or any case, $f(\underline{x})$ is not a function of $\underline{\theta}$.

$$\frac{\partial}{\partial \underline{\theta}} \int f(\underline{x}) \log(f(\underline{x})) d\underline{x} = \underline{0}.$$

The derivation of AIC occurs in the context of probability distributions and expectations of functions of random variables. Such a derivation involves various types of integrals, but the notation and “machinery” of stochastic calculus are more convenient to use here than the explicit notation of integration. A particular aspect of this matter that needs to be noted is that the order of taking two expectations of the form $E_{\underline{x}} E_{\underline{y}} [h(\underline{x}, \underline{y})]$ is not important. \underline{x} and \underline{y} denote random variables. The function $h(\cdot, \cdot)$ is a scalar. The calculus of integrals as linear operators, $E_{\underline{x}} E_{\underline{y}} [h(\underline{x}, \underline{y})] = E_{\underline{y}} E_{\underline{x}} [h(\underline{x}, \underline{y})]$.

Another aspect of preliminaries concerns (7.2). If $h(\cdot)$ is $\log(g(\underline{x} | \underline{\theta}))$, then (7.2) is

$$\left\{ \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \theta_i \partial \theta_j} \right\}_{|\underline{\theta}=\underline{\theta}_o},$$

which is related to the Fisher information matrix

$$\mathcal{I}(\underline{\theta}_o) = E_g \left\{ - \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \theta_i \partial \theta_j} \right\}_{|\underline{\theta}=\underline{\theta}_o}$$

(expectation here is with respect to $g(\cdot)$). If $g(\cdot)$ is the true model (it is if f is a special case of g , or if $g = f$), then the covariance matrix Σ of the MLE is (for large samples) $\Sigma \approx [\mathcal{I}(\underline{\theta}_o)]^{-1}$. If g is not the true model, $\Sigma = E(\hat{\underline{\theta}} - \underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)'$ is $[\mathcal{I}(\underline{\theta}_o)]^{-1}$. If g is not the true model, Σ may be less general than the true model, or otherwise different. In general we must expect that $\Sigma \neq [\mathcal{I}(\underline{\theta}_o)]^{-1}$. In fact, we should take expectations with respect to f , not g . Hence, we define

$$I(\underline{\theta}_o) = E_f \left\{ - \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \theta_i \partial \theta_j} \right\}_{|\underline{\theta}=\underline{\theta}_o}$$

In the case that $f = g$ or f is a special case of g , then and only then $\mathcal{I}(\underline{\theta}_o) = I(\underline{\theta}_o)$. We will not generally make this distinction between $\mathcal{I}(\underline{\theta}_o)$ and $I(\underline{\theta}_o)$ to whether the situation allows $\mathcal{I}(\underline{\theta}_o) = I(\underline{\theta}_o)$ or not. It is important, however, to be always cognizant of whether the expectation in $\mathcal{I}(\underline{\theta}_o)$ is with respect to f or g .

Additional notation useful here is the empirical, but un-

$$\hat{I}(\underline{\theta}_o) = \left\{ - \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \theta_i \partial \theta_j} \right\}_{|\underline{\theta}=\underline{\theta}_o}.$$

For simpler notation we will use

$$I(\underline{\theta}_o) = E_f \left[- \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}_o))}{\partial \theta^2} \right],$$

which means exactly the same as (7.7), and hence simply $\hat{I}(\underline{\theta}_o)$ is

$$\hat{I}(\underline{\theta}_o) = - \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}_o))}{\partial \theta^2}.$$

It is obvious that $E_f[\hat{I}(\underline{\theta}_o)] = I(\underline{\theta}_o)$. When \underline{x} is a random vector, $\hat{I}(\underline{\theta}_o)$ converges to $I(\underline{\theta}_o)$ as $n \rightarrow \infty$. We can express this as

$$\hat{I}(\underline{\theta}_o) = I(\underline{\theta}_o) + Re, \text{ and usually } Re \text{ is } O(n^{-1/2}).$$

$$\hat{I}(\hat{\theta}) = -\frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \underline{\theta}^2}.$$

Because $\hat{\theta}$ is the MLE under the model $g(\underline{x} | \underline{\theta})$, $\hat{\theta}$ converges and hence $\hat{I}(\hat{\theta})$ converges to $I(\underline{\theta}_o)$. Thus, $\hat{I}(\hat{\theta}) \approx I(\underline{\theta}_o)$ approximation is at worst $O(1/\sqrt{n})$, and in most common cases be $O(1/n)$. If we could determine the analytical form (using the data) of $I(\underline{\theta}_o)$, an alternative estimator would be $I(\hat{\theta})$, at the MLE; $I(\hat{\theta})$ is often not the same as $\hat{I}(\hat{\theta})$. Note also that the used estimator $\mathcal{I}(\hat{\theta})$ (i.e., (7.6) evaluated at the MLE) is not the same as either $I(\hat{\theta})$ or $\hat{I}(\hat{\theta})$ and may not converge to $I(\underline{\theta}_o)$.

There are two ways to compute the Fisher information matrix when $f = g$. This additional material, and more, is needed and is given here. Because the model is a probability distribution

$$\int g(\underline{x} | \underline{\theta}) d\underline{x} = 1,$$

and therefore (under the same mild regularity conditions

$$\int \frac{\partial g(\underline{x} | \underline{\theta})}{\partial \underline{\theta}} d\underline{x} = \underline{0}.$$

Next, we use in the above the result

$$\frac{\partial \log(g(\underline{x} | \underline{\theta}))}{\partial \underline{\theta}} = \frac{1}{g(\underline{x} | \underline{\theta})} \left[\frac{\partial g(\underline{x} | \underline{\theta})}{\partial \underline{\theta}} \right]$$

and hence we get

$$\int g(\underline{x} | \underline{\theta}) \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] d\underline{x} = \underline{0}.$$

Now take the partial derivative vector of (7.9) with respect to $\underline{\theta}$; this derivation uses the chain rule of differentiation and some algebraic results:

$$\begin{aligned} \int g(\underline{x} | \underline{\theta}) \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]' d\underline{x} \\ + \int g(\underline{x} | \underline{\theta}) \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \underline{\theta}^2} d\underline{x} \end{aligned}$$

(O is a $K \times K$ matrix of zero elements). We can rewrite

$$\mathbb{E}_g \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]' \right] = \mathbb{E}_g \left[- \right]$$

$$E_g \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]^T \right]$$

We will denote the left-hand side of the above by $\mathcal{J}(\underline{\theta})$; hence

$$\mathcal{J}(\underline{\theta}) = E_g \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]^T \right]$$

Thus, $\mathcal{I}(\underline{\theta}) = \mathcal{J}(\underline{\theta})$, but the expectations underlying $\mathcal{I}(\underline{\theta})$ are with respect to $g(\underline{x} | \underline{\theta})$, not with respect to $f(\underline{x})$. One consequence is that the inverse Fisher information matrix may not be the correct conditional variance–covariance matrix of the MLE if the model is misspecified.

What we need more than (7.11) is

$$J(\underline{\theta}) = E_f \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta})) \right]^T \right]$$

We can expect $J(\underline{\theta}) = \mathcal{J}(\underline{\theta})$ only when $f = g$, or f is truth. Although $\mathcal{I}(\underline{\theta}) = \mathcal{J}(\underline{\theta})$, there is no such general equality for $J(\underline{\theta})$ when g is only an approximation to f , hence when $I(f, g) > 0$. Heuristically, however, the inequalities of the sort $I(\underline{\theta}_o) \approx J(\underline{\theta}_o)$, $\mathcal{I}(\underline{\theta}_o) \approx I(\underline{\theta}_o)$, and $J(\underline{\theta}_o) \approx I(f, g) \approx 0$, hence when a good approximating model is used.

There is a large-sample relationship among $I(\underline{\theta}_o)$, $J(\underline{\theta}_o)$, and Σ knowing, and perhaps should be used more:

$$I(\underline{\theta}_o)\Sigma = J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1},$$

and hence

$$\Sigma = [I(\underline{\theta}_o)]^{-1} J(\underline{\theta}_o) [I(\underline{\theta}_o)]^{-1},$$

where Σ is the true large-sample variance–covariance matrix derived from model g when f is truth. It suffices to derive (7.13) that we will use more directly in deriving AIC.

Expanding the likelihood equations evaluated at $\underline{\theta}_o$ as a series about the MLE, we have

$$\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \approx \frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \hat{\underline{\theta}})) + \left[\frac{\partial^2 \log(g(\underline{x} | \hat{\underline{\theta}}))}{\partial \underline{\theta}^2} \right] (\underline{\theta}_o - \hat{\underline{\theta}})$$

The MLE satisfies

$$\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \hat{\underline{\theta}})) = \underline{0};$$

$$\frac{\partial \underline{\theta}}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \sim \left[-\frac{\partial^2}{\partial \underline{\theta}^2} \right] (\underline{\theta} - \underline{\theta}_o) \\ = \hat{I}(\hat{\underline{\theta}})(\hat{\underline{\theta}} - \underline{\theta}_o) \approx I(\underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)$$

From the above we get

$$[I(\underline{\theta}_o)]^{-1} \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right] \approx (\hat{\underline{\theta}} - \underline{\theta}_o)$$

Transpose (7.15) and use that transposed result along with (7.14) to derive

$$[I(\underline{\theta}_o)]^{-1} \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right]' [I(\underline{\theta}_o)]^{-1}$$

Now take the expectation of the above with respect to $f(\underline{x})$ to get

$$[I(\underline{\theta}_o)]^{-1} J(\underline{\theta}_o) [I(\underline{\theta}_o)]^{-1} \approx E_f(\hat{\underline{\theta}} - \underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)'$$

hence, we have (7.14) as a large-sample result.

The above likelihood-based results under either a true data-generating process or under model misspecification (i.e., truth is f , the model is g) are standard in the statistical literature. For very rigorous derivations see, for example, [10].

To take expectations of the quadratic forms that are in (7.14) and (7.15) we will need to use an equivalent expression of that form:

$$\underline{z}' A \underline{z} = \text{tr} [A \underline{z} \underline{z}'] .$$

Here “tr” stands for the matrix trace function, the sum of the diagonal elements of a square matrix. The trace function is a linear operator: $\text{tr}(A+B) = \text{tr}(A) + \text{tr}(B)$. If \underline{z} is a quadratic form in a stochastic variable in \underline{z} , its expectation can be written as

$$E_{\underline{z}} [\underline{z}' A \underline{z}] = \text{tr} [E_{\underline{z}} [A \underline{z} \underline{z}']] .$$

If A is fixed (or stochastic but independent of \underline{z}), then $E_{\underline{z}} [A \underline{z} \underline{z}'] = A E_{\underline{z}} [\underline{z} \underline{z}']$. If \underline{z} has mean $\underline{0}$ (such as $\underline{z} = \hat{\underline{\theta}} - E(\hat{\underline{\theta}})$), then $E_{\underline{z}} [\underline{z} \underline{z}'] = \Sigma$, the covariance matrix of \underline{z} ; hence then

$$E_{\underline{z}} [\underline{z}' A \underline{z}] = \text{tr} [A \Sigma] .$$

If A is stochastic but independent of \underline{z} , then we can use

$$E_A E_{\underline{z}} [\underline{z}' A \underline{z}] = \text{tr} [E_A E_{\underline{z}} [A \underline{z} \underline{z}']] = \text{tr} [E_A (A \Sigma)]$$

A final aspect of notation, and of concepts, reemphasized at the start of this section: The notation for a random variable is arbitrary in taking expectations over the sample space. What matters for such an expectation (i.e., integration) is the model under consideration, its form, its assumptions, its parameters, and the distribution of the random variable.

data analysis framework. Because an expectation is a type of average over a well-defined space, the result of the integration is not dependent on the space used in the integrand. Thus

$$\begin{aligned} E_f \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] &= \int f(\underline{x}) \log(g(\underline{x} | \hat{\theta}(\underline{x}))) d\underline{x} \\ &\equiv \int f(\underline{y}) \log(g(\underline{y} | \hat{\theta}(\underline{y}))) d\underline{y} \\ &= E_f \left[\log(g(\underline{y} | \hat{\theta}(\underline{y}))) \right] \end{aligned}$$

Changing notation for the integrand (i.e., \underline{x} to \underline{y}) has no effect on the result. This type of useful notation change is required in derivations of AIC. In places, at a conceptual level, we recognize two independent derivations. We have two notations, \underline{x} and \underline{y} . In fact, these derivations are based on the frequentist properties of data-analysis methods, but there is no need for being used in these derivations. Rather, in these theoretical derivations, possible “data” are just points in an n -dimensional sample space, in accordance with some true probability distribution $f(\cdot)$.

7.2 A General Derivation of AIC

We now give a general conceptual and then mathematical derivation of AIC, starting from K-L information for the best approximation of models $g(\underline{x} | \underline{\theta})$:

$$I(f, g(\cdot | \underline{\theta}_o)) = \int f(\underline{x}) \log \left(\frac{f(\underline{x})}{g(\underline{x} | \underline{\theta}_o)} \right) d\underline{x}$$

Note that while for the model we do not know $\underline{\theta}$, the target value for the class of models is appropriately taken as $I(f, g(\cdot | \underline{\theta}_o))$ (i.e., 7.16), because the parameter value we will be estimating is based on the expanded notation in (7.16), so we can represent $I(f, g(\cdot | \underline{\theta}_o))$ in general, on the unknown parameter value, given the model $g(\cdot | \underline{\theta})$. $I(f, g)$ does not involve any data, nor any value of \underline{x} , since it is a constant value out.

Given that we have data \underline{y} as a sample from $f(\cdot)$, the likelihood function to find the MLE $\hat{\theta} = \hat{\theta}(\underline{y})$ and compute an estimate of $I(f, g(\cdot | \hat{\theta}(\underline{y})))$ is

$$I(f, g(\cdot | \hat{\theta}(\underline{y}))) = \int f(\underline{x}) \log \left(\frac{f(\underline{x})}{g(\underline{x} | \hat{\theta}(\underline{y}))} \right) d\underline{x}$$

approaches, both will lead to the same basis for AIC (the difference from K-L to AIC).

If we could find the $\underline{\theta}_o$ that minimizes K-L (for a given f) that our target for a perfect model would be $I(f, g) = 0$. How good any model is relative to this absolute value of K-L does not change when we have only an estimate of $\underline{\theta}$. Even if our estimate (miraculously) truth, hence $g(\underline{x} | \underline{\theta}_o) = f(\underline{x})$, our estimates are equal $\underline{\theta}_o$ almost surely for continuous parameters and distributions. For some discrete distributions the equality would be with probability 1. Any value of $\hat{\underline{\theta}}(\underline{y})$ other than $\underline{\theta}_o$ results in $I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) > 0$. Thus, even if we had the correct model structure, because of estimation error $\underline{\theta}$ we should think in terms of the (essentially estimated) K-L information, on average, a value > 0 . This motivates us to revise our idea of perfect agreement to be as a measure of perfect agreement of fitted model to the true model.

In the context of repeated sampling properties as a criterion, we would expect our estimated K-L to have on average the same value as the true K-L. We should therefore readjust our idea of perfect agreement to be the minimizing of $I(f, g(\cdot | \underline{\theta}_o))$ (given f) to a larger value, on average, given by

$$E_{\underline{y}} \left[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) \right] > I(f, g(\cdot | \underline{\theta}_o))$$

(and repeating ourselves because it is an important point: A and B are with respect to f regardless of the notation for random variables such as \underline{x} , \underline{y} , or $\hat{\underline{\theta}}$). Thus, given the reality that we must adopt the criterion

“select the model g to minimize $E_{\underline{y}} \left[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) \right]$ ”

Hence our goal must be to minimize the expected value of the estimated K-L information value. (If we could compute the K-L for each model, we could stay with the goal of minimizing K-L. Curious we note here that the large-sample difference is $E_{\underline{y}} \left[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) \right] - I(f, g(\cdot | \underline{\theta}_o)) = \frac{1}{2} \text{tr} \left[J(\underline{\theta}_o) I(\underline{\theta}_o)^{-1} \right]$, which does not depend on n .)

Rewriting the basis of this new target to be minimized,

$$E_{\underline{y}} \left[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) \right] = \int f(\underline{x}) \log(f(\underline{x})) d\underline{x} - E_{\underline{y}} \left[\int f(\underline{x}) \log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))) d\underline{x} \right]$$

hence

$$E_{\underline{y}} \left[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) \right] = \text{constant} - E_{\underline{y}} E_{\underline{x}} \left[\log[g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))] \right]$$

value given by (7.18). In most of our writing here about T , it is much simpler just to say that we are selecting an estimated model by use of AIC.

There is a second, less compelling, approach that we call K-L to AIC: Start with

$$I(f, g(\cdot | \underline{\theta}_o)) = \text{constant} - E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))]$$

and see whether we can compute (or estimate) $E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))]$ by Taylor series expansions. As will be made evident below, the result

$$\begin{aligned} E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))] &\approx E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \frac{1}{2} \text{tr} \\ &\quad - \frac{1}{2} (\hat{\underline{\theta}}(\underline{y}) - \underline{\theta}_o)' I(\underline{\theta}_o) (\hat{\underline{\theta}}(\underline{y}) - \underline{\theta}_o) \end{aligned}$$

On the right-hand side above, the only component that is estimated or computed (in any useful way) is the quadratic form $(\hat{\underline{\theta}}(\underline{y}) - \underline{\theta}_o)$ (and it is pointless therein to use $\hat{\underline{\theta}}_o = \hat{\underline{\theta}}(\underline{y})$). Taking the expectation of both sides above with respect to \underline{y} , we get an estimate:

$$E_{\underline{y}} E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))] \approx E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \text{tr} \left[\frac{1}{2} I(\underline{\theta}_o) \right]$$

Thus, either line of derivation demonstrates that we can avoid the subjective from model selection based on minimum K-L value, by selecting the model with estimated $\underline{\theta}$ based on minimum K-L information measure. It is still the case that only a relative value can be found based on $E_{\underline{y}} E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))]$ as the target objective function to be maximized; the constant $E_{\underline{x}} [f(\underline{x}) \log(f(\underline{x}))]$ cannot be computed.

Only some of the literature is clear that AIC model selection is based on the concept of minimizing the expected K-L criterion $E_{\underline{y}} E_{\underline{x}} [f(\underline{x}) \log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))]$ (see, e.g., Sawa 1978, Sugiura 1978, Bozdogan 1987 (page 10), Akaike 1973, Milhaud 1994). It is the relative value of this criterion that is to be compared across a set of models. That is, we want to estimate without bias, as closely as possible, the criterion (denoted below by T for target) for each approach. The value of

$$T = \int f(\underline{y}) \left[\int f(\underline{x}) \log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))) d\underline{x} \right] d\underline{y}$$

The change from conceptual model selection based on minimizing the expected K-L criterion to model selection based on maximizing an estimate of T is

Sometimes the criterion given by (7.17), hence AIC, is interpreted as a concept of Akaike's predictive likelihood $E_p[\log(\mathcal{L}(\hat{\theta}))] = E_p[\log(g(\underline{x} | \hat{\theta}))] \equiv T$, which has a heuristic interpretation in terms of the expected log-likelihood of independent random variables \underline{x} and \underline{y} . However, the criterion by maximizing \hat{T} (or minimizing $-2\hat{T}$), does arise from an approach to the problem of model selection without ever using cross-validation.

In a slightly simplified, but obvious, notation, the K-L-based model selection problem is now to find a useful expression for, and estimate,

$$T = E_{\hat{\theta}} E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\theta})) \right],$$

where it is understood that the MLE $\hat{\theta}$ is based on sample means, so that the expectations are for \underline{x} and \underline{y} (hence $\hat{\theta}$) both with respect to the same samples because T is also a double expectation based, conceptually, on the same samples that AIC-based model selection is asymptotically equivalent to cross-validation (see, e.g., Stone 1977); cross-validation is a well-known model selection method.

Step 1 is an expansion of the form (7.3) applied to $\log(g(\underline{x} | \hat{\theta}))$ for any given \underline{x} :

$$\begin{aligned} \log(g(\underline{x} | \hat{\theta})) &\approx \log(g(\underline{x} | \underline{\theta}_o)) + \left[\frac{\partial \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}} \right]_{\underline{\theta} = \underline{\theta}_o} [\hat{\theta} - \underline{\theta}_o] \\ &\quad + \frac{1}{2} [\hat{\theta} - \underline{\theta}_o]' \left[\frac{\partial^2 \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}^2} \right]_{\underline{\theta} = \underline{\theta}_o} [\hat{\theta} - \underline{\theta}_o]. \end{aligned}$$

Truncation at the quadratic term entails an unknown degree of error (but it is an error of approximation that goes to zero as the sample size increases). In going from (7.21) to (7.20) we first take the expected value of (7.21)

$$\begin{aligned} E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\theta})) \right] &\approx E_{\underline{x}} \left[\log(g(\underline{x} | \underline{\theta}_o)) \right] + E_{\underline{x}} \left[\frac{\partial \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}} \right]_{\underline{\theta} = \underline{\theta}_o} [\hat{\theta} - \underline{\theta}_o] \\ &\quad + \frac{1}{2} [\hat{\theta} - \underline{\theta}_o]' \left[E_{\underline{x}} \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}^2} \right]_{\underline{\theta} = \underline{\theta}_o} [\hat{\theta} - \underline{\theta}_o]. \end{aligned}$$

The vector multiplier of $[\hat{\theta} - \underline{\theta}_o]$ in the linear term above is $E_{\underline{x}} \left[\frac{\partial \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}} \right]$ as (7.5). It is just that for clarification $E_{\underline{x}}$ is used to mean the expectation over the random variable \underline{x} (and keep remembering that $\hat{\theta} \equiv \hat{\theta}(\underline{x})$ of \underline{x}). Therefore, upon taking this expectation, the linear term in (7.5) applies:

$$E_{\underline{x}} \left[\frac{\partial \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}} \right] = \underline{0}.$$

$$E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\underline{\theta}})) \right] \approx E_{\underline{x}} \left[\log(g(\underline{x} | \underline{\theta}_o)) \right] - \frac{1}{2} [\hat{\underline{\theta}} - \underline{\theta}_o]' I(\underline{\theta}_o)$$

Now we can take the expectation of (7.23) with respect to $\underline{\theta}$ where the trace function is used, yielding

$$E_{\hat{\underline{\theta}}} E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\underline{\theta}})) \right] \approx E_{\underline{x}} \left[\log(g(\underline{x} | \underline{\theta}_o)) \right] - \frac{1}{2} \text{tr} \left[I(\underline{\theta}_o) E_{\hat{\underline{\theta}}} \right]$$

The left-hand side above is T from (7.20), and $E_{\hat{\underline{\theta}}} \left[[\hat{\underline{\theta}} - \underline{\theta}_o] [\hat{\underline{\theta}} - \underline{\theta}_o]' \right]$ is the correct large-sample theoretical sampling variance of $\hat{\underline{\theta}}$. The expectation herein is taken with respect to truth f , not \hat{f} . Thus we have

$$T \approx E_{\underline{x}} \left[\log(g(\underline{x} | \underline{\theta}_o)) \right] - \frac{1}{2} \text{tr} \left[I(\underline{\theta}_o) \Sigma \right]$$

Step 2 starts with the realization that we have not yet derived a relationship between T and $E_{\underline{x}} \left[\log[g(\underline{x} | \hat{\underline{\theta}}(\underline{x}))] \right]$, which is the log-likelihood at the MLE. We now do a second-order Taylor expansion of $\log(g(\underline{x} | \underline{\theta}_o))$ about $\hat{\underline{\theta}}(\underline{x})$, treating \underline{x} as the sample point and $\hat{\underline{\theta}}(\underline{x})$ as the MLE of $\underline{\theta}$ for this \underline{x} . This procedure is acceptable, because $\hat{\underline{\theta}}(\underline{x})$ is an expected value, which means taking an integral over the sample space. Therefore, it does not matter what sample point we use: \underline{x} or \underline{y} . Applying the Taylor series (but with the roles of $\hat{\underline{\theta}}$ and $\underline{\theta}_o$ switched; also note well that $\hat{\underline{\theta}}$ is a function of \underline{x}) we obtain

$$\begin{aligned} \log(g(\underline{x} | \underline{\theta}_o)) &\approx \log(g(\underline{x} | \hat{\underline{\theta}})) + \left[\frac{\partial \log(g(\underline{x} | \hat{\underline{\theta}}))}{\partial \underline{\theta}} \right]_{\underline{\theta} = \hat{\underline{\theta}}} (\underline{\theta}_o - \hat{\underline{\theta}}) \\ &\quad + \frac{1}{2} (\underline{\theta}_o - \hat{\underline{\theta}})' \left[\frac{\partial^2 \log(g(\underline{x} | \hat{\underline{\theta}}))}{\partial^2 \underline{\theta}} \right]_{\underline{\theta} = \hat{\underline{\theta}}} (\underline{\theta}_o - \hat{\underline{\theta}}) \end{aligned}$$

The MLE $\hat{\underline{\theta}}$ is the solution of, hence satisfies, the equation

$$\frac{\partial \log(g(\underline{x} | \hat{\underline{\theta}}))}{\partial \underline{\theta}} = 0.$$

Therefore, the linear term in (7.25) vanishes. Taking the expectation over $\underline{\theta}$ we can write

$$E_{\underline{x}} \left[\log(g(\underline{x} | \underline{\theta}_o)) \right] \approx E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\underline{\theta}})) \right] - \frac{1}{2} \text{tr} \left[E_{\underline{x}} \left[\hat{I}(\hat{\underline{\theta}}) \right] \right]$$

See (7.8) for $\hat{I}(\hat{\underline{\theta}})$, the Hessian of the log-likelihood evaluated at $\hat{\underline{\theta}}$.

$$\begin{aligned}
E_{\underline{x}} \left[\hat{I}(\hat{\theta}) \right] [\underline{\theta}_o - \hat{\underline{\theta}}][\underline{\theta}_o - \hat{\underline{\theta}}]' &\approx [I(\underline{\theta}_o)] \left[E_{\underline{x}}[\underline{\theta}_o - \hat{\underline{\theta}}] \right] \\
&= [I(\underline{\theta}_o)] \left[E_{\underline{x}}[\hat{\underline{\theta}} - \underline{\theta}_o] \right] \\
&= [I(\underline{\theta}_o)] \Sigma.
\end{aligned}$$

The approximation made in (7.27) is often good to $O(1/n)$. However, there are circumstances where the approximation is not good, and the overall approximation in (7.27) is equivalent to $I(\underline{\theta}_o)$ after first writing the approximation

$$\begin{aligned}
E_{\underline{x}} \left[\hat{I}(\hat{\theta}) \right] [\underline{\theta}_o - \hat{\underline{\theta}}][\underline{\theta}_o - \hat{\underline{\theta}}]' &\approx \left[E_{\underline{x}} \left[\hat{I}(\hat{\theta}) \right] \right] \left[E_{\underline{x}}[\underline{\theta}_o - \hat{\underline{\theta}}] \right] \\
&= [[I(\underline{\theta}_o)]] \Sigma
\end{aligned}$$

to arrive at the same result as (7.27). In any case, (7.27) is good with sample size, but the overall error involved in this approximation $E_{\underline{x}} \left[\hat{I}(\hat{\theta}) \right] [\underline{\theta}_o - \hat{\underline{\theta}}][\underline{\theta}_o - \hat{\underline{\theta}}]'$ is hard to assess, in general. We will revisit below for the exponential family of distributions, but to be there a good approximation to $O(1/n)$.

Using either (7.27) or (7.28), along with (7.26), we have

$$E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))] \approx E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \frac{1}{2} \text{tr} [I(\underline{\theta}_o) \Sigma]$$

Recall (7.24):

$$T \approx E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))] - \frac{1}{2} \text{tr} [I(\underline{\theta}_o) \Sigma]$$

Substituting (7.29) into (7.24) we have a key result that appears in the literature:

$$T \approx E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \text{tr} [I(\underline{\theta}_o) \Sigma]$$

The literature usually presents not (7.30), but rather an alternative form based on (7.13):

$$T \approx E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o) \Sigma]]$$

The notation $\hat{\underline{\theta}}(\underline{x})$ rather than just $\hat{\underline{\theta}}$ is used above only to emphasize that on the right-hand side of (7.31) only one random variable is involved. It can be taken to refer to the actual data. From (7.30) or (7.31) it follows that a criterion for model selection (i.e., a nearly unbiased estimator of the structural form of the form

$$\hat{T} \approx \log(g(\underline{x} | \hat{\underline{\theta}})) - \widehat{\text{tr}} [I(\underline{\theta}_o) \Sigma],$$

Simple, direct estimation of Σ from one sample is not possible if only one $\hat{\theta}$ is available (a bootstrap estimator of Σ is possible). $J(\theta_o)$ and $I(\theta_o)$ are directly estimable from the single sample (7.33), but not (7.32), requires a parametrization wherein $J(\theta_o)$ whence its inverse exists. There is no loss in generality if we restrict probability distribution models have fully identifiable parameters that are of full rank.

The maximized log-likelihood $\log(g(\underline{x} | \hat{\theta}))$ in (7.31) is not a estimator of its own expectation $E_{\underline{x}}[\log(g(\underline{x} | \hat{\theta}))]$ (but is biased and inconsistent). Hence, the only problem left is to get a reliable (low, or near zero) estimate of the trace term, or at least an estimator with small mean squared error. The best model to use is the one with the largest value of \hat{T} , because it will produce a model with the smallest estimated expected K-L divergence. By convention the criterion is often stated as that of minimizing

$$-2 \log(g(\underline{x} | \hat{\theta})) + 2 \widehat{\text{tr}} [J(\theta_o)[I(\theta_o)]^{-1}]$$

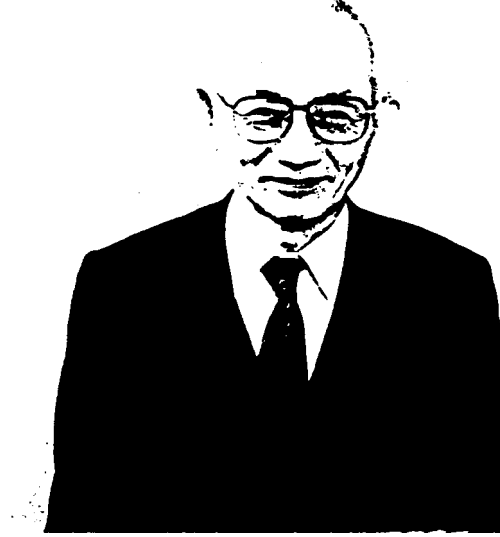
If f is a subset of g (i.e., if $g = f$ or f is contained within g), then of nested models), then $I(\theta_o) \equiv \mathcal{I}(\theta_o) = \mathcal{J}(\theta_o) = J(\theta_o)$ and $\widehat{\text{tr}} [I(\theta_o)\Sigma] = K$. Even if g is just a good model (i.e., a good approximation for f), the literature supports the idea that our best estimator is $\widehat{\text{tr}} [I(\theta_o)\Sigma] = K$ (Shibata 1989).

When the model is too restrictive to be good, the term $\widehat{\text{tr}} [I(\theta_o)\Sigma]$ will be much inflated (compared to this same term for a “good” model). It will not select that model. In this case having a good estimator of the trace term should not matter. The practical key to making AIC work is then to have some good models considered, but not too many good, but over parametrized, models. The model we mean one that is close to f in the sense of having small K-L divergence in which case such “closeness” also means that the use of g is itself a parsimonious estimator. This matter of estimation error and closeness of g to f is explored further in Section 7.6. The 7.6 derivations, and the above ideas in this paragraph, then the criterion which is seen as a special case of (7.34):

$$\text{AIC} = -2 \log(g(\underline{x} | \hat{\theta})) + 2K.$$

The generalization given by (7.34) leads to Takeuchi’s criterion (TIC) for model selection (Shibata 1989). The result is that we might use the bootstrap to compute $\widehat{\text{tr}} [I(\theta_o)\Sigma]$ and then use the TIC criterion via

$$-2 \log(g(\underline{x} | \hat{\theta})) + 2 \widehat{\text{tr}} [I(\theta_o)\Sigma]$$



Kei Takeuchi was born in 1933 in Tokyo, Japan, and graduated in 1956 from the University of Tokyo. He received a Ph.D. in economics in 1966 (Keizaigaku Haka). His research interests include mathematical statistics, econometrics, global environment, history of civilization, and Japanese economy. He is the author of many books on statistics, and the impacts of science and technology on society. He is currently on the Faculty of International Studies at Meiji Gakuin University and the University of Tokyo (recent photograph).

or even use more exact forms for the trace term. These ideas will be discussed in the next section.

First, however, there is one more crucial point on which we need to be clear: It is not required that truth f be in the set of models considered in AIC model selection. Many derivations of AIC are quite misleading under the assumption (often implicitly, hence without realizing) that $f \subset g$. Such derivations lead directly to AIC, hence by the general result of (7.33), which does not require $f \subset g$. Of course, then it is possible to see how a proper philosophy of hypothesis testing approximating models to complex truth in conjunction with the choice of $\widehat{\text{tr}}[I(\underline{\theta}_o)\Sigma] = K$ justifies use of AIC.

There are a few odds and ends worth considering at this point, and the result

$$E_{\underline{y}}[I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y}))) - I(f, g(\cdot | \underline{\theta}_o))] = \frac{1}{2} \text{tr}[J(\underline{\theta}_o)]$$

Now substitute (7.23) for $E_{\underline{x}}[\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))]$ in the above

$$\frac{1}{2} E_{\underline{y}} \left[\hat{\underline{\theta}}(\underline{y}) - \underline{\theta}_o \right]' I(\underline{\theta}_o) [\hat{\underline{\theta}}(\underline{y}) - \underline{\theta}_o],$$

which becomes $\frac{1}{2} \text{tr} [I(\underline{\theta}_o) \Sigma] = \frac{1}{2} \text{tr} [J(\underline{\theta}_o) I(\underline{\theta}_o)^{-1}]$.

It should be almost obvious (and it is true) that $\text{tr} [J(\underline{\theta}_o) I(\underline{\theta}_o)^{-1}]$, does not depend upon sample size. Rather it is about equal to K (these matters are explored in sections below). In stark contrast, quantities such as the expected log-likelihood, and both of K-L $I(f, g(\cdot | \underline{\theta}_o))$ and $E_{\underline{y}} [I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y})))]$ increase linearly in sample size n . As sample sizes, and K/n small, the ratio

$$\frac{E_{\underline{y}} [I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y})))]}{I(f, g(\cdot | \underline{\theta}_o))}$$

is essentially 1 even though the difference between expected and actual is > 0 . Thus, on an absolute scale TIC and AIC (when the model selection are producing the model estimated to produce the K-L model from the set of models considered if sample size is small.

The reason that the criterion for practical model selection is minimum K-L to minimum expected K-L as $E_{\underline{y}} [I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y})))]$ we must estimate $\underline{\theta}$ by the model-based MLE. This seems a fact has deep ramifications. It is why the K-L-based concept (the start of this section) virtually forces us to adopt $E_{\underline{y}} [I(f, g(\cdot | \hat{\underline{\theta}}(\underline{y})))]$ minimized, hence T , i.e., (7.20), to be maximized.

In this regard there is a nominally puzzling result: If we have no actual data as

$$I(f, g(\cdot | \underline{\theta}_o)) = \text{constant} - E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))]$$

and no actual data in hand, hence no estimate of $\underline{\theta}$, we may use a Taylor series expansion of $\log(g(\underline{x} | \underline{\theta}_o))$ about what would be any value of the variable of integration \underline{x} (which is *not* data) to get an expectation over the sample space of the random variable

$$E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))] = E_{\underline{x}} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x})))] - \frac{1}{2} \text{tr} [J(\hat{\underline{\theta}}(\underline{x})) I(\hat{\underline{\theta}}(\underline{x}))^{-1}]$$

The above would suggest that K-L model selection is maximizing $\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x}))) - \frac{1}{2} \text{tr} [\hat{J}(\underline{\theta}_o) \hat{I}(\underline{\theta}_o)^{-1}]$; it cannot

tion has already been integrated over the sample space, no \underline{x} and no data involved in K-L. Data cannot be manually series expansion on a random variable. Thus, the intricacy is mathematically correct, but conceptually wrong for what we are doing, hence misleading.

7.3 General K-L-Based Model Selection

7.3.1 Analytical Computation of TIC

There are other alternatives to estimation of relative K-L that try to provide a data-based estimator of the trace term. These are computationally much more intense, and the results of the trace term can be so variable, and may have its own bias, that questionable whether such approaches are worth applying (this is huge). Takeuchi (1976) proposed TIC (see also Shibata 1980, Kitagawa 1996): Select the model that minimizes (7.34) the sum of $J(\underline{\theta}_o)$ and $I(\underline{\theta}_o)$, hence getting an estimator of $\text{tr}[J(\underline{\theta}_o)]$. The estimator of $I(\underline{\theta}_o)$ is (7.8), the empirical Hessian:

$$\hat{I}(\underline{\theta}_o) = \hat{I}(\hat{\underline{\theta}}) = -\frac{\partial^2 \log(g(\underline{x} | \hat{\underline{\theta}}))}{\partial \underline{\theta}^2}.$$

General estimation of $J(\underline{\theta}_o)$ relies on recognizing the trace on n independent units of information. In the simplest case, as an *iid* sample, x_1, \dots, x_n . It is required only that the sample as having n conditionally independent components so that the trace can be computed as the sum of n terms; hence we have

$$\log(g(\underline{x} | \hat{\underline{\theta}})) = \sum_{i=1}^n \log(g_i(x_i | \hat{\underline{\theta}})).$$

For the *iid* sample case, $g_i(x_i | \hat{\underline{\theta}}) \equiv g(x_i | \hat{\underline{\theta}})$. Using here the basic sample-size one pdf and for the probability distribution, the full sample of size n is a minor abuse of notation. However, the reader will understand the meaning of the formulas as they are. We minimize notation to facilitate comprehension of concepts.

A general estimator of $J(\underline{\theta}_o)$ for TIC can be derived from

$$J(\underline{\theta}_o) = E_f \left[\left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(\underline{x} | \underline{\theta}_o)) \right]^T \right]$$

$$J(\underline{\theta}_o) = E_f \left[\left[\sum_{i=1}^n \frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \underline{\theta}_o)) \right] \left[\sum_{i=1}^n \frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \underline{\theta}_o)) \right]^T \right]$$

$$= \sum_{i=1}^K E_f \left[\frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \underline{\theta}_o)) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \underline{\theta}_o)) \right]^T$$

Therefore, we are led to use

$$\hat{J}(\underline{\theta}_o) = \sum_{i=1}^K \left[\frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \hat{\underline{\theta}})) \right] \left[\frac{\partial}{\partial \underline{\theta}} \log(g(x_i | \hat{\underline{\theta}})) \right]^T$$

A general version of TIC can be defined based on (7.36) (Shibata 1989:222):

$$\text{TIC} = -2 \log(g(\underline{x} | \hat{\underline{\theta}})) + 2 \text{tr} \left[\hat{J}(\underline{\theta}_o) [\hat{I}(\underline{\theta}_o)]^{-1} \right]$$

One selects the model that produces the smallest TIC. But for each model an asymptotically unbiased estimator of $I(f, g)$ – constant, the underlying optimization criterion is that we on average (over the set of models) minimizes this expected loss. For large n this expected criterion is almost the same as the criterion $I(f, g)$ – constant; thus using (7.38), we are selecting the K-L best model of the set of models, and whether or not f is in the model set.

The estimator $\hat{J}(\underline{\theta}_o)$ converges to $J(\underline{\theta}_o)$, and $\hat{I}(\underline{\theta}_o)$ converges to $I(\underline{\theta}_o)$. TIC is asymptotically unbiased (i.e., consistent) as a selector of the minimum expected K-L model. In practice this estimator is so variable (and is not unbiased), even for large n , that it is better to use the parsimonious “estimator” $\widehat{\text{tr}} \left[J(\underline{\theta}_o) [I(\underline{\theta}_o)]^{-1} \right] = \text{tr} \left[\hat{J}(\underline{\theta}_o) [\hat{I}(\underline{\theta}_o)]^{-1} \right]$ (we will consider the matter further in later sections). This is appropriate if we have done a good job of specifying our model, and which to select a best-fitting model.

7.3.2 *Bootstrap Estimation of TIC*

The primary value of the bootstrap method herein is to assess the uncertainty based on applying an analytical model selection criterion, AIC, AIC_c, QAIC_c, or TIC based on formulas (7.33), (7.34), (7.35), (7.36). However, a second and quite different use of the bootstrap can be made. The bootstrap method to estimate directly the quantity $T = E[\log(g(\underline{x} | \hat{\underline{\theta}}))]$ the K-L best model is the one that maximizes \hat{T} . Variations of this involve more direct bootstrap estimation of the key quantities.

From (7.24) and (7.25) (wherein $\underline{\theta}$ denotes $\underline{\theta}(\underline{x})$) we have

$$T \approx E_{\underline{x}} \left[\log(g(\underline{x} | \hat{\underline{\theta}})) \right] - \frac{1}{2} \text{tr} \left[I(\underline{\theta}_o) \Sigma \right] - \frac{1}{2} \text{tr} \left[E_{\underline{x}} \left[\hat{I}(\hat{\underline{\theta}}) \right] \right]$$

Hence, a model selection criterion can be based on

$$\hat{T} = \log(g(\underline{x} | \hat{\underline{\theta}})) - \frac{1}{2} \text{tr} \left[\hat{I}(\underline{\theta}_o) \hat{\Sigma} \right] - \frac{1}{2} \text{tr} \left[\hat{E}_{\underline{x}} \left[\hat{I}(\hat{\underline{\theta}}) \right] \right]$$

Additional approximations applied to (7.39), or to the basis to

$$\hat{T} = \log(g(\underline{x} | \hat{\underline{\theta}})) - \text{tr} \left[\hat{I}(\underline{\theta}_o) \hat{\Sigma} \right],$$

which could also be the basis for a bootstrap estimator (a

We assume that the sample structure allows a meaningful procedure (easily done in the *iid* sample case). Let a bootstrap sample be denoted by \underline{x}^* with corresponding bootstrap MLE $\hat{\underline{\theta}}^*$. The second partial derivatives will have to be determined either numerically. To avoid more notation, we do not index the derivatives but rather just note that needed summations are over B bootstrap samples.

In the bootstrap estimators, the MLE $\hat{\underline{\theta}}$ plays the role of the true parameter. The bootstrap estimators of $I(\underline{\theta}_o)$, Σ , and $E_{\underline{x}}[\hat{I}(\hat{\underline{\theta}})]$ are

$$\begin{aligned} \hat{I}(\underline{\theta}_o) &= -\frac{1}{B} \left[\sum_B \frac{\partial^2 \log(g(\underline{x}^* | \hat{\underline{\theta}}))}{\partial \underline{\theta}^2} \right], \\ \hat{\Sigma} &= \frac{1}{B} \left[\sum_B [\hat{\underline{\theta}}^* - \hat{\underline{\theta}}][\hat{\underline{\theta}}^* - \hat{\underline{\theta}}]' \right], \\ \hat{E}_{\underline{x}} \left[\left[\hat{I}(\hat{\underline{\theta}}) \right] [\underline{\theta}_o - \hat{\underline{\theta}}][\underline{\theta}_o - \hat{\underline{\theta}}]' \right] \\ &= \frac{1}{B} \left[\sum_B \left[-\frac{\partial^2 \log(g(\underline{x}^* | \hat{\underline{\theta}}^*))}{\partial \underline{\theta}^2} \right] [\hat{\underline{\theta}}^* - \hat{\underline{\theta}}][\hat{\underline{\theta}}^* - \hat{\underline{\theta}}]' \right] \end{aligned}$$

These estimators mimic the expectation over f , because the bootstrap resamples the sample, and under any model g , the bootstrap MLE $\hat{\underline{\theta}}$ (note that $\underline{\theta}_o$ varies by model g). One should bootstrap samples with every model in the set of models considered. This is made.

The above suffices to compute TIC as

$$\text{TIC} = -2 \log(g(\underline{x} | \hat{\underline{\theta}})) + 2 \text{tr} \left[\hat{I}(\underline{\theta}_o) \hat{\Sigma} \right]$$

based model selection, base the estimation of its second order term on (7.40), (7.41), and (7.42); or in a form analogous to AIC selection criterion to minimize is

$$-2 \log(g(\underline{x} | \hat{\theta})) + \text{tr} \left[\hat{I}(\underline{\theta}_o) \hat{\Sigma} \right] + \text{tr} \hat{E}_{\underline{x}} \left[\left[\hat{I}(\hat{\theta}) \right] [\underline{\theta}_o - \hat{\theta}] \right]$$

It may well be that (7.43), i.e., TIC, would suffice and \hat{T} estimator of $-2T$.

Recent work on this use of the bootstrap to find \hat{T} for model selection is found in Ishiguro, et al. (1997), Cavanaugh and Shao (1996) and Chung et al. (1996). Shibata (1997a) in a general context, theoretical properties of many alternative uses of the bootstrap to estimate the needed model selection criterion, there is no unique way to do this bootstrapping to estimate the model selection criterion, but that all reasonable bootstrap estimates are asymptotically equivalent to TIC. This use of the bootstrap bypassing concerns about all approximations used to get \hat{T} and this apparent advantage, Shibata (1997a, page 393) concludes is a reason to use the bootstrap this way to compute \hat{T} . It provides a simple nonbootstrap computation of \hat{T} (in particular, AIC).

It should thus be clear that there are two very different uses of the bootstrap in model selection. Not much used is the case of estimate of T for each model based on the full set of bootstrap estimates. More common (and more useful) use of the bootstrap in model selection is to accept some easily computable model selection criterion, then to apply that criterion to all models considered for all times created (and tabulate results like frequency of selection of each model). Use of the bootstrap leads to information about inference for model selection. [There is also a large literature on use of the bootstrap in non-K-L-based model selection; see, e.g., Breiman 1992; Hjorth 1994; Linhart and Zucchini 1986; and Shao 1996.]

7.4 AIC_c: A Second-Order Improvement

7.4.1 Derivation of AIC_c

The results above are completely general, and as such do not contradict the more specific results in the literature. In particular, for a multivariate linear structural model with homogeneous, normal errors, conditional on any regressor variables, we can get the results

$$\mu_i = E(x_i | \underline{z}) = \sum_{j=1}^K z_{ij} \beta_j, \quad i = 1, \dots, n$$

More specifically (but without explicitly denoting the “regressors” \underline{z}_i),

$$x_i = \sum_{j=1}^{K-1} z_{ij} \beta_j + \epsilon_i, \quad i = 1, \dots, n$$

where the ϵ_i are *iid* normal($0, \sigma^2$). There are thus K parameters (σ^2 is the K th one), and $g(\underline{x} | \underline{\theta})$ is given by the multivariate normal distribution (MVN); I is the $n \times n$ identity matrix. If we then we can derive the AIC_c results of Hurvich and Tsai (1989) means that either g is the true data-generating “model,” or g is a distribution and structural form as *model* g but with one or more parameters set to 0 (hence there are superfluous parameters). The superfluous parameters serve only to increase K ; hence the simplest way to get a regression model g and assume that $f \equiv g$. The derivation of AIC_c is given in some detail because of the importance of AIC_c .

Matrix notation is simpler to use, and hence $\underline{X} = Z\underline{\beta}$. Without loss of generality we assume that Z (n by $K - 1$) is full rank. The log-likelihood is

$$g(\underline{x} | \underline{\theta}) = \left[\frac{1}{\sqrt{2\pi}} \right]^n \left[\frac{1}{\sigma^2} \right]^{n/2} \exp \left[-\frac{1}{2} \frac{(\underline{X} - Z\underline{\beta})'(\underline{X} - Z\underline{\beta})}{\sigma^2} \right]$$

and we are here taking $f \equiv g$. Ignoring additive constants, the log-likelihood can be taken as

$$\log(g(\underline{x} | \underline{\theta})) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2} \frac{(\underline{X} - Z\underline{\beta})'(\underline{X} - Z\underline{\beta})}{\sigma^2}$$

The MLEs are well known here:

$$\begin{aligned} \hat{\underline{\beta}} &= (Z'Z)^{-1} Z' \underline{X}, \\ \hat{\sigma}^2 &= \frac{(\underline{X} - Z\hat{\underline{\beta}})'(\underline{X} - Z\hat{\underline{\beta}})}{n}. \end{aligned}$$

Therefore,

$$\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x}))) = -\frac{n}{2} \log(\hat{\sigma}^2) - \frac{1}{2} \frac{(\underline{X} - Z\hat{\underline{\beta}})'(\underline{X} - Z\hat{\underline{\beta}})}{\hat{\sigma}^2}$$

hence, the maximized log-likelihood is

$$\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{x}))) = -\frac{n}{2} \log(\hat{\sigma}^2) - \frac{n}{2}$$

our target

$$T = E_{\underline{x}} E_{\hat{\underline{\theta}}(\underline{y})} \left[\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))) \right],$$

where \underline{x} and \underline{y} are two independent random samples of the evaluation here we actually use the specified form of the model to take expectations with respect to $f \equiv g$). Hence, we want (the form is used here)

$$\begin{aligned} T &= E_{\underline{x}} E_{\hat{\underline{\theta}}(\underline{y})} \left[\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))) \right] \\ &= E_{\hat{\underline{\theta}}(\underline{y})} E_{\underline{x}} \left[-\frac{n}{2} \log(\hat{\sigma}_y^2) - \frac{1}{2} \frac{(\underline{X} - Z\hat{\underline{\beta}}_y)'(\underline{X} - Z\hat{\underline{\beta}}_y)}{\hat{\sigma}_y^2} \right] \end{aligned}$$

The order of integration was reversed for the right-hand side. The first task is to evaluate

$$\begin{aligned} E_{\underline{x}} \left[(\underline{X} - Z\hat{\underline{\beta}}_y)'(\underline{X} - Z\hat{\underline{\beta}}_y) \right] &= E_{\underline{x}} \left[((\underline{X} - Z\hat{\underline{\beta}}) + (Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y))'((\underline{X} - Z\hat{\underline{\beta}}) + (Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)) \right] \\ &= E_{\underline{x}} \left[(\underline{X} - Z\hat{\underline{\beta}})'(\underline{X} - Z\hat{\underline{\beta}}) \right] + E_{\underline{x}} \left[2(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)'(\underline{X} - Z\hat{\underline{\beta}}) \right] \\ &\quad + E_{\underline{x}} \left[(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)'(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y) \right] \\ &= E_{\underline{x}} \left[(\underline{X} - Z\hat{\underline{\beta}})'(\underline{X} - Z\hat{\underline{\beta}}) \right] + \left[2(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)' E_{\underline{x}}(\underline{X} - Z\hat{\underline{\beta}}) \right] \\ &\quad + \left[(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)'(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y) \right]. \end{aligned}$$

The middle term above vanishes because $E_{\underline{x}}(\underline{X}) = Z\hat{\underline{\beta}}$. The first and third terms above is identical to $E_{\underline{x}}(\underline{\epsilon}'\underline{\epsilon}) = n\sigma^2$. So we have

$$E_{\underline{x}} \left[(\underline{X} - Z\hat{\underline{\beta}}_y)'(\underline{X} - Z\hat{\underline{\beta}}_y) \right] = n\sigma^2 + \left[(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)'(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y) \right]$$

Using this partial result we have

$$T = E_{\hat{\underline{\theta}}(\underline{y})} \left[-\frac{n}{2} \log(\hat{\sigma}_y^2) \right] - \frac{1}{2} E_{\hat{\underline{\theta}}(\underline{y})} \left[\frac{n\sigma^2 + \left[(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y)'(Z\hat{\underline{\beta}} - Z\hat{\underline{\beta}}_y) \right]}{\hat{\sigma}_y^2} \right]$$

The first term above does not need to be evaluated further as it is the leading term in the expected log-likelihood. Also, at this point the designation of $\underline{\theta}$ as being based on sample \underline{y} . The designation of $\hat{\underline{\theta}}$ as being based on sample \underline{y} is just dummy arguments in integrals. Consequently, in the

$$T = E_{\hat{\theta}} \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{1}{2} E_{\hat{\theta}} \left[\frac{[(Z\underline{\beta} - Z\underline{\hat{\beta}})'(Z\underline{\beta} - Z\underline{\hat{\beta}})]}{\hat{\sigma}^2} \right]$$

Now we make use of another well-known result in statistics. Under a linear model structure with errors as *iid* normal(0, σ^2) and $\hat{\sigma}^2$ are independent random variables. Therefore, the term in (7.45) partitions into two multiplicative parts, as follows

$$T = E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{1}{2} E_{\hat{\beta}} \left[n\sigma^2 + [(Z\underline{\beta} - Z\underline{\hat{\beta}})'(Z\underline{\beta} - Z\underline{\hat{\beta}})] \right]$$

As a next step, rewrite the needed expectation of the quadratic term above as

$$E \left[(Z\underline{\beta} - Z\underline{\hat{\beta}})'(Z\underline{\beta} - Z\underline{\hat{\beta}}) \right] = \text{tr} \left[(Z'Z) E \left[(\underline{\hat{\beta}} - \underline{\beta})(\underline{\hat{\beta}} - \underline{\beta})' \right] \right]$$

The expectation on the right-hand side above, i.e., $E[(\underline{\hat{\beta}} - \underline{\beta})(\underline{\hat{\beta}} - \underline{\beta})']$ is the sampling variance–covariance matrix of $\underline{\hat{\beta}}$, which is known to be $\sigma^2(Z'Z)^{-1}$. Thus, for the $K - 1$ square identity matrix I ,

$$E \left[(Z\underline{\beta} - Z\underline{\hat{\beta}})'(Z\underline{\beta} - Z\underline{\hat{\beta}}) \right] = \text{tr}[\sigma^2 I] = \sigma^2 \text{tr}[I] = \sigma^2(K - 1)$$

Putting it all together to this point in the derivation, we have

$$T = E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{1}{2} [(n + K - 1)\sigma^2] E_{\hat{\sigma}^2} \left[\frac{1}{\hat{\sigma}^2} \right]$$

To finish the process we relate $\hat{\sigma}^2$ to a central chi-square distribution, namely χ_{df}^2 on $n - (K - 1)$ degrees of freedom, df. These results are well known in statistical theory:

$$\frac{n\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-K+1}^2.$$

So we now rearrange (7.46) to be

$$T = E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{1}{2} [(n + K - 1)n] E \left[\frac{1}{n\hat{\sigma}^2} \right]$$

$$T = E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{n}{2} (n + K - 1) E \left[\frac{1}{\chi_{n-K+1}^2} \right]$$

Yet another known exact result is

$$E \left[\frac{1}{\chi_{df}^2} \right] = \frac{1}{df - 2}$$

(assuming $df > 2$).

$$T = E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{n}{2} (n + K - 1) \left[\frac{1}{n - K - 1} \right]$$

This result is exact. No approximations were made in its derivation, it applies only to the particular context of its derivation, constraint $f \subseteq g$. Some more simplification of (7.47):

$$\begin{aligned} T &= E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{n}{2} \left[\frac{n + K - 1}{n - K - 1} \right] \\ &= E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{n}{2} \left[1 + \frac{2K}{n - K - 1} \right] \\ &= E \left[-\frac{n}{2} \log(\hat{\sigma}^2) \right] - \frac{n}{2} - \frac{nK}{n - K - 1} \\ &= E \left[-\frac{n}{2} \log(\hat{\sigma}^2) - \frac{n}{2} \right] - \frac{nK}{n - K - 1}. \end{aligned}$$

The term above within the expectation operator is the log-likelihood. Thus we have

$$\begin{aligned} T &= E \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] - \frac{nK}{n - K - 1} \\ &= E \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] - \frac{(n - K - 1 + K)}{n - K - 1} \\ &= E \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] - K - \frac{K(K + 1)}{n - K - 1} \end{aligned}$$

If we convert this to an AIC result, we have, as an exact result

$$\begin{aligned} -2T &= -2E \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] + 2K + \frac{2K(K + 1)}{n - K - 1} \\ &= E(\text{AIC}) + \frac{2K(K + 1)}{n - K - 1} = E(\text{AIC}_c). \end{aligned}$$

This result thus motivates use of the term $2K(K + 1)$ as a small-sample-size bias-correction term added to AIC. This correction term in the fixed-effects linear model with normal errors and constant variance is well known. Under different sorts of models, a different small-sample-size correction may arise (the matter is explored some in the next subsection). The correction given by (7.48) seems useful in other contexts, especially in the case where n is also large relative to n . Without exception, if sample size of “AIC_c” is required for good model selection results, then the correction (7.48) unless a more exact small-sample correction to AIC is available.

so for multivariate) models, with $g = f$ to facilitate a derivation of a small-sample adjustment to AIC, will lead to the same terms. This section is just a brief elaboration of this idea.

The simplest case to present arises for a situation analogous to ANOVA, but we let the within-subgroup variance differ. This can be generalized to having m subsets of data, each of size n_i , and the full model is as used in Section 7.4.1 above, but with $\{\theta\} = \{\beta_i, \sigma_i^2, i = 1, \dots, m\}$ (this might be a global model). Let each parameter subset be of size K_i ; hence $K = K_1 + \dots + K_m$. It will be almost obvious, after some thought, that for this situation the correction to AIC is

$$-2T = -2E \left[\log(g(\underline{x} | \hat{\theta}(\underline{x}))) \right] + 2 \sum_{i=1}^m \left[K_i + \frac{2K_i}{n_i} \right]$$

hence,

$$AIC_c = -2 \log(g(\underline{x} | \hat{\theta}(\underline{x}))) + 2K + \sum_{i=1}^m \left[\frac{2K_i}{n_i} \right]$$

The reason that there are m “correction” terms is that we have m different variance parameters. One can thus envision many other forms of AIC_c must differ from that of the simple normal model where only one estimated σ^2 .

Another informative exact calculation of the bias term, T , can be obtained for the case of the model and truth being the one-parameter exponential distribution (hence $K = 1$):

$$g(x | \lambda) = \frac{1}{\lambda} e^{-x/\lambda}.$$

For an *iid* sample from $g(x | \lambda)$, let $S = x_1 + \dots + x_n$. Then

$$\log(g(\underline{x} | \lambda)) = -n \log(\lambda) - S/\lambda.$$

The MLE is $\hat{\lambda} = S/n$, so

$$\log(g(\underline{x} | \hat{\lambda})) = -n \log(\hat{\lambda}) - n.$$

The target to be unbiasedly estimated is

$$T = E_{\underline{x}} E_{\underline{y}} \left[-n \log(\hat{\lambda}) - S_{\underline{y}}/\hat{\lambda} \right],$$

where the sum $S_{\underline{y}}$ is based on an independent sample of size n on sample \underline{x} . It is easy to evaluate the above T to be

$$T = E \left[\log(g(\underline{x} | \hat{\lambda})) \right] + n - n^2 \lambda E \left[\frac{1}{S} \right].$$

$$-2T = -2E(\log(g(\underline{x} | \hat{\lambda}))) + 2 + \frac{2}{n -$$

Recall that here $K = 1$, so the corresponding total bias-correction term would be $2 + 4/(n - 2)$. The point is that this would be $2K$ plus a small-sample correction term that would depend on the model assumed. It is reasonable to think that this small-sample correction term should be $O(1/n)$.

Theoretically, when $f \subset g$ the error in using K as a correction to $\hat{T} = \log(g(\underline{x} | \hat{\theta}))$ is always $O(1/n)$, and Hurvich and Tsai (1989) like a good general choice. There is, however, considerable research on improved bias terms for AIC_c -type criteria. In this research, one opportunity is that of when the random variable is discrete (Sugiura 1978, Shibata 1997b), such as Poisson, binomial, or multinomial (also logistic regression), because then we can get parameter estimates that are exactly the value 0. This creates a problem in evaluating the theoretical properties of the selection criterion because we encounter the need to compute expectations that are not defined (see, e.g., Burnham et al. 1994). AIC is not a good choice. Small-sample properties are now more problematic, as the small-sample bias-correction term needed to define an AIC_c . Operating on the basis of AIC -based model selection for count-type data need more research on small sample sizes.

7.5 Derivation of AIC for the Exponential Family of Distributions

A generalization of normality-based models is found in the theory of distributions. The realizations that (1) many common statistical analyses are based on exponential family models, and (2) logistic regression is in the exponential family and leads to exact results (i.e., AIC_c) motivated us to show the derivation of AIC for this restricted but very useful case. The canonical representation of the exponential family pdf involves sums of functions of the sample values. We denote these sums by S_j .

A suitable canonical representation for the exponential family distributions is

$$\begin{aligned} g(\underline{x} | \underline{\theta}) &= \exp \left[\left[\sum_{j=1}^K S_j \theta_j \right] + H(\underline{\theta}) + G(\underline{\theta}) \right] \\ &= \exp \left[\underline{S}' \underline{\theta} + H(\underline{\theta}) + G(\underline{\theta}) \right]. \end{aligned}$$

A K -element vector of sufficient statistics is $\underline{S} = (S_1, \dots, S_K)$. In the canonical representation of (7.50) the parameter $\underline{\theta}$ is a 1-to-1 transformation of another K -dimensional parameter $\underline{\theta}_o$. There is no loss of generality in allowing any such 1-to-1 transformation. We will revisit this matter and show why it is so at the end of this section.

Our goal is to evaluate

$$T = E_{\underline{x}} E_{\underline{y}} \left[\log(g(\underline{x} | \hat{\underline{\theta}}_y)) \right] = E_{\underline{x}} E_{\underline{y}} \left[\underline{S}'_x \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y) + G(\underline{S}_x) \right]$$

Here, \underline{S}_x and $\hat{\underline{\theta}}_y$ are thought of as based on independent samples. We have also simplified the notation, now using $\hat{\underline{\theta}}_y$ rather than $\hat{\underline{\theta}}(\underline{y})$. Formula (7.51) above can be rewritten as

$$\begin{aligned} T &= E_{\underline{x}} E_{\underline{y}} \left[(\underline{S}_x - \underline{S}_y + \underline{S}_y)' \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y) + G(\underline{S}_x) \right] \\ &= E_{\underline{x}} E_{\underline{y}} \left[(\underline{S}_x - \underline{S}_y)' \hat{\underline{\theta}}_y + \underline{S}'_y \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y) + G(\underline{S}_x) \right] \\ &= E_{\underline{x}} E_{\underline{y}} \left[\underline{S}'_y \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y) + G(\underline{S}_x) \right] + E_{\underline{x}} E_{\underline{y}} \left[(\underline{S}_x - \underline{S}_y)' \hat{\underline{\theta}}_y \right] \\ &= \left[E_{\underline{y}} (\underline{S}'_y \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y)) + E_{\underline{x}} (G(\underline{S}_x)) \right] + E_{\underline{y}} \left[(E_{\underline{x}} (\underline{S}_x - \underline{S}_y))' \hat{\underline{\theta}}_y \right] \end{aligned}$$

The interchangeability of integration arguments now is useful. It is possible because both expectations are with respect to f ; $E_{\underline{x}} (G(\underline{S}_x)) = E_{\underline{y}} (G(\underline{S}_x))$. Also, for simplicity we will use $E_{\underline{x}} (\underline{S}_x) = E(\underline{S})$.

$$T = E_{\underline{y}} \left[\underline{S}'_y \hat{\underline{\theta}}_y + H(\hat{\underline{\theta}}_y) + G(\underline{S}_y) \right] + E_{\underline{y}} \left[(E(\underline{S}) - \underline{S}_y)' \hat{\underline{\theta}}_y \right]$$

Changing the argument from y to x in the first part above and using the result for this exponential family case, we have

$$\begin{aligned} T &= E_{\underline{x}} (\log(g(\underline{x} | \hat{\underline{\theta}}))) + E_{\underline{y}} \left[(E(\underline{S}) - \underline{S}_y)' \hat{\underline{\theta}}_y \right] \\ &\equiv E_{\underline{x}} (\log(g(\underline{x} | \hat{\underline{\theta}}))) - E_{\underline{y}} \left[(\underline{S}_y - E(\underline{S}))' \hat{\underline{\theta}}_y \right] \end{aligned}$$

Formula (7.52) is an exact result and clearly shows the bias. We will change from $E_{\underline{x}} (\log(g(\underline{x} | \hat{\underline{\theta}})))$ to get T :

$$\text{Bias} = E_{\underline{y}} \left[(\underline{S}_y - E(\underline{S}))' (\hat{\underline{\theta}}_y - \underline{\theta}_*) \right].$$

The notation used here is $E(\hat{\underline{\theta}}_y) = \underline{\theta}_*$ to denote the exact MLE for the given sample size n and model g ; $\underline{\theta}_o \approx \underline{\theta}_*$ for equality.

$$\text{Bias} = \text{E} \left[(\underline{S} - \text{E}(\underline{S}))'(\hat{\underline{\theta}} - \underline{\theta}_*) \right],$$

$$\text{Bias} = \text{tr} \text{E} \left[(\hat{\underline{\theta}} - \underline{\theta}_*)(\underline{S} - \text{E}(\underline{S}))' \right] = \text{tr} \left[\text{COV}(\hat{\underline{\theta}}, \underline{S}) \right]$$

Hence for the exponential family *an exact result is*

$$T = \text{E}_{\underline{x}}(\log(g(\underline{x} | \hat{\underline{\theta}}))) - \text{tr} \left[\text{COV}(\hat{\underline{\theta}}, \underline{S}) \right]$$

(something similar appears in Bonneu and Milhaud 1994), of covariance elements, $\text{COV}(\hat{\underline{\theta}}, \underline{S})$, can be approximated by various methods. If the exact covariance matrix can be found, then the result for the needed bias term above (Hurvich and Tsai 1999) is such an exact evaluation for the normal distribution case. This may seem not very useful because it seems to apply only to the normal form of the exponential family. This is not true; the matter of a canonical result will be addressed below.

Before further evaluation of the bias term, we consider the Hessian. First,

$$\log(g(\underline{x} | \underline{\theta})) = \underline{S}'\underline{\theta} + H(\underline{\theta}) + G(\underline{S}),$$

so

$$\frac{\partial \log(g(\underline{x} | \underline{\theta}))}{\partial \underline{\theta}} = \underline{S} + \frac{\partial H(\underline{\theta})}{\partial \underline{\theta}},$$

$$\frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \underline{\theta}^2} = \frac{\partial^2 H(\underline{\theta})}{\partial \underline{\theta}^2},$$

and thus

$$I(\underline{\theta}_o) = \text{E}_f \left[-\frac{\partial^2 \log(g(\underline{x} | \underline{\theta}_o))}{\partial \underline{\theta}^2} \right] = -\frac{\partial^2 H(\underline{\theta}_o)}{\partial \underline{\theta}^2}$$

It follows that the MLE satisfies

$$\underline{S} = -\frac{\partial H(\hat{\underline{\theta}})}{\partial \underline{\theta}}.$$

It is worth noting here that $\underline{\theta}_o$ satisfies

$$\text{E}_f(\underline{S}) = -\frac{\partial H(\underline{\theta}_o)}{\partial \underline{\theta}}.$$

This is an exact result, whereas $\text{E}(\hat{\underline{\theta}}) \approx \underline{\theta}_o$ is (in general)

The formula for $J(\underline{\theta}_o)$, based on (7.12), becomes

$$J(\underline{\theta}_o) = \text{E}_f \left[\underline{S} - \text{E}_f(\underline{S}) \right] \left[\underline{S} - \text{E}_f(\underline{S}) \right]'$$

$$\hat{J}(\underline{\theta}_o) = \frac{n}{n-1} \left[\sum_{i=1}^n [s_i - \bar{s}] [s_i - \bar{s}]' \right]$$

Returning now to the evaluation of the bias term, a first-order Taylor expansion gives us

$$-\frac{\partial H(\hat{\underline{\theta}})}{\partial \underline{\theta}} \approx -\frac{\partial H(\underline{\theta}_o)}{\partial \underline{\theta}} - \frac{\partial^2 H(\underline{\theta}_o)}{\partial \underline{\theta}^2} (\hat{\underline{\theta}} - \underline{\theta}_o)$$

hence

$$\underline{S} \approx E(\underline{S}) + I(\underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o), \quad O_p(1/\sqrt{n})$$

Inserting (7.58) into the exact result (7.53) as well as also replacing the exact $\underline{\theta}_*$ (inasmuch as we are now replacing an exact result by an approximation result anyway), we have

$$\begin{aligned} \text{Bias} &\approx E \left[\left[I(\underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o) \right]' (\hat{\underline{\theta}} - \underline{\theta}_o) \right] \\ &= E \left[(\hat{\underline{\theta}} - \underline{\theta}_o)' I(\underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o) \right] \\ &= E \text{tr} \left[I(\underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)' \right] \\ &= \text{tr} \left[I(\underline{\theta}_o) E \left[(\hat{\underline{\theta}} - \underline{\theta}_o)(\hat{\underline{\theta}} - \underline{\theta}_o)' \right] \right] = \text{tr} [I(\underline{\theta}_o) \Sigma_{\hat{\underline{\theta}}} - I(\underline{\theta}_o) \Sigma_{\underline{\theta}_o}] \end{aligned}$$

Thus we have shown that in this common case of a linear model,

$$T \approx E_{\underline{x}}(\log(g(\underline{x} | \hat{\underline{\theta}}))) - \text{tr} [I(\underline{\theta}_o) \Sigma_{\hat{\underline{\theta}}} - I(\underline{\theta}_o) \Sigma_{\underline{\theta}_o}]$$

(the approximation is to $O(1/n)$). Note that this derivation avoids any problems like those in approximation (7.28) in the general case of AIC in Section 7.2.

These results can be extended to any parametrized family model, because then we just have a 1-to-1 transformation $\underline{\beta}$ via some set of K functions, denoted here by $\underline{W}(\underline{\theta}) = \underline{\beta}$, and let $\Sigma_{\underline{\theta}}$ and $\Sigma_{\underline{\beta}}$ be the variance–covariance matrices for the two parametrizations. An expected matrix of mixed second-order derivatives as per (7.55) exists for the $\underline{\beta}$ parametrization; denote it by $\Sigma_{\underline{\beta}}$. The Jacobian of \underline{W} , evaluated at $\underline{\theta}_o$, be

$$J_w = \left\{ \frac{\partial W_i(\underline{\theta}_o)}{\partial \theta_j} \right\}.$$

Then

$$J_w \Sigma_{\underline{\theta}} J_w' = \Sigma_{\underline{\beta}},$$

Both the K-L-based target and the expected log-likelihood are invariant to 1-to-1 parameter transformations, so this must also be true for the K-L-based correction. That is, any likelihood and MLE-based model selection criterion ought to be invariant to 1-to-1 reparametrizations of the model. This is the case here:

$$\begin{aligned}\text{tr}(I(\underline{\beta}_o)\Sigma_\beta) &= \text{tr}[(J'_w)^{-1}I(\underline{\theta}_o)(J_w)^{-1}J_w\Sigma_\theta J'_w] \\ &= \text{tr}[(J'_w)^{-1}I(\underline{\theta}_o)\Sigma_\theta J'_w] \\ &= \text{tr}[I(\underline{\theta}_o)\Sigma_\theta J'_w(J'_w)^{-1}] = \text{tr}[I(\underline{\theta}_o)\Sigma_\theta]\end{aligned}$$

Note, however, that if we were to estimate this trace term, it would perform better under some parametrizations than under others.

One last point here: It is certainly still true that

$$\text{tr}[I(\underline{\theta}_o)\Sigma] = \text{tr}[J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}].$$

So an alternative to (7.59) is

$$T \approx E_{\underline{x}}(\log(g(\underline{x}|\hat{\underline{\theta}}))) - \text{tr}[J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}].$$

This could be directly proven here, based on the simple result

$$\frac{\partial \log(g(\underline{x}|\underline{\theta}))}{\partial \underline{\theta}} = \underline{S} - E(\underline{S})$$

and (7.58) to derive $[I(\underline{\theta}_o)]^{-1}J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1} = \Sigma$.

For TIC we can use $\hat{J}(\underline{\theta}_o)$ from (7.57) and from (7.55)

$$\hat{I}(\underline{\theta}_o) = -\frac{\partial^2 H(\hat{\underline{\theta}})}{\partial \theta^2},$$

getting an estimator of $\text{tr}[J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}]$ that can be used to estimate the variance even if the parametrization of interest (and used for the model selection) is $\underline{\beta}_o = \underline{W}(\underline{\theta}_o)$, not $\underline{\theta}_o$.

Working with exponential family cases facilitates some simplification of both $\text{tr}[J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}]$, relative to the value K , and the estimator $\text{tr}[\hat{J}(\underline{\theta}_o)[\hat{I}(\underline{\theta}_o)]^{-1}]$. These topics, and others, are discussed in the next section.

7.6 Evaluation of $\text{tr}(J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1})$ and

The general derivation of a K-L-based model selection criterion is given in (7.31) and hence (7.33). By “a general derivation,” we mean

L -based model selection does not appear to be AIC. Kaufman's bias correction term subtracted from the expected maximum likelihood estimate T is $\text{tr}(J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1})$ (Takeuchi 1976). In deriving this and need not be, any assumption that any of the candidate models is truth. However, in general we know with certainty that $\text{tr}(J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}) = 0$ only when $J(\underline{\theta}_o) = I(\underline{\theta}_o)$ (this is sufficient but not necessary for equality is certain, in general, only when f is a special case of model g equals or is a generalization of "truth." This correction is expected, so how good is the approximation $\text{tr}(J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1})$ to the truth is more general than the model, but the model is not a good approximation to truth? We make here some limited, but useful, progress. Extensive theory, simulation studies, and experience (e.g., Takeuchi 1986:176–182, especially results such as in their Table 1) suggest to give us full confidence in when we can expect reliable results versus when we might have to use TIC. Below, we establish results for some models within the exponential family of distributions.

7.6.1 Comparison of AIC Versus TIC in a Very Simple Situation

We consider two simple one-parameter distributions: standard normal and half-normal. For each distribution we can compare the two selection criteria, assuming either that the distribution is truth or that the other distribution is truth f and the given distribution is model g (so $f \neq g$). We will also examine the estimators of the true distribution in TIC model selection and contrast TIC selection with AIC selection in the two distributions as models. This is a convenient situation because both distributions are in the exponential family.

For the negative-exponential distribution let $S = x_1 + \dots + x_n$

$$\begin{aligned} g(\underline{x} | \lambda) &= \frac{1}{\lambda^n} e^{-S/\lambda}, \\ \log(g(\underline{x} | \lambda)) &= -n \log(\lambda) - S/\lambda, \\ \frac{\partial}{\partial \lambda} \log(g(\underline{x} | \lambda)) &= -n/\lambda + S/\lambda^2, \end{aligned}$$

so $\hat{\lambda} = S/n = \bar{x}$; also here $E(x) = \lambda$.

Direct verification yields the following results:

$$\begin{aligned} I(\lambda) &= \frac{n}{\lambda^2}, \\ \hat{I}(\hat{\lambda}) &= \frac{n}{\hat{\lambda}^2}, \end{aligned}$$

$$J(\lambda) = \sum_{i=1}^n \left[\frac{1}{\lambda} + \frac{1}{\lambda^2} \right] = \left[\sum_{i=1}^n \frac{1}{\lambda^4} \right]$$

The true $J(\lambda) = \hat{E}_f(\hat{J}(\lambda))$, and $\hat{J}(\hat{\lambda})$ is the empirical estimator of $J(\lambda)$. We assume that the negative-exponential model is truth, (because $E(x^2) = 2\lambda^2$),

$$J(\lambda) = \frac{n}{\lambda^2}.$$

Clearly, if $f = g$, then here $\text{tr}[J(\lambda)[I(\lambda)]^{-1}] = 1$. The estimator of this trace is

$$\text{tr}[\hat{J}(\hat{\lambda})[\hat{I}(\hat{\lambda})]^{-1}] = \frac{1}{n\bar{x}^2} \left[\sum_{i=1}^n (x_i)^2 \right] -$$

or

$$\text{tr}[\hat{J}(\hat{\lambda})[\hat{I}(\hat{\lambda})]^{-1}] = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n\bar{x}^2}.$$

This estimator of the trace is the same as $\frac{n-1}{n}(\widehat{\text{cv}})^2$ (where $\widehat{\text{cv}}$ is the sample standard deviation, the denominator of the sample s^2 as per convention), which is quite variable. The trace estimator of (7.60) does converge to 1 if $f = g$. This trace estimator is scale-invariant, so we can obtain conditional properties by Monte Carlo methods with a single run at each n , and for each case of truth being either the negative-exponential distribution (any value of λ can be used) or half-normal distribution (any σ can be used). In the case where the negative-exponential model is truth, the following is a revised, nearly unbiased, version of (7.60):

$$\widehat{\text{tr}}[J(\lambda)[I(\lambda)]^{-1}] = \frac{n}{n-1}(\widehat{\text{cv}})^2.$$

The estimated mean and standard deviation of the estimator based on (7.60) and (7.61) are given below, based on 1000 simulations.

n	Eq. (7.60)		Eq. (7.61)	
	mean	st.dev.	mean	st.dev.
20	0.90	0.37	1.00	0.41
50	0.96	0.26	1.00	0.27
100	0.98	0.19	1.00	0.20
500	1.00	0.09	1.00	0.09

Notice the substantial standard deviation of either trace estimator for small n . The estimator based on (7.61) has less variability than the estimator based on (7.60). The direct estimator has bias, but the bias is small for large n . We also see that the direct estimator has bias.

$$f(x | \sigma^2) = \sqrt{\frac{2}{\pi \sigma^2}} \exp \left[-\frac{1}{2} \left(\frac{x}{\sigma} \right)^2 \right]$$

Under this distribution $E(x^2) = \sigma^2$, and direct integration

$$E(x) = \sigma \sqrt{\frac{2}{\pi}}.$$

If the half-normal distribution is truth and the negative-exponential is model, then λ_o as a function of σ^2 is found from

$$E_f(x) = \lambda_o = \sigma \sqrt{\frac{2}{\pi}}.$$

For example, with $\sigma^2 = 1$, $\lambda_o = 0.79788$ is the K-L best

Note the usage and concepts here: We denote the K-L best negative-exponential model based on λ_o to distinguish that the corresponding negative-exponential model based on λ_o is the K-L best negative-exponential model for the underlying truth. By denoting this value as λ_o we are emphasizing that all we have is the K-L best negative-exponential model; it may be a poor model; it certainly may not be truth.

The above expectations producing $J(\lambda)$ and $I(\lambda)$ were for the negative-exponential as g (ignoring what f might be), but those expectations with respect to f as half-normal. Directly, $I(\lambda_o) = n/\lambda_o^2$, whereas

$$\begin{aligned} J(\lambda_o) &= n E_f \left[\frac{x}{\lambda_o^2} - \frac{1}{\lambda_o} \right]^2 \\ &= I(\lambda_o) \left[\frac{\sigma^2}{\lambda_o^2} - 1 \right] \\ &= I(\lambda_o) \left[\frac{\pi}{2} - 1 \right]. \end{aligned}$$

Hence, when f is half-normal and g is negative-exponential

$$\text{tr} [J(\lambda_o)[I(\lambda_o)]^{-1}] = \frac{\pi}{2} - 1 = 0.5708$$

This trace term is not very close to 1, the number AIC big relative difference (i.e., 0.5708 versus 1) results because the negative-exponential model is a very poor approximation to the half-normal. Note that this trace term is < 1 (K is 1 here).

$$\begin{aligned}\hat{\sigma}^2 &= \frac{\sum_{i=1}^n (x_i)^2}{n}, \\ I(\sigma^2) &= \frac{n}{2\sigma^4}, \\ \hat{J}(\sigma^2) &= \sum_{i=1}^n \left[\frac{-1}{2\sigma^2} + \frac{(x_i)^2}{2\sigma^4} \right]^2 \\ &= I(\sigma^2) \left[\frac{1}{2n} \right] \sum_{i=1}^n \left[\frac{(x_i)^2}{\sigma^2} - 1 \right]^2.\end{aligned}$$

The key part of the empirical estimator $\hat{J}(\hat{\sigma}^2)$ is again a sample variation, but here it is for the variable x^2 . Denote this $\widehat{\text{cv}}$, we can use the notation

$$\hat{J}(\hat{\sigma}^2) = \hat{I}(\hat{\sigma}^2) \left[\frac{n-1}{2n} \right] [\widehat{\text{cv}}(x^2)]^2.$$

Thus, for the half-normal distribution being the model, the trace is

$$\text{tr}[\hat{J}(\hat{\sigma}^2)[\hat{I}(\hat{\sigma}^2)]^{-1}] = \left[\frac{n-1}{2n} \right] [\widehat{\text{cv}}(x^2)].$$

If truth is the half-normal model, this quantity will converge to 1. (7.62) is scale-invariant. For the half-normal model as truth, the mean and standard deviation of (7.62) by Monte Carlo method for 500 samples at each n . This led to a nearly unbiased version of

$$\widehat{\text{tr}}[J(\sigma^2)[I(\sigma^2)]^{-1}] = \frac{1}{2} \left[\frac{n}{n-1} \right]^2 [\widehat{\text{cv}}(x^2)]^2$$

that we then also used in the simulations. The results are

n	Eq. (7.62)		Eq. (7.63)	
	mean	st.dev.	mean	st.dev.
20	0.87	0.39	1.01	0.45
50	0.94	0.30	1.00	0.32
100	0.97	0.23	1.00	0.23
500	0.99	0.11	1.00	0.11

The main point from the above is how variable the trace of

To complete a set of analytical results (useful for validation of the results) we computed the value of $\text{tr}[J(\sigma_o^2)[I(\sigma_o^2)]^{-1}]$ for the model when truth is the negative-exponential distribution with $\sigma_o^2 = \text{E}(x^2) = 2\lambda^2$, because the expectation of x must be finite with respect to the negative-exponential distribution. We find that

$$J(\sigma_o) = I(\sigma_o) \left[\frac{1}{2} \right] E \left[\frac{1}{\sigma_o^2} - 1 \right] = 2.5$$

Thus when truth is the negative-exponential, and the model is the half-normal (the terrible model in this case), $\text{tr}[J(\sigma_o^2)[I(\sigma_o^2)]^{-1}] = 2.5$ (not 1). Done the other way around we had the trace < 1 . It turns out that the trace function under model misspecification can be either above or below 1, by situation, and in some situations the trace function can be significantly above a misspecified model.

We can now compare AIC versus TIC model selection when truth is negative-exponential or half-normal.

For the negative-exponential model:

$$\text{AIC} = 2n[\log(\bar{x}) + 1] + 2,$$

$$\text{TIC} = 2n[\log(\bar{x}) + 1] + 2\frac{n-1}{n}(\widehat{\text{CV}}(x))$$

$$\text{TICu} = 2n[\log(\bar{x}) + 1] + 2\frac{n}{n-1}(\widehat{\text{CV}}(x))$$

For the half-normal model (using $\hat{\sigma}^2 = \text{mean of the } x_i^2$):

$$\text{AIC} = n [\log(\hat{\sigma}^2) + 1 - \log(2/\pi)] + 2,$$

$$\text{TIC} = n [\log(\hat{\sigma}^2) + 1 - \log(2/\pi)] + \left[\frac{n-1}{n} \right]$$

$$\text{TICu} = n [\log(\hat{\sigma}^2) + 1 - \log(2/\pi)] + \left[\frac{n}{n-1} \right]$$

In both cases here TICu means just that the estimator of σ^2 is almost unbiased, as opposed to the direct, biased, plug-in estimator. It needed coefficients of variation.

Table 7.1 shows some results. The point of this brief comparison is to show something about AIC versus TIC in a simple setting, especially where the two models they would give greatly different results. The context here is that two models are compared. Moreover, one or the other model is the true data-generating distribution (i.e., truth). We did not consider the case where so the only possible criterion to use to compare performance is that TIC is rate of selection of the true model. We did not wish to do any serious evaluation of AIC versus TIC under full-blown misspecification of complex truth, and a set of approximating models, where the criterion of evaluation is how well a selection procedure does at selecting the true model (technically, we would be selecting the expected Kullback-Leibler divergence).

Several inferences supported by Table 7.1, and by all the simulations examined for this situation, surprised us. For the case where the negative-exponential model is true, the selection results based on

average percent correct is based on equal weighting of the two cases.

sample size, n	selection criterion	truth		per co
		negative expon.	half- normal	
20	AIC	64	85	
	TIC	73	77	
	TICu	75	75	
50	AIC	82	92	
	TIC	87	87	
	TICu	87	87	
100	AIC	93	97	
	TIC	95	95	
	TICu	95	95	
500	AIC	100	100	1
	TIC	100	100	1
	TICu	100	100	1

as good or better than those under AIC. The improvement is at small sample sizes, wherein an “AIC_c” should be used. For the case that the half-normal model is true, the selection criteria under AIC were uniformly as good or better than those under TIC. The bias-correction of the trace estimator makes no real difference. To know a priori which (if either) model was true. If we do not know the true model, percent-correct selection based on the idea that we have no way to justify any weighting other than a 50:50 weighting of these two models, on average, no advantage at all for TIC over AIC. Clearly, the extent to which these results would generalize.

7.6.2 Evaluation Under Logistic Regression

Logistic regression is used often, therefore we illustrate the exponential family model, and we explore the above trace estimator. Let x_i be a Bernoulli random variable with true probability μ_i (and probability $1 - \mu_i$ of being 0). For a sample of n independent observations, we perform a logistic regression analysis on some assumed model for the μ_i . In order to do this, we adopt the notation for the model as $p_i \equiv p_i(\theta)$, where θ is a K -dimensional vector of parameters imposed on these p_i , as a function of a K -dimensional vector of parameters. The relevant pdf, or likelihood (the same notation continues)

$$g(\underline{x} | \underline{\theta}) = \prod_{i=1}^n (p_i)^{x_i} (q_i)^{1-x_i}.$$

We assume that known covariates \underline{z}_i , as $K \times 1$ column vector with each observation, x_i , and an explanatory structural parameter

$$p_i = \frac{1}{1 + e^{-\underline{z}_i' \underline{\theta}}}, \text{ or } q_i = \frac{e^{-\underline{z}_i' \underline{\theta}}}{1 + e^{-\underline{z}_i' \underline{\theta}}},$$

which is equivalent to

$$\log [p_i / (1 - p_i)] = \underline{z}_i' \underline{\theta}.$$

A modest amount of algebra gives the result

$$g(\underline{x} | \underline{\theta}) = \exp \left[\left[\sum_{i=1}^n (x_i \underline{z}_i) \right]' \underline{\theta} + \left[\sum_{i=1}^n (-\log (1 + e^{\underline{z}_i' \underline{\theta}})) \right] \right]$$

which is in the canonical form of the exponential family

$$H(\underline{\theta}) = \sum_{i=1}^n (-\log (1 + e^{\underline{z}_i' \underline{\theta}}))$$

and

$$\underline{S} = \sum_{i=1}^n (x_i \underline{z}_i) = \sum_{i=1}^n s_i$$

($G(\underline{\theta}) = 0$). We will need the true expectation of \underline{S} :

$$E_f(\underline{S}) = \left[\sum_{i=1}^n (\mu_i \underline{z}_i) \right].$$

Also, from (7.64), $H(\underline{\theta}) = \sum_{i=1}^n \log(1 - p_i)$ is an equivalent

Two key quantities we need are

$$I(\underline{\theta}) = -\frac{\partial^2 H(\underline{\theta})}{\partial \underline{\theta}^2}$$

(see 7.55) and $-\partial H(\underline{\theta}) / \partial \underline{\theta} = E_f(\underline{S})$. Some straightforward algebra leads to the results

$$-\frac{\partial H(\underline{\theta})}{\partial \underline{\theta}} = \sum_{i=1}^n p_i \underline{z}_i$$

and

$$I(\underline{\theta}) = -\frac{\partial^2 H(\underline{\theta})}{\partial \underline{\theta}^2} = \sum_{i=1}^n p_i q_i \underline{z}_i \underline{z}_i'.$$

with full diagonal element $p_i q_i$, and an $n \times K$ matrix Z with columns \underline{z}_i' . Then

$$\begin{aligned} E_f(\underline{S}) &= Z' \underline{\mu}, \\ -\frac{\partial H(\underline{\theta})}{\partial \underline{\theta}} &= Z' \underline{P}(\underline{\theta}), \end{aligned}$$

and

$$I(\underline{\theta}) = Z' V_p Z.$$

The MLE $\hat{\underline{\theta}}$ is found by setting (7.65) to \underline{S} and solving the resulting equations for $\underline{\theta}$, hence solving $\underline{S} = Z' \underline{P}(\hat{\underline{\theta}})$. The true parameter value that applies here, given truth $\underline{\mu}$ and the model, is found by solving the same equations but with \underline{S} replaced by its true expectation, hence

$$Z' \underline{\mu} = Z' \underline{P}(\underline{\theta}_o),$$

or

$$Z'(\underline{\mu} - \underline{P}(\underline{\theta}_o)) = \underline{0}.$$

In partly nonmatrix notation, we solve

$$\sum_{i=1}^n (\mu_i - p_i(\underline{\theta}_o)) \underline{z}_i = \underline{0}.$$

If truth $\underline{\mu}$ is not given exactly by the assumed model \underline{P} , then $\underline{\mu} = \underline{P}(\underline{\theta}_o)$ will not hold even though the above equation has a solution in $\underline{\theta}_o$, just as the MLE equations will have a unique solution in $\underline{\theta}_o$.

To proceed we also need to know the general formula for the variance of the MLE. We have

$$J(\underline{\theta}_o) = E_f \left[[\underline{S} - Z' \underline{P}(\underline{\theta}_o)] [\underline{S} - Z' \underline{P}(\underline{\theta}_o)]' \right]$$

In partly nonmatrix form this formula is

$$\begin{aligned} J(\underline{\theta}_o) &= E_f \left[\sum_{i=1}^n (\underline{s}_i - p_i \underline{z}_i) (\underline{s}_i - p_i \underline{z}_i)' \right] \\ &= \sum_{i=1}^n \sum_{j=1}^n E_f (\underline{s}_i - p_i \underline{z}_i) (\underline{s}_j - p_j \underline{z}_j)' \end{aligned}$$

Here, using $\underline{s}_i = x_i \underline{z}_i$ and $E_f(x_i) = \mu_i$ the above becomes

$$J(\underline{\theta}_o) = \sum_{i=1}^n \sum_{j=1}^n E_f (x_i - p_i)(x_j - p_j) [\underline{z}_i \underline{z}_j']$$

$$\begin{aligned}
&= \sum_{i=1}^n [(\mu_i(1 - \mu_i)) + (\mu_i - p_i)^2] [\underline{z}_i \underline{z}_i'] \\
&\quad + \sum_{i \neq j}^n \sum_{j=1}^n (\mu_i - p_i)(\mu_j - p_j) [\underline{z}_i \underline{z}_j'].
\end{aligned}$$

Completing the square in the trailing term above, we get

$$\begin{aligned}
J(\underline{\theta}_o) &= \sum_{i=1}^n [(\mu_i(1 - \mu_i)) + (\mu_i - p_i)^2] [\underline{z}_i \underline{z}_i'] \\
&\quad + \left[\sum_{i=1}^n (\mu_i - p_i(\underline{\theta}_o)) \underline{z}_i \right] \left[\sum_{i=1}^n (\mu_i - p_i(\underline{\theta}_o)) \underline{z}_i' \right] \\
&\quad - \sum_{i=1}^n [(\mu_i - p_i)^2] [\underline{z}_i \underline{z}_i'].
\end{aligned}$$

The middle term of the above is zero because of the equality $\sum_{i=1}^n (\mu_i - p_i(\underline{\theta}_o)) = 0$. The third term cancels with part of the first term, so we have

$$J(\underline{\theta}_o) = \sum_{i=1}^n \mu_i(1 - \mu_i) [\underline{z}_i \underline{z}_i'],$$

or in pure matrix terms,

$$J(\underline{\theta}_o) = \underline{Z}' V_\mu \underline{Z}.$$

Here, V_μ is an $n \times n$ diagonal matrix with i th diagonal element $\mu_i(1 - \mu_i)$. Contrast (7.67) to $I(\underline{\theta}_o) = \underline{Z}' V_p \underline{Z}$.

It is easy, but not very informative, now to write

$$\begin{aligned}
\text{tr} [J(\underline{\theta}_o) [I(\underline{\theta}_o)]^{-1}] &= \text{tr} [(Z' V_\mu Z) (Z' V_p Z)^{-1}] \\
&= K + \text{tr} [(Z' (V_\mu - V_p) Z) (Z' V_p Z)^{-1}].
\end{aligned}$$

The above makes it easier to realize that the trace term is $\mu_i(1 - \mu_i) = p_i(\underline{\theta}_o)(1 - p_i(\underline{\theta}_o))$. However, these equalities are only true if the model is correct. If not, yet we can still get $\text{tr} [J(\underline{\theta}_o) [I(\underline{\theta}_o)]^{-1}] = K$; hence this last term is a measure of the discrepancy between a model that does not match truth, i.e., where $g \subset f$ v.

The above results are totally general, so they apply to w replicate observations are taken at each of $j = 1, \dots, r$. The total sample size is then $n = r * w$, but we will have w values of μ_j to specify for truth and only r values of p_j for any model. Hence, to gain some insights here we use $\log [p_j / (1 - p_j)] = a + bj$ for $j = 1, \dots, r$, with w replicates at each j . Thus $K = 2$, $\underline{\theta} = (a, b)'$, and $\underline{z}_j = (1, j)'$. In

j . However, the results so derived apply reasonably well. $n = r * w$ would be “large,” say 100 or more (given $K =$ the trace term under this use of a simple logistic regression) only specify a set of μ_1, \dots, μ_r , solve

$$\sum_{i=1}^r (\mu_i - (a_o + b_o i)) \underline{z}_i = 0$$

for $\underline{\theta}_o = (a_o, b_o)'$, and compute $J(\underline{\theta}_o)$ (7.67), $I(\underline{\theta}_o)$, $\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}]$ (also denoted by “bias”). In doing this, if $\underline{\mu}$ that were near to fitting the logistic structural model, or a $\underline{p}(\underline{\theta})$ vector that fit the model, then perturbing some (or all) of $\underline{\mu}$ by starting with $\mu_j = j/(r + 1)$, which is not a logistic model but is not too far from fitting such a structural model

With at most modest deviation of truth from any actual regression model structure we found that the trace term value varied (between about 1.8 and 2, sometimes going a little above 2) in the case of truth being the simple linear model ($\mu_j = j/(r + 1)$) varied monotonically from 1.98 at $r = 5$, to 1.91 at $r = 10$. Some results for $r = 10$, based on truth being perturbed from the logistic model $\text{logit}(p_j) = 3.0 - 0.5j$. The first line of Table 7.2 shows μ_j computed from this model (scaled by 1,000).

What one can see in Table 7.2 (and other computations of this type) is that the true μ_j have to be here a very poor approximation to the simple logistic model before the trace term deviates much from 2. If the data seem at all well fit by a logistic model, then the trace (as opposed to any attempted estimation of the trace) seems to be a very important here because V_μ , hence $J(\underline{\theta}_o)$, can be estimated all unless there is replication at each \underline{z}_i , and there would not be such replication; this condition rarely occurs with logistic models.

Formula (7.68) was corroborated by direct Monte Carlo estimation of the target bias (trace term) for a few cases in Table 7.2. The most direct, and hence most direct, way to do this is to evaluate using

$$\text{bias} = E_{\underline{x}} E_{\hat{\underline{\theta}}(\underline{y})} [\log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y})))] - E_{\underline{x}} [\log(g(\underline{x} | \underline{\theta}_o))]$$

Hence for one Monte Carlo replicate (generating *iid* \underline{x} and \underline{y})

$$\widehat{\text{bias}} = \log(g(\underline{x} | \hat{\underline{\theta}}(\underline{y}))) - \log(g(\underline{x} | \underline{\theta}_o))$$

Averaged over many samples (m), if large sample size n is used, the bias will equal $\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}]$.

For many models the first term on the right-hand side of (7.68) can be evaluated in \underline{x} , so we can analytically take the expectation with re-

$S = 0.5J$, values of μ_j are shown, scaled by 1,000, case one (i.e., the logistic model, but none of the other cases are a perfect fit to the data). The results are reasonably applicable if w is at least 10 or 20.

μ_1	μ_2	μ_3	μ_4	μ_5	μ_6	μ_7	μ_8	μ_9
924	881	818	731	622	500	378	269	182
900	900	900	900	622	500	378	269	182
924	881	970	757	725	530	410	231	186
900	944	898	773	666	522	349	264	176
864	993	971	821	564	514	366	236	150
894	990	870	826	583	457	372	252	216
924	720	818	796	749	632	492	188	260
924	881	768	831	622	400	378	269	282
924	952	650	838	638	448	278	189	204
874	874	874	558	558	558	558	190	190
924	881	818	731	400	600	378	269	182
924	881	818	731	300	700	378	269	182
924	881	818	731	622	500	378	269	182
924	881	568	831	622	100	378	269	432
924	881	818	731	200	800	378	269	182
674	981	818	731	622	500	128	669	182
924	881	818	731	622	500	900	900	900

logistic example we thus get, expressed in basic form,

$$\widehat{\text{bias}} = \sum_{i=1}^n \text{E}_{\underline{y}} \Big[n_i \mu_i \log(p_i(\hat{\underline{\theta}})) + n_i (1 - \mu_i) \log(1 - p_i(\hat{\underline{\theta}})) \Big] \\ - \sum_{i=1}^n \text{E}_{\underline{y}} \Big[y_i \log(p_i(\hat{\underline{\theta}})) + (n_i - y_i) \log(1 - p_i(\hat{\underline{\theta}})) \Big]$$

Here the MLE $\hat{\underline{\theta}}$ is based on data \underline{y} .

For direct Monte Carlo evaluation we used replicates noted above with the same number of replicates w for each \underline{y} . For each \underline{y} we generated a large number m of independent samples from the model, fit the model-based MLE to each sample, compute the bias by sample, and got its average and empirical standard error.

All this is quite obvious; where we are going here is that Monte Carlo evaluation is poor in the sense of needing a huge number of replicates. The problem is that as w (i.e., n) increases, the number of replicates needed to get a small standard error (like 0.005) on the estimate bias increases. The variance of $\widehat{\text{bias}}$, for one sample, increases with increasing w . That variance can be quite large. For a large sample size n (like 1000) the variance can be quite large.

example, for the case in Table 7.2 where the trace is computed with $w = 100$ (hence sample size $n = 1,000$ Bernoulli trials) for a set of 10,000 ($= m$) Monte Carlo samples we got average bias = 1.982 with an estimated standard error of 0.135. Other runs verified that for one million Monte Carlo samples to get a standard error of about 0.014 with $w = 100$ in this example. But we might apply the trace formula to apply exactly; for $w = 1,000$ and 10,000 samples we got the average bias = 1.982 with $\hat{se} = 0.42$. This is the reverse of what we expect; i.e., we expect to get increased bias (the same number, m , of Monte Carlo samples) as sample size n increases. This reverse phenomenon occurs here because the expected difference between the two terms involved in direct computation of bias (i.e., (7.70)) is constant for sample size n , but the variance of each of those two likelihood terms is proportional to n and the two terms are not highly correlated. So as sample size n increases, the precision of the estimated bias, given by the Monte Carlo samples (m), actually decreases. So to evaluate the bias by a brute-force approach, the adequacy of the trace term approaches zero. For sample sizes, it takes a huge number of Monte Carlo samples to get a small standard error.

With models that are in the exponential family there is a way to do exact Monte Carlo evaluation of the bias that must be computed for the maximized log-likelihood for exact K-L based models. The result (7.54) is an exact result for any sample size:

$$\text{bias} = \text{tr} \left[\text{COV}(\hat{\underline{\theta}}, \underline{S}) \right].$$

While $\underline{\theta}$ and \underline{S} are only for the canonical form of the model, the result can apply for any parametrization of the assumed model because of the result to 1-to-1 transformations of $\underline{\theta}$ (see end of Section 7.2). An alternative Monte Carlo evaluation is simply to take for each parameter the computed MLE and minimal sufficient statistic and, from a set of Monte Carlo samples of size m , estimate the covariances $\text{cov}(\hat{\theta}_i, S_i)$, $i = 1, \dots, K$, for these K estimates. The result is bias, and this approach is simple. For the same case in Table 7.2 (i.e., trace = 1.774), using $w = 100$ (hence $n = 1,000$) and 10,000 Monte Carlo samples we got average bias = 1.982 with $\hat{se} = 0.019$ using the covariance approach. Based on this result there was a clear suggestion that $w = 100$ was not quite adequate to get the trace (7.68) to apply reliably to three digits (it was then recommended to use $w = 1,000$ and 10,000 Monte Carlo samples we got average bias = 1.982 with $\hat{se} = 0.017$. This result held up on more study: (7.68) seen to be accurate for $w = 1,000$ (which here meant $n = 10,000$).

As another example consider the last case in Table 7.2, where the trace is computed with $w = 100$ (hence sample size $n = 1,000$ Bernoulli trials) for a set of 10,000 ($= m$) Monte Carlo samples we got average bias = 1.982 with an estimated standard error of 0.135. Other runs verified that for one million Monte Carlo samples to get a standard error of about 0.014 with $w = 100$ in this example. But we might apply the trace formula to apply exactly; for $w = 1,000$ and 10,000 samples we got the average bias = 1.982 with $\hat{se} = 0.42$. This is the reverse of what we expect; i.e., we expect to get increased bias (the same number, m , of Monte Carlo samples) as sample size n increases. This reverse phenomenon occurs here because the expected difference between the two terms involved in direct computation of bias (i.e., (7.70)) is constant for sample size n , but the variance of each of those two likelihood terms is proportional to n and the two terms are not highly correlated. So as sample size n increases, the precision of the estimated bias, given by the Monte Carlo samples (m), actually decreases. So to evaluate the bias by a brute-force approach, the adequacy of the trace term approaches zero. For sample sizes, it takes a huge number of Monte Carlo samples to get a small standard error.

exponential family; Monte Carlo or bootstrap evaluation of the trace term should be based on (7.54).

It is worth noting a basis for the estimated standard error $\widehat{\text{tr}}[\text{COV}(\hat{\theta}, \underline{S})]$. For the point estimate, use all the simulated means; then for component i ,

$$\widehat{\text{cov}}(\hat{\theta}_i, S_i) = \frac{\sum_{j=1}^m (\hat{\theta}_{i,j} - \bar{\hat{\theta}}_i)(S_{i,j} - \bar{S}_i)}{m - 1}$$

and

$$\widehat{\text{bias}} = \sum_{i=1}^K \widehat{\text{cov}}(\hat{\theta}_i, S_i).$$

However, to estimate the standard error we must partition the simulated means into 25 equal-sized subsets (for $m = 10,000$ then $m/25 = 400$). Compute by the above formulas $\widehat{\text{bias}}_s$ for each subset s . Then estimate the standard error of bias from these 25 independent estimates. The estimates will almost equal bias, but will not be equal due to nonlinearities.

The standard error of bias from this covariance approach is a function of data sample size n because of how the product $(\hat{\theta}_i - \theta_{o,i})(S_i - E_f(S_i))$ is computed. It suffices to consider the product $(\hat{\theta}_i - \theta_{o,i})(S_i - E_f(S_i))$ for component i . This product has variance virtually independent of n because $\hat{\theta}_i$ converges (in n) at rate proportional to $1/\sqrt{n}$, while the standard error of S_i is at rate proportional to \sqrt{n} . As a result, the standard error of the bias-based bias (hence trace) estimator is almost independent of n . This is much better behavior (as a function of n) than the standard error estimator of bias based directly on the likelihood function (i.e., (7.69)) also requires more calculations beyond first order.

7.6.3 Evaluation Under Multinomially Distributed Data

We here assume that we have count data n_1, \dots, n_r that sum to n . Truth is the multinomial distribution $\text{mult}(n, \mu_1, \dots, \mu_r)$ with probabilities μ_i summing to 1, and $0 < \mu_i < 1$. To know truth we only need to know the true μ_i (assuming that the counts are multinomially distributed; they could have overdispersion, which violates the multinomial assumption). We might totally fail to know how these true probabilities relate to any explanatory variables, or what would happen if the probabilities were defined in some other way. Thus, deeper truth may exist, but it is irrelevant to model selection purposes once we have observed the data to a particular multinomial setting.

not to be considered, since then the fitted model matches the data (being a perfect match to the data). The theory in Section 7.5 is used; note that here

$$\log(g(\underline{n} | \underline{\theta})) = \sum_{i=1}^r n_i \log(p_i(\underline{\theta})).$$

First, $\underline{\theta}_o$ is determined as the solution to (7.5), which here

$$\sum_{i=1}^r \frac{\mu_i}{p_i(\underline{\theta}_o)} \frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} = \underline{0}.$$

In (7.71) if we replace μ_i by n_i , we have the likelihood can treat the μ_i as data and find $\underline{\theta}_o$ by MLE methods. Equivalently, we can find $\underline{\theta}_o$ by MLE when the data are replaced by their true expected values.

Second, applying (7.7) we directly get

$$\begin{aligned} I(\underline{\theta}_o) = n & \left[\sum_{i=1}^r \frac{\mu_i}{[p_i(\underline{\theta}_o)]^2} \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right) \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right)^T \right. \\ & \left. - n \left[\sum_{i=1}^r \frac{\mu_i}{p_i(\underline{\theta}_o)} \left(\frac{\partial^2 p_i(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right) \right] \right]. \end{aligned}$$

Finally, applying the definition in (7.12), we have

$$J(\underline{\theta}_o) = E_f \left[\sum_{i=1}^r \frac{n_i}{p_i(\underline{\theta}_o)} \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right) \right] \left[\sum_{i=1}^r \frac{n_i}{p_i(\underline{\theta}_o)} \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right)^T \right]$$

The evaluation of $J(\underline{\theta}_o)$ does take some algebra and knowledge of the multinomial distribution, but it is mostly a straightforward exercise. The result is:

$$J(\underline{\theta}_o) = n \left[\sum_{i=1}^r \frac{\mu_i}{[p_i(\underline{\theta}_o)]^2} \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right) \left(\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} \right)^T \right]$$

Define the matrix A as

$$A = n \left[\sum_{i=1}^r \frac{\mu_i}{p_i(\underline{\theta}_o)} \left(\frac{\partial^2 p_i(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right) \right],$$

and we have $I(\underline{\theta}_o) = J(\underline{\theta}_o) - A$. Furthermore, if the true distribution is $\mu_i = p_i(\underline{\theta}_o)$, and A reduces to the null matrix; hence then

Using these results we can write

$$\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] = \text{tr} [(J(\underline{\theta}_o) - A + A)[I(\underline{\theta}_o)]^{-1}]$$

For a long time a nagging question for us was whether K always be either $> K$ or $< K$ when the model did not yet the model is logically known to be simpler than true (in some general sense). Stated differently, if the Kullback-Leibler divergence is positive, i.e., $K-L = I(f, g) > 0$, then must $\text{tr}[J(\underline{\theta}_o)]$ maybe $< K$ always occur when the model is some form of approximation (hence the model can be said to approximate, but not equal to the truth) is no, as was indicated by the logistic regression example. However, a more convincing answer is given here: The trace of $J(\underline{\theta}_o)$ is $> K$ or $< K$ and there need be no consistency as to which way the question also explored below is, If $\text{tr}[J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] = K$. That answer is also no.

Because the cell probabilities sum to 1, the sum of the second partials is the null matrix, O , of all zeros (the vector of first partials is a null vector). Therefore, an equivalent expression for A is

$$A = n \left[\sum_{i=1}^r \frac{\mu_i - p_i(\underline{\theta}_o)}{p_i(\underline{\theta}_o)} \left(\frac{\partial^2 p_i(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right) \right]$$

The weights in this linear combination of second partials must be either identically zero (hence K-L is 0), or some are positive and some are negative. This would suggest that A might not always be positive or negative, unless the second partials are very strangely related to the truth. But a more detailed case is needed to get an example, so we use $K = 1$ for an example, such as by using a binomial model.

Let us further assume that the data arise from n independent integer random variables, y , taking values 0 to $r - 1$. Then the frequency counts n_i of times $y = i - 1$. A very simple model is to assume that this underlying random variable is a binomial random variable. This corresponds to imposing multinomial cells, without loss of generality, and thus the probability is

$$p_i(\theta) = \binom{r-1}{i-1} \theta^{i-1} (1-\theta)^{r-i}, \quad i = 1, \dots, r$$

Thus we have, as our model, an assumed underlying binomial random variable, $y \sim \text{bin}(r-1, \theta)$ and a random sample of size n of this variable. In fact, y ($= 0, 1, \dots, r-1$) has the distribution given by the true distribution. We will need the functions below, involving partial derivatives:

$$P1_i = \frac{1}{p_i(\theta)} \left(\frac{\partial p_i(\theta)}{\partial \theta} \right) = \frac{(i-1) - \theta(r-1)}{\theta(1-\theta)},$$

$$- \left[\frac{(i-1)(1-2\theta)}{(\theta(1-\theta))} \right]$$

We solve (7.71), which is $\sum \mu_i P 1_i = 0$, to find θ_o same process as finding an MLE (again, the only tricky aspect assumed here):

$$\theta_o = \frac{\sum_{i=1}^r \mu_i (i-1)}{r-1}.$$

This θ_o is the true expected value of $y/(r-1)$ regardless of θ . We compute (7.72) as $n \sum \mu_i (P 1_i)^2$:

$$J(\theta_o) = n \left[\sum_{i=1}^r \mu_i \left(\frac{(i-1) - \theta_o(r-1)}{\theta_o(1-\theta_o)} \right)^2 \right] \equiv n E_f$$

We find matrix A as $n \sum \mu_i P 2_i$:

$$A = J(\theta_o) - \frac{n(r-1)}{\theta_o(1-\theta_o)};$$

hence

$$I(\theta_o) = \frac{n(r-1)}{\theta_o(1-\theta_o)}.$$

It is now easy to find the trace:

$$\text{tr} [J(\theta_o)[I(\theta_o)]^{-1}] = \sum_{i=1}^r \mu_i \frac{[(i-1) - \theta_o(r-1)]}{(r-1)\theta_o(1-\theta_o)}$$

For the case of $\mu_i = p_i(\theta_o)$, then (7.73) is 1 (this can be seen) hence using $\mu_i \equiv p_i(\theta_o) + (\mu_i - p_i(\theta_o))$ in (7.73) we obtain

$$\text{tr} [J(\theta_o)[I(\theta_o)]^{-1}] = 1 + \sum_{i=1}^r (\mu_i - p_i(\theta_o)) \frac{[(i-1) - \theta_o(r-1)]}{(r-1)\theta_o(1-\theta_o)}$$

whereupon it should be essentially obvious that the term in the sum can be either positive or negative. However, we will give numerical examples mostly for $r = 3$ because this is the smallest r we can use here, and small r is desirable when we need to display truth.

Our model is thus $\text{bin}(2, \theta)$; hence $p_1 = (1-\theta)^2$, $p_2 = 2\theta(1-\theta)$, $p_3 = \theta^2$. The approach is to specify the μ_i and compute $J(\theta_o)$ and from (7.73), for $r = 3$,

$$\text{tr} [J(\theta_o)[I(\theta_o)]^{-1}] = \frac{4\mu_1(\theta_o)^2 + \mu_2(1-2\theta_o)^2 + 4\mu_3\theta_o^2}{2\theta_o(1-\theta_o)}$$

cases of truth (the μ_i) in relationship to the K-L best approximating binomial model. Case one exactly fits a binomial model. Cases two and three are approximated reasonably by even the K-L best approximating binomial model with $\theta_o = 0.5$. In what is below, “Trace” means the value computed for the trace term $\theta_o = 0.5$, and K-L is the Kullback–Leibler information content of the truth and the best approximating binomial model (“Bias-Correction” is below):

μ_1	μ_2	μ_3	Trace	K-L	Bias-L
0.25	0.50	0.25	1.000	0.000	1.000
0.05	0.90	0.05	0.200	0.368	0.196
0.45	0.10	0.45	1.800	0.368	1.804

(we note that in this situation the trace term (7.73) seems to be off by a factor of 2). Clearly, this bias-correction trace term can be either greater or less than 1 when the model does not match truth. This is because the variance of y can be either larger or smaller than the variance implied by the K-L best-fitting binomial model.

We build on this example by doing an exact Monte Carlo simulation of the expected log-likelihood and the K-L–based target model. We use this to verify the asymptotic derivation of the bias as being the difference between the above, “Bias-MC” denotes the results (accurate to two decimal places) for sample size $n = 200$, from one million Monte Carlo simulations. The bias that the trace term measures based on asymptotic theory is shown.

For a truth that cannot be well approximated here by a binomial model, it is clear that the trace (equation 7.74) can be far from 1. Rather than using this example for models that are arbitrarily poor (like cases two and three), we should consider models that are closer to truth, because the term $-2 \log(\mathcal{L})$ will prevent the selection of really poor models. For those models a choice between the use of K or trace is irrelevant; the models has some good candidates.

So we looked at one set of cases where a binomial model is not the best wrong to use. We chose a θ , generated $p_1 = (1 - \theta)^2$, $p_2 = 2\theta(1 - \theta)$, $p_3 = \theta^2$, then perturbed these cell probabilities to get a model that is close to a binomial model by setting $\mu_i = p_i + \epsilon_i$, where $\epsilon_i \sim \text{uniform}(-h, h)$. Inadmissible sets of μ_i were not generated. The trace term (7.74) was evaluated; thus, this is not a Monte Carlo simulation. We use Monte Carlo methods only as a convenience in generating models μ_i that are close to a binomial model.

For $\theta = 0.5$ and $h = 0.1$ (and 1,000 generated sets of truth and model) the following results for the trace given by (7.74): min = 0.196 and mean = 0.996. These results support practical use of the trace term $\text{tr}[J(\theta_o)[I(\theta_o)]^{-1}]$. However, it is fair to ask about estimating the trace term

general, the true trace value is known for these cases, the mean and standard deviation of the $\hat{\text{tr}}$ values (7.75) are given based on one million samples.

μ_1	μ_2	μ_3	n	Trace	Mean	SD
0.25	0.50	0.25	50	1.000	0.990	0.005
0.25	0.50	0.25	100	1.000	0.995	0.004
0.25	0.50	0.25	200	1.000	0.997	0.003
0.04	0.32	0.64	200	1.000	0.997	0.003
0.20	0.55	0.25	50	0.897	0.887	0.005
0.20	0.60	0.20	50	0.800	0.790	0.005
0.30	0.40	0.30	50	1.200	1.190	0.005
0.30	0.45	0.25	50	1.098	1.088	0.005
0.05	0.90	0.05	200	0.200	0.199	0.003
0.45	0.10	0.45	200	1.800	1.799	0.003

here (hence using TIC), as can be done by plugging $\hat{\theta}_o$ (7.74); after simplification,

$$\widehat{\text{tr}}\left[J(\theta_o)[I(\theta_o)]^{-1}\right] = \frac{4\hat{\mu}_1(\hat{\theta}_o)^2 + \hat{\mu}_2(1 - 2\hat{\theta}_o)^2 + 4\hat{\mu}_3}{2\hat{\theta}_o(1 - \hat{\theta}_o)}$$

A small Monte Carlo evaluation of this estimator was done to see if it was badly biased or highly variable. Variables in this study were the true μ values, the sample size n . Results, given in Table 7.3 based on one million simulations, are the theoretical trace value, and the mean and standard deviation of the $\hat{\text{tr}}$ values evaluated by simulation, accurate to two decimal places. If the true trace is 1, then the binomial distribution is truth; otherwise, it is not.

From Table 7.3 it appears that the trace estimator has good properties. It is reasonable to consider using TIC rather than AIC; at least, for the binomial model, of the two seems worth doing here. For the sets of true μ values, we compared AIC to TIC for the binomial model ($K = 2$) and the parameter-saturated general model ($K = 3$). The corresponding maximized likelihoods be \mathcal{L}_R and \mathcal{L}_G . Here

$$\begin{aligned} \text{AIC}_R &= -2\log(\mathcal{L}_R) + 2, \\ \text{AIC}_G &= -2\log(\mathcal{L}_G) + 4, \\ \text{TIC}_R &= -2\log(\mathcal{L}_R) + 2\widehat{\text{tr}}\left[J(\theta_o)[I(\theta_o)]^{-1}\right], \\ \text{TIC}_G &= -2\log(\mathcal{L}_G) + 4, \end{aligned}$$

where $\widehat{\text{tr}}\left[J(\theta_o)[I(\theta_o)]^{-1}\right]$ is given by (7.75). Because there are many methods here, and because we want to keep matters simple, we just use the methods based on how often they selected the same model. Results are given in Table 7.3. They selected the correct data-generating model. Results are given in Table 7.3. Monte Carlo samples, which suffices here to get standard deviations.

generating model varied (R for the binomial, G for the multinomial) on 10,000 samples; column AIC (or TIC) denotes the proportion of times AIC (or TIC) selected the correct data-generating model; column “Match” denotes the proportion of times selected the same model whether or not it was the data-generating model (for details).

μ_1	μ_2	μ_3	n	Truth	AIC	TIC
0.25	0.50	0.25	50	R	0.83	0.84
0.25	0.50	0.25	100	R	0.85	0.84
0.25	0.50	0.25	200	R	0.84	0.84
0.04	0.32	0.64	200	R	0.84	0.83
0.20	0.55	0.25	50	G	0.69	0.76
0.20	0.60	0.20	50	G	0.57	0.53
0.30	0.40	0.30	50	G	0.46	0.53
0.30	0.45	0.25	50	G	0.24	0.28
0.05	0.90	0.05	200	G	1.00	1.00
0.45	0.10	0.45	200	G	1.00	1.00
means					0.73	0.75

for estimated proportions. “Truth” denotes the correct generating model; “AIC” and “TIC” columns denote the proportion of samples for which AIC or TIC selected the correct model. “Match” denotes the proportion of samples for which both methods selected the same model, regardless of whether the correct model was selected.

We looked at many more results than are given above. In general, AIC and TIC performed similarly. In this limited context; there was then no change from the previous results. The obvious conclusion: no meaningful difference in performance between AIC and TIC. More study is surely warranted; this limited study was in the spirit that maybe something dramatic would result. It did not. The obvious conclusion (based on all the considerations we have discussed) is that AIC is just those of this section), it seems that simplicity strongly favors AIC over TIC.

We return to an interesting theoretical question posed above. Suppose $f = g$. If $f = g$, then $\text{tr} [J(\theta_o)[I(\theta_o)]^{-1}] = K$. However, if $\text{tr} [J(\theta_o)] = K$, does this mean $f = g$?

A counterexample shows that the assertion is false; hence, the answer is no. wherein truth is more complex than the models used for approximation. AIC is appropriate to use (as opposed to TIC, which would tend to overfit) just estimating K).

For this multinomial context we have shown that $I(\underline{\theta}_o) = J(\underline{\theta}_o)$ (i.e. the same) for any value of r . So if $A = O$ (i.e. the origin), then $I(\underline{\theta}_o) = J(\underline{\theta}_o)$ and $\text{tr} [J(\theta_o)[I(\theta_o)]^{-1}] = K$ regardless of the value of r .

$$A = n \left[\sum_{i=1}^r \frac{p_i}{p_i(\underline{\theta}_o)} \left(-\frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}^2} \right) \right];$$

therefore, if all second partial derivatives of the model are zero, we do get $A = O$. This will occur for any linear probabilities; that is, $p_i(\underline{\theta}) = \underline{x}_i' \underline{\theta}$ for a set of known vectors such models are discouraged because they can generate values out of range.

As an example we revert to the case of $r = 3$ and use $\mu_1 = \mu_3 = \theta/2$ and $\mu_2 = 1 - \theta$; so $\log(g(n|\underline{\theta})) = (n/2) \log(\theta/2) + n_3 \log(1 - \theta)$. Here, $\theta_o = 1 - \mu_2$. Upon computing $I(\underline{\theta}_o)$ from the basic definitions, we do in fact get $I(\underline{\theta}_o) = J(\underline{\theta}_o) = n/[\theta_o]$ of the values of the μ_i . This means that here is a situation where AIC rather than TIC is the correct selection procedure even if the model does differ from truth (i.e., $f \subseteq g$ is not true, yet this conclusion is cited as always required for the theoretical validity of AIC).

For the case of general r and the binomial model we will investigate this trace term and AIC versus TIC. But even if we cannot derive any insights from it directly, and numerical methods are used, all we really need are usable computational formulas to compute the estimate $\text{tr}[J(\theta_o)[I(\theta_o)]^{-1}]$. We can get the needed formulas for the model for multinomial data. First, we can find $\hat{\underline{\theta}}_o$ by solving

$$\sum_{i=1}^r \frac{n_i/n}{p_i(\underline{\theta}_o)} \frac{\partial p_i(\underline{\theta}_o)}{\partial \underline{\theta}} = \underline{0},$$

which we do anyway, since this is just our MLE of $\underline{\theta}$ for the model. We do have to compute the set of first and second partial derivatives of the model cell structures evaluated at the MLE, but even if we do this numerically. Thus we can get, hence use and explore, TIC

$$\begin{aligned} \hat{I}(\underline{\theta}_o) &= n \left[\sum_{i=1}^r \frac{n_i/n}{[p_i(\hat{\underline{\theta}}_o)]^2} \left(\frac{\partial p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}} \right) \left(\frac{\partial p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}} \right)' \right. \\ &\quad \left. - n \left[\sum_{i=1}^r \frac{n_i/n}{p_i(\hat{\underline{\theta}}_o)} \left(\frac{\partial^2 p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}^2} \right) \right] \right], \\ \hat{J}(\underline{\theta}_o) &= n \left[\sum_{i=1}^r \frac{n_i/n}{[p_i(\hat{\underline{\theta}}_o)]^2} \left(\frac{\partial p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}} \right) \left(\frac{\partial p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}} \right)' \right. \\ &\quad \left. - \sum_{i=1}^r \frac{n_i/n}{p_i(\hat{\underline{\theta}}_o)} \left(\frac{\partial^2 p_i(\hat{\underline{\theta}}_o)}{\partial \underline{\theta}^2} \right) \right]. \end{aligned}$$

Clearly, we can also compute theoretical values of these quantities for the postulated truth and model. Such studies would be informative for the intention of this book.

7.6.5 Evaluation for Fixed-Effects Normality-Based Models

The fixed-effects linear model based on n iid normally distributed errors is so common that it seems almost mandatory that we consider it. We consider it under this model for some tractable “truth.” The model is $\underline{y} = X\underline{\beta} + \underline{\epsilon}$, where $\underline{\epsilon} \sim \text{multivariate-normal}(\underline{0}, \sigma^2 I)$, and without loss of generality, the design matrix X is assumed of full rank. Truth has a structural component $\underline{\beta}$ (which can be estimated by \underline{Y}) and a stochastic component $\underline{\epsilon}$. If $\underline{\epsilon}$ is distributed in some unknown way, the properties of which we do not know without strong assumptions. If the model is truth, then $\underline{\mu} = X\underline{\beta}$ and distributional assumptions of the model are not needed. In reality, the ϵ_i may not be independent, may not be iid, and may not be normally distributed. In fact, they may not even be zero; some or all ϵ_i are zero with probability 1. In this latter case, the model is not a formula; that is, there is some sufficiently complex computing formula (e.g., a formula with many covariates) such that if we knew that formula, we could predict \underline{Y} with certainty (measurement error would be zero before this level of model accuracy was reached). Hence, for real data, we cannot, for cases of real data, evaluate K-L-based model accuracy for continuous random variables.

We can, however, derive informative results under general assumptions that are better approximations to reality (by assumption, the model can be used for data analysis. Therefore, we assume here that the model is a good approximation to reality). We assume $\underline{\epsilon} \sim \text{multivariate-normal}(\underline{0}, \tau^2 I)$, where τ^2 may be zero. We can even drop the full distributional assumption, as we will see, because the relevant evaluations require only the first four moments of the distribution. More generally, results could be gotten under the assumption $\underline{\epsilon} \sim \text{multivariate-normal}(\underline{0}, \Sigma)$ for given Σ , but the more general case will suffice.

We first need basic notation and results: $\underline{\theta}$ denotes the K parameters for $\underline{\beta}$ a $(K - 1) \times 1$ vector of the structural parameters. \underline{Y} is the parameter to estimate. The model pdf for the data is

$$g(\underline{y} | \underline{\theta}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left[-\frac{1}{2\sigma^2} (\underline{Y} - X\underline{\beta})' (\underline{Y} - X\underline{\beta}) \right]$$

and we take, without loss of generality,

$$\log(g(\underline{y} | \underline{\theta})) = -\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} (\underline{Y} - X\underline{\beta})' (\underline{Y} - X\underline{\beta})$$

$$\frac{\partial \log(g(\underline{y} | \underline{\theta}))}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{1}{2(\sigma^2)^2}(\underline{Y} - X\underline{\beta})'(\underline{Y} - X\underline{\beta})$$

As per theory, we take the expectations of (7.76) and (7.77) what is here just assumed truth, f (evaluation under absolute truth, impossible) to get the equation

$$\frac{1}{\sigma_o^2} X'(\underline{\mu} - X\underline{\beta}_o) = \underline{0}$$

from (7.76), and then from (7.77) we derive

$$\begin{aligned} & -\frac{n}{2\sigma_o^2} + \frac{1}{2(\sigma_o^2)^2} E_f(\underline{Y} - \underline{\mu} + \underline{\mu} - X\underline{\beta}_o)'(\underline{Y} - \underline{\mu} + \underline{\mu} - X\underline{\beta}_o) \\ & = -\frac{n}{2\sigma_o^2} + \frac{1}{2(\sigma_o^2)^2} \left[E_f(\underline{\epsilon}'\underline{\epsilon}) + (\underline{\mu} - X\underline{\beta}_o)'(\underline{\mu} - X\underline{\beta}_o) \right] \\ & = -\frac{n}{2\sigma_o^2} + \frac{1}{2(\sigma_o^2)^2} \left[n\tau^2 + \|\underline{\mu} - X\underline{\beta}_o\|^2 \right] = 0 \end{aligned}$$

It is now a simple matter to find

$$\begin{aligned} \underline{\beta}_o &= (X'X)^{-1} X' \underline{\mu}, \\ \sigma_o^2 &= \tau^2 + \frac{\|\underline{\mu} - X\underline{\beta}_o\|^2}{n}. \end{aligned}$$

These parameter values define the vector $\underline{\theta}_o = (\underline{\beta}_o', \sigma_o^2)'$. From (7.76) and (7.77) we get the expected second mixed partial derivatives of the log-likelihood function, that lack-of-fit variation, from the assumed structural model and of residual (unexplained) variation.

To find $I(\underline{\theta}_o)$ we need the expected second mixed partial derivatives of the log-likelihood function, (7.77), as follows:

$$\begin{aligned} E_f \left[-\frac{\partial^2 \log(g(\underline{y} | \underline{\theta}_o))}{\partial \underline{\beta}^2} \right] &= \frac{1}{\sigma_o^2} X'X, \\ E_f \left[-\frac{\partial^2 \log(g(\underline{y} | \underline{\theta}_o))}{\partial \sigma^2 \partial \sigma^2} \right] &= -\frac{n}{2(\sigma_o^2)^2} + \frac{1}{(\sigma_o^2)^3} n\sigma_o^2 \\ \text{and} \\ E_f \left[-\frac{\partial^2 \log(g(\underline{y} | \underline{\theta}_o))}{\partial \underline{\beta} \partial \sigma^2} \right] &= \frac{1}{\sigma_o^2} X'(\underline{\mu} - X\underline{\beta}_o) = \underline{0} \end{aligned}$$

$$I(\underline{\theta}_o) = \begin{bmatrix} \frac{1}{\sigma_o^2} X'X & \underline{0} \\ \underline{0}' & \frac{n}{2(\sigma_o^2)^2} \end{bmatrix}.$$

The evaluation of $J(\underline{\theta}_o)$ is harder, and more dependent on the distribution of $\underline{\epsilon}$. Evaluation of $I(\underline{\theta}_o)$ required only the second moment of $\underline{\epsilon}$. Evaluation of $J(\underline{\theta}_o)$ also requires third and fourth moments. Both derivations depend on the independence of the ϵ_i . The upper left $(K-1) \times (K-1)$ part of $J(\underline{\theta}_o)$ is

$$E_f \left[\frac{1}{(\sigma_o^2)^2} X'(\underline{Y} - X\underline{\beta}_o)(\underline{Y} - X\underline{\beta}_o)'X \right] = \frac{\tau}{(\sigma_o^2)^2}$$

The last $(K-1) \times 1$ column vector (of the first $K-1$ rows) is

$$E_f \left[\frac{1}{\sigma_o^2} X'(\underline{Y} - X\underline{\beta}_o) \left[-\frac{n}{2\sigma_o^2} + \frac{1}{2(\sigma_o^2)^2} (\underline{Y} - X\underline{\beta}_o)'(\underline{Y} - X\underline{\beta}_o) \right] \right]$$

Making use of $X'(\underline{\mu} - X\underline{\beta}_o) = \underline{0}$, and some algebra similar to the above to

$$E_f \left[\frac{1}{2(\sigma_o^2)^3} X'(\underline{Y} - \underline{\mu}) \left[(\underline{Y} - \underline{\mu})'(\underline{Y} - \underline{\mu}) + 2(\underline{Y} - \underline{\mu})'(\underline{Y} - X\underline{\beta}_o) \right] \right]$$

and then to

$$\begin{aligned} E_f \left[\frac{1}{2(\sigma_o^2)^3} \left[X'(\underline{Y} - \underline{\mu})[(\underline{Y} - \underline{\mu})'(\underline{Y} - \underline{\mu})] + 2X'(\underline{Y} - \underline{\mu})(\underline{Y} - X\underline{\beta}_o)'(\underline{Y} - \underline{\mu}) \right] \right] \\ = E_f \left[\frac{1}{2(\sigma_o^2)^3} \left[X'(\underline{Y} - \underline{\mu})[(\underline{Y} - \underline{\mu})'(\underline{Y} - \underline{\mu})] \right] \right] \end{aligned}$$

Now write the needed expectation in terms of the hypothesis that the ϵ_i are iid $N(0, \tau^2)$; hence

$$E_f \left[(\underline{Y} - \underline{\mu})[(\underline{Y} - \underline{\mu})'(\underline{Y} - \underline{\mu})] \right] = E_f \left[\underline{\epsilon}' \left[\sum_{i=1}^n \epsilon_i^2 \right] \right]$$

The j th element of this vector is $E_f((\epsilon_j)^3 + \sum_{i \neq j} \epsilon_j(\epsilon_i)^2)$. Because of the mutual independence of the ϵ_j , the expectation of the second term is $E_f(\epsilon_j)^3$. Because the ϵ_j are assumed to be normally distributed, their third central moment is 0. Hence, we have

$$E_f \left[\underline{\epsilon}' \left[\sum_{i=1}^n (\epsilon_i)^2 \right] \right] = \underline{0},$$

and the desired part of $J(\underline{\theta}_o)$ is $\underline{0}$.

$$J_{KK}(\underline{\theta}_o) = E_f \left[-\frac{1}{2\sigma_o^2} + \frac{1}{2(\sigma_o^2)^2} (\underline{Y} - X\underline{\beta}_o)' (\underline{Y} - X\underline{\beta}_o) \right]$$

Several straightforward steps reduce the above to

$$J_{KK}(\underline{\theta}_o) = \frac{1}{(2\sigma_o^2)^2} \left[-n^2 + \frac{1}{(\sigma_o^2)^2} E_f \left[(\underline{Y} - X\underline{\beta}_o)' (\underline{Y} - X\underline{\beta}_o) \right] \right]$$

Define the i th row vector of X as \underline{x}'_i . Then

$$\begin{aligned} E_f \left[(\underline{Y} - X\underline{\beta}_o)' (\underline{Y} - X\underline{\beta}_o) \right]^2 &= E_f \left[\sum_{i=1}^n (y_i - \underline{x}'_i \underline{\beta}_o)^2 \right]^2 \\ &= E_f \left[\sum_{i=1}^n \sum_{j=1}^n (y_i - \underline{x}'_i \underline{\beta}_o)^2 (y_j - \underline{x}'_j \underline{\beta}_o)^2 \right] \\ &= E_f \left[\sum_{i=1}^n (y_i - \underline{x}'_i \underline{\beta}_o)^4 + \sum_{i \neq j} (y_i - \underline{x}'_i \underline{\beta}_o)^2 (y_j - \underline{x}'_j \underline{\beta}_o)^2 \right] \end{aligned}$$

By virtue of mutual independence, the expectation of the above is easily found, giving

$$\begin{aligned} E_f \left[(\underline{Y} - X\underline{\beta}_o)' (\underline{Y} - X\underline{\beta}_o) \right]^2 &= E_f \left[\sum_{i=1}^n (y_i - \underline{x}'_i \underline{\beta}_o)^4 \right] + (n\sigma_o^2)^2 - \sum_{i=1}^n \left[\tau^2 + \right] \end{aligned}$$

For the case $\tau^2 = 0$ note that $\underline{Y} = \underline{\mu}$, and so the above

$$E_f \left[(\underline{Y} - X\underline{\beta}_o)' (\underline{Y} - X\underline{\beta}_o) \right]^2 = (n\sigma_o^2)^2; \text{ hence}$$

$$J_{KK}(\underline{\theta}_o) = \frac{1}{(2\sigma_o^2)^2} \left[-n^2 + \frac{1}{(\sigma_o^2)^2} (n\sigma_o^2)^2 \right]$$

It is thus clear that if $\tau^2 = 0$, then $J(\underline{\theta}_o) = O$.

The next steps are valid only if $\tau^2 > 0$. Let

$$\sqrt{\lambda_i} = \frac{\mu_i - \underline{x}'_i \underline{\beta}_o}{\tau}$$

and

$$z_i = \frac{y_i - \mu_i}{\tau}.$$

$$E_f \left[\sum_{i=1}^n (y_i - \underline{x}_i' \underline{\beta}_o) \right] = \tau \left[\sum_{i=1}^n E_f [z_i + v_i] \right]$$

The needed expectation is now easily found because it is just that of a normal random variable with a nonzero mean; or it can be found as a function of the first four moments of a standard normal random variable. We find it easier to note that the needed expectation is that of a noncentral chi-square random variable on 1 df and noncentral parameter τ^2 . The result is

$$E_f \left[\sum_{i=1}^n (y_i - \underline{x}_i' \underline{\beta}_o)^4 \right] = \tau^4 \left[\sum_{i=1}^n [3 + 6\lambda_i + \tau^2] \right]$$

Now, by carefully constructing the full result from all the pieces and simplifying it, we get

$$J_{KK}(\underline{\theta}_o) = \frac{n}{2(\sigma_o^2)^2} \left[\frac{2\tau^2\sigma_o^2 - \tau^4}{(\sigma_o^2)^2} \right].$$

While derived only for $\tau > 0$, the above result can also be used in the case of $\tau^2 = 0$.

Finally,

$$J(\underline{\theta}_o) = \begin{bmatrix} \frac{\tau^2}{(\sigma_o^2)^2} X'X & \underline{0} \\ \underline{0}' & \frac{n}{2(\sigma_o^2)^2} \left[\frac{2\tau^2\sigma_o^2 - \tau^4}{(\sigma_o^2)^2} \right] \end{bmatrix}$$

The result we sought can now be found:

$$\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] = \frac{\tau^2}{\sigma_o^2} \left[K + 1 - \frac{\tau^2}{\sigma_o^2} \right]$$

If model equals truth, we have $\underline{\mu} = X\underline{\beta}$, so that $\sigma_o^2 = \tau^2$ (and the trace term equals K). By continuity in τ^2 we may use the trace term as K when $\tau^2 = 0$ if the model is true. However, for philosophical issues and problems associated with a truth concept, we will consider only the situation wherein even for truth there is unexplainable uncertainty. In particular, if true replication is not possible (or study), we suggest that it is most useful to consider the variance within true replicates (assuming variance homogeneousity) as a definition of convenience, as even truth can be at different levels. We are mostly interested in structural truth of our models in the face of irreducible uncertainty inherent in data for finite sample sizes.

Surprisingly enough, we see from (7.79) that for this case, and context, $\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] < K$ under a misspecified

average, on the side of parsimony. Moreover, estimation (hence, TIC) seems very problematic, since τ^2 cannot be lacking true replication (cf. Linhart and Zucchini 1986) call true replication in an experiment provides only a definition, if we restrict our concept of truth to what was the design structure and independent variables used in the experiment. Philosophically, we might be able to predict some (hence, some $\tau^2 = 0$) of the observed differences among replicates from the true ultimate truth.

This example can be easily generalized; that is, we can assume a normal model and generalize truth somewhat. Whereas we assume $\epsilon \sim \text{normal}(0, \tau^2)$, the only way this entered the derivations was through the central moments of ϵ . If we retain the *iid* assumption, we can generalize the results; note that we retain $E(\epsilon) = 0$ with loss of generality, we assume an asymmetric distribution for ϵ ; we will not do so: We assume $\epsilon \sim \text{normal}(0, \tau^2)$ all we need is the fourth moment of ϵ , which we will express in the following form as

$$\gamma = E(\epsilon^4)/[E(\epsilon^2)]^2.$$

For assumed normal truth, $\gamma = 3$. For f as a logistic distribution, $\gamma = 15$; for a Laplace distribution, $\gamma = 6$; and for a uniform $(-h, h)$ distribution, $\gamma = 9/5$. The last two are extreme cases; one might think that γ lies in the range 2 to 4.

Redoing the derivations for this more general way of expressing the error distribution is straightforward; the results are the same for $I(\underline{\theta}_o)$; but for

$$J(\underline{\theta}_o) = \begin{bmatrix} \frac{\tau^2}{(\sigma_o^2)^2} X'X & \underline{0} \\ \underline{0}' & \frac{n}{2(\sigma_o^2)^2} \left[\frac{2\tau^2\sigma_o^2 + \tau^4 \left(\frac{\gamma-3}{2} \right)}{(\sigma_o^2)^2} \right] \end{bmatrix}$$

For the trace function we get

$$\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] = \frac{\tau^2}{\sigma_o^2} \left[K + 1 + \frac{\tau^2}{\sigma_o^2} \left[\frac{\gamma - 1}{2} \right] \right]$$

Hence, for these fixed-effects linear models assuming normal errors, the effect of structural misspecification appears in σ_o^2 as manifest in the trace function. However, the effect of error distribution misspecification is manifest in the fourth moment, γ (assuming symmetric errors). We have stressed models wherein K will not be trivially small; from (7.80) it is clear that the large the effect of error distribution misspecification upon the trace function becomes trivial. In contrast, the effect of structural mis-

misspecification of the error distribution becomes trivial. The use of the approximation $\text{tr} [J(\underline{\theta}_o)[I(\underline{\theta}_o)]^{-1}] = K$, hence the use of AIC rather than TIC. These musings seem likely to apply to linear models.

One last point: How good is use of $[I(\underline{\theta}_o)]^{-1}$ for the matrix of $\hat{\underline{\theta}}$ under model misspecification here (ignoring σ_o^2 that is)? As was shown (and is known in general) in Section 7.7.1, the asymptotic variance–covariance matrix is $V(\hat{\underline{\theta}}) = [I(\underline{\theta}_o)]^{-1}$. Thus, here, $V(\hat{\underline{\beta}}) = \tau^2(X'X)^{-1}$, and

$$V(\hat{\sigma}^2) = \frac{2(\sigma_o^2)^2}{n} \left[\frac{2\tau^2\sigma_o^2 + \tau^4 \left(\frac{\gamma-1}{2} - 2 \right)}{(\sigma_o^2)^2} \right]$$

If a good structural fit has been achieved (so $\tau^2 \approx \sigma_o^2$), then

$$V(\hat{\sigma}^2) = \frac{2(\sigma_o^2)^2}{n} \left[\frac{\gamma-1}{2} \right].$$

The drastic bias induced by a γ not near 3 (but assumed to be one to use an estimator of γ , and with the same estimator for AIC. This approach can be recommended only weakly, since the estimator of the fourth moment is so highly variable.

7.7 Additional Results and Considerations

7.7.1 Selection Simulation for Nested Models

The detailed stochastic characteristics of the model selection process can be studied mostly by Monte Carlo simulation methods. Clearly, this is not if the only completely general approach is to specify a data-generating process and a set of models to be fit to each generated sample, generate many samples and do all the calculations associated with model selection. This is very useful, but is not a study of properties of model selection in the abstract. Rather, each application has some underlying process and type of truth (as generated data), and may require extensions that are peripheral to the heart of the model selection problem.

An exception arises if we restrict ourselves to a single class of models and to selection methods based on log-likelihood differences between models. This scenario includes simulation of AIC, AIC_c, BIC, ratio testing-based methods (and can be easily adapted to QAIC_c model selection). All we need to generate are the independent and identically distributed log-likelihood chi-square random variables.

but no parameters need actually be postulated or estimated. The important restriction here is that we can correctly generate variables only for what would be a single chain of nested models. The advantage of the method is speed and generality; this allows us to study the properties of some model selection procedures.

At the heart of this procedure we have (conceptually) a sequence of models g_1, g_2, \dots, g_R , where g_i is nested in model g_{i+1} , and the difference in the number of parameters is j . The method can be developed in general for any j , and we use it, with $j = 1$. That is, our conceptual sequence of models is $g_1 \subset g_2 \subset \dots \subset g_R$, and each incremented model (i.e., g_{i+1}) has only 1 added parameter. We assume that large sample theoretical results for parameter estimation conceptually underlie such Monte Carlo simulation. The usual likelihood ratio test statistic has, in general, a noncentral χ^2 distribution on 1 df; denote that random variable by $\chi_1^2(\lambda_i)$, where λ_i is the parameter for model g_i versus g_{i+1} is λ_i .

For our set of R models we have $R - 1$ noncentrality parameters $\lambda_1, \dots, \lambda_{R-1}$. These λ_i would be functionally related to the true data-generating process, the model structures assumed $g_i(\underline{x} | \cdot)$, and the specific parameters of the data that specify the actual best approximating model in each case. However, we will be able to bypass all of those specifications in the method below. We do need to be able to interpret sets of noncentrality parameters. If we had $\lambda_{i-1} > 0$ and $\lambda_i = \dots = \lambda_{R-1} = 0$, we would interpret g_{i-1} as being the true data-generating model (we will ignore parameters of g_{i-1} that invalidate this interpretation). Also, it is possible to have noncentrality parameters for example, $\lambda_1 = 0$ but $\lambda_2 > 0$. Then both models g_1 and g_2 are bad approximations to truth, because model g_3 improves over g_2 and closer to truth compared to model g_2 , but g_2 does not improve over g_1 . In this situation we would have the set of λ_i monotonically decreasing. The model is that of which model provides the AIC best model, i.e., the best model, when parameter estimation occurs. We continue with the simulation how to do simulation in this context.

Let the fitted model log-likelihoods be $\log(\mathcal{L}_i)$. From the

$$2 \log(\mathcal{L}_{i+1}) - 2 \log(\mathcal{L}_i) \sim \chi_1^2(\lambda_i).$$

Let the number of parameters in the simplest model be K , and let i be the above as

$$-(\text{AIC}_{i+1} - 2K_1 - 2i) + (\text{AIC}_i - 2K_1 - 2i + 2),$$

or

$$\text{AIC}_i - \text{AIC}_{i+1} \sim \chi_1^2(\lambda_i) - 2.$$

$$\text{AIC}_i = \text{AIC}_{i+1} + \chi_1^2(\lambda_i) - 2.$$

Thus for purposes of a simulation study, if we know λ_i generate AIC_i . We just need to be able to generate a noncentral chi-square random variable (there are routines for this, such as CINV).

Based on (7.82), we can do a backwards recursive generate AIC_1 starting with $i = R - 1$, given a value for AIC_R . The formula

$$\text{AIC}_i = \text{AIC}_R + \sum_{j=i}^{R-1} (\chi_1^2(\lambda_j) - 2);$$

we just need a value for AIC_R . Because everything we care about in model selection depends only on the relative differences, say Δ_i or Δ_p , it suffices to set $\text{AIC}_R = R$ (any constant will do, one has advantages) for every sample of AICs generated. This corresponds to a realization of a set of independent $\chi_1^2(\lambda_i)$.

We have used this approach to do simulation studies of model selection with AIC and other likelihood-based methods. To evaluate $\text{AIC}_i - 2i + i \log(n)$. Because n should now vary, one must vary the noncentrality parameters on a per-unit sample size basis. Compute $\lambda_{i,n} = n\lambda_{i,1}$. The $\lambda_{i,1}$ should be very small, but can be chosen as arbitrary. To mimic AIC_c selection, use

$$\text{AIC}_{c,i} = \text{AIC}_i + 2 \frac{K_i(K_i + 1)}{n - K_i - 1},$$

where $K_i = K_1 + i$. Now one must specify K_1 , the number of parameters envisioned in model g_1 , as well as use $\lambda_{i,n} = n\lambda_{i,1}$.

To mimic QAIC model selection is a little more involved. One needs a value of c (variance inflation factor; $c \geq 1$) and its df (independent of n , vary over samples). Generate $\hat{c} = \chi_{\text{df}}^2(\text{df}(c - 1))/\text{df}$ for each sample. $\text{df}(c - 1)$ is the noncentrality parameter for this chi-square. Then use

$$\text{QAIC}_i = \text{QAIC}_R + \sum_{j=i}^{R-1} \left[\frac{\chi_1^2(c - 1 + c \cdot \lambda_j)}{\hat{c}} - 2 \right],$$

and $\text{QAIC}_R = R$. Also,

$$\text{QAIC}_{c,i} = \text{QAIC}_i + 2 \frac{K_i(K_i + 1)}{n - K_i - 1}.$$

A few results are given below using this simulation approach applied to model selection. For the most part, however, it is not possible to present extensive tables of simulation results. We encourage interested

From this setup which allows simulating model selection, we can also compute theoretical expected AIC *differences*, expected AIC best model exactly. From (7.83) we get, for n by just one parameter,

$$E(\text{AIC}_i) = E(\text{AIC}_R) + \sum_{j=i}^{R-1} (E[\chi_1^2(\lambda_j)] - 1)$$

$$E(\text{AIC}_i) = E(\text{AIC}_R) + \sum_{j=i}^{R-1} (\lambda_j - 1).$$

Let $\lambda_{i+} = \lambda_i + \cdots + \lambda_{R-1}$, $i = 1, \dots, R-1$, and $\lambda_{R+} = 0$.

$$E(\text{AIC}_i) = (E(\text{AIC}_R) - R) + (\lambda_{i+} + i).$$

Then compute the set of values V_i given by

$$V_i = \lambda_{i+} + i, \quad i = 1, \dots, R,$$

find their minimum, V_{\min} , and then compute

$$E(\Delta_i) = V_i - V_{\min}.$$

As an example, if $R = 10$ and (in order) we have λ_i as 2, 0.4, and 0.2, then the $E(\Delta_i)$ are, in order, 22.5, 21.5, 16.5, 0.8, and 1.6. Thus the best expected AIC selected model is the one with the smallest variance. The expected AIC variances can also be computed, but nothing directly useful like Δ_p .

7.7.2 Simulation of the Distribution of Δ_p

The random variable $\Delta_p = \text{AIC}_{\text{best}} - \text{AIC}_{\min}$ was introduced in (Section 4.5). For a set of models indexed $i = 1, \dots, R$ of size n , and a conceptually well-defined repeated sampling (over a sample space), we let model g_{best} represent the best model to fit under the AIC selection criterion. Monte Carlo simulation can determine this actual best model (sometimes theory suffices). We are not saying that model g_{best} is truth; it is just that it must be the best model, on average, to use for all possible samples. The truth that model g_{best} represents.

For each simulation-generated sample we can compute AIC_{\min} . This AIC_{\min} and the value of AIC_{best} vary by sample; the value of best is fixed for all samples; for example, model g_4 is the best model to use (hence $\text{best} = 4$). If model g_{best} is selected in the sample, then $\Delta_p = 0$; otherwise, $\Delta_p > 0$. We can compute the distribution, hence percentiles, of this pivotal under the

for R and $\lambda_1, \dots, \lambda_{R-1}$, if we use AIC_g , we also need λ_R . Results about Δ_p under AIC model selection, for a few $\lambda_i = 0$, are feasible to show. In this case model g_1 is the model ($k = 1$ in Δ_p). This scenario is clearly at odds with what applies to real data analysis (all models for data analysis are equally close to truth). However, it can be used as a benchmark for purposes to the extent that this situation is too simple, it may serve only as a guide on the percentiles of the cumulative distribution function for nested models, or for real problems where there is substantial uncertainty of the models considered.

For the case of $R = 2$ and large sample size, the q th percentile ($0 < q < 1$), $\Delta_{p,q}$, is

$$\Delta_{p,q} = \max\{0, (\chi^2_{1,q} - 2)\}$$

(easily derivable from 7.81). Here, $\chi^2_{1,q}$ is the q th percentile of a square random variable on 1 df. For example, $\chi^2_{1,0.95} = 3.84$. We used Monte Carlo simulation to determine some values of $R > 2$. One million samples were used for each of 100 independent runs of 250,000 samples, so we can estimate the values below for $R > 2$ have a cv of about 0.5%:

R	percentiles of Δ_p			
	80%	90%	95%	99%
2	0.00	0.71	1.84	4.63
3	0.11	1.37	2.67	5.77
4	0.35	1.71	3.33	6.40
5	0.49	1.93	3.40	6.86
10	0.75	2.34	3.97	7.61
20	0.82	2.47	4.15	8.05

We have done many of these simulations to find the percentiles for sets of noncentrality parameters wherein $\lambda_i > 0$; the results are somewhat larger as compared to the case where all $\lambda_i = 0$. For example, for $R = 10$ and $\lambda_1, \dots, \lambda_9$ be 2, 6, 10, 6, 3, 1.5, 0.8, 0.4, 0.2. In this case the set of models (actually, g_R could be truth; we cannot be sure) was found on 20,000 Monte Carlo samples (two sets of 10,000), we found the expected K-L (i.e., AIC) best model. The averages of the K-L values (rescaled so their minimum is 0), in order, are 22.6, 21.5, 20.5, 19.5, 18.5, 17.5, 16.5, 15.5, 14.5, 13.5, 12.5, 11.5, 10.5, 9.5, 8.5, 7.5, 6.5, 5.5, 4.5, 3.5, 2.5, 1.5, 0.5, 0.2, 0.8, 1.6 (reliable to $\pm 0.1 = 2$ se); compare these values with the expectations from the end of Section 7.7.1: 22.5, 21.5, 20.5, 19.5, 18.5, 17.5, 16.5, 15.5, 14.5, 13.5, 12.5, 11.5, 10.5, 9.5, 8.5, 7.5, 6.5, 5.5, 4.5, 3.5, 2.5, 1.5, 0.2, 0.8, and 1.6.

Based on $\Delta_p = \text{AIC}_7 - \text{AIC}_{\min}$ over these 20,000 samples, the percentiles of Δ_p are 3.3 (80%), 4.6 (90%), 6.4 (95%), 10.6 (99%). From these sorts of simulations, and others with explicit noncentrality

sequences, a model g_i for which $\Delta_i = 1$ is implausible under the model structure, and $\Delta_i \approx 7$ is strong evidence against g_i being the K-L best model (and $\Delta_i \geq 10$ is very strong evidence against g_i). The Akaike weights provide a refined interpretation of these results.

7.7.3 Does AIC Overfit?

The conceptual framework underlying valid use of AIC is that there are infinitely many parameters. Overfitting is often defined in terms of a true model: if there is a simple true model, with a finite number of parameters, and this model is in the set of models considered. Then if the true model is nested within the selected model structure, the selected model is overfitted. One has estimated more parameters than are in the true model.

This simplistic concept of overfitting does not apply in the selection framework. However, there is a best expected K-L model. If the model we should use as our basis for data analysis is nested within the selected model, might we claim that the selected model is an overfit model? We decline to use this definition because it is based on a lack of ability in the model selected. If we miss the target model (or, what is the same, that the structure of the selected model is the same as that of the target model), we should not say that the selected model is overfitted. It is philosophically the same as being concerned that a point estimate may sometimes give a point estimate far away from the true value. We are upset only by cases where $\hat{\theta}$ is quite far from θ . But if we consider $\hat{\theta}$ as an acceptable estimate with suitably small probability, we consider $\hat{\theta}$ as an acceptable estimate. The argument becomes somewhat circular at this point because we have accustomed ourselves to being satisfied if $\hat{\theta}$ is within a certain range, hence $\hat{\theta}$ is unacceptable with a probability of only about 0.05. Practitioners of null hypothesis testing typically are willing to accept (as a type I error) a 0.05 probability of type I error (and probably a smaller probability in most applications).

Something similar should apply to possible overfitting in the selection framework (underfitting is at best a minor concern with AIC). We need some idea of how far from the actual K-L best model the selected model is before it is regarded as an overfit model. We have to allow for the fact that we do well to select a model within, say, 1 or 2 parameters of the K-L best model. In contrast, if the selected model has 10 or 20 parameters more than the target model, we think that most people will agree that the selected model is overfitted. We can use simulation (in some cases theory exists, see Yu and Yu 1993) to find the probability distribution of the selected model. We can say *best*, for a nested sequence of models g_1 to g_R and accept the model with the smallest AIC value.

Here is a typical example of what can happen. For models with $R = 10$ and the λ_i as 2, 6, 10, 6, 3, 1.5, 0, the best model was found to be g_7 . We extend this to R beyond 10 with all additional $\lambda_i = 0$ (and assume that sample size n is large enough). This means that the true data-generating model is model g_7 and models g_{30} are overfit, if selected, in the sense that they do contain the true structure. The K-L best model remains model g_7 even for $R = 30$. We give the model selection frequencies, in model order g_1 to g_{30} , for 11,000 Monte Carlo samples:

R	g_7																													
10:	11	5	50	746	2282	2635	1924	1161	735	451																				
20:	16	3	20	405	1557	2399	1970	1373	816	542	274	192	125	93																
30:	15	3	23	407	1689	2351	1954	1360	800	522	237	155	121	90																

11

The long-tailed nature of the distribution of selected AIC when there are many “big” models (models with many parameters) containing the K-L best model structure. For $R = 10$ we have a probability of about 0.01 of selecting a model with more than 10 parameters (models g_{19} to g_{30}), hence having estimated 10 or more parameters. For both $R = 20$ and 30 there is about a 1% chance of selecting a model with five or more unneeded parameters (models g_{25} to g_{30} higher). In general, if we say that we can accept a procedure with a probability of its cases a bit misleading, then we should not be upset if the K-L best model by about 5 or more parameters with probability 0.06 if many such too-general models are in the set of models considered.

Note that these results are effectively for very large sample sizes, where AIC_c was used, not AIC_c . For not-large n , use of AIC_c would lead to a much longer long tail at $R = 20$ and 30 (say $n = 100$). Hence, the results are for the worst-case scenarios.

Here is a worst-case scenario for one linear sequence of models. All $\lambda_i = 0$, so model g_1 is the K-L best model. Based on 11,000 Monte Carlo samples, the estimated model selection probabilities (with error $\leq 0.05\%$) are below:

R	$\hat{\pi}_i \times 100\%$, in order $i = 1, 2, \dots$																													
3	78.7	13.3	8.0																											
4	76.0	12.5	6.7	4.8																										
5	74.4	12.0	6.4	4.1	3.1																									
10	71.8	11.4	5.8	3.5	2.4	1.7	1.2	.9	.7	.6																				
20	71.2	11.2	5.7	3.4	2.3	1.6	1.1	.8	.6	.5	.4	.3	.2	.2																

(For this case of a single nested sequence of models, these probabilities are known theoretically for large n and all $\lambda_i = 0$; see Akaike (1989). The long tail is disturbing, yet with probability about 0.06 of selecting a model with 5 or more unneeded parameters.)

of realistic AIC model selection, we claim that this is a type of sample AIC (about 0.06 probability of overfitting by five parameters if such general, overparametrized models are in use). Given models with a large amount of unneeded structure, we can select overfitted models, but the probability of a serious overfit is arguably less than the total error probabilities (type I plus type II) in hypothesis testing.

What will reduce the probability of getting a badly overfitted model? AIC_c, which helps considerably when sample size, relative to the number of parameters, is small. Otherwise, the only recommendation we have to avoid the selection of a much overfit AIC-selected model is to be very thoughtful in the set of models considered. In particular, do not casually include a great many parameters more than you think are really needed. In variable selection this would mean do not simply consider all possible regressor variable and include it for possible selection. If you have a large number of variables that have no explanatory power (a $\lambda \approx 0$), and that leaves you with a small but real probability of selecting a model with many worthless variables.

7.7.4 Can Selection Be Improved Based on All the Information?

Given the potentially long-tailed nature of K-L-based model selection (depends on the set of models), it seems natural to ask whether we can use information in the full set of Δ_i values that would allow us to avoid cases where we have selected a very overfit model. If so, we can select to a better model based on information in the entire set of Δ_i . We have explored this matter for a single series of nested models (as considered in the above three subsections). The idea was that perhaps the Δ_i would be like that below in the event of selecting a model (line one is model number i ; line two is Δ_i):

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
30	10	5	0.1	1	2.5	2	1.5	2.7	3	0.8	2.3	1	0	2.5

Here, AIC has selected model g_{14} , but we might suspect that a better model to use is g_4 ($\Delta_4 = 0.1$ and none of models g_5 to g_{13} have $\Delta_i < 0.1$). So we could change our choice to model g_4 .

There is no theory to help here regarding properties of the Δ_i . So we looked at a large number of simulated results for the worst-case scenario of all $\lambda_i = 0$. Hence, the simplest model is the true data-generating model and the K-L best model. Table 7.7.4 shows $R = 10$, some selected cases of $\Delta_1, \dots, \Delta_{10}$ (a case is one set of selection results, for a large sample size). Cases 1, 2, and 3 show model g_1 is selected in about 72% of all samples here. Case 4 shows

which informative results can be computed, for a certain regression under constant error variance. Specifically, we use linear regression based on orthogonal regressors, but only use a subset of the regressors (the global model may include z_1, \dots, z_m be *iid* normal(0,1). Independently, let ϵ be noise response variable x , based on the regressors, is given by

$$x = \beta_o + \sum_{j=1}^m \beta_j z_j + \epsilon,$$

so

$$E(x | \underline{z}) = \beta_o + \sum_{j=1}^m \beta_j z_j.$$

Results below are scale-invariant in terms of the regressors z_j occur in the models only as a product. To keep notation simple, Section 7.7.1, we define here the base model, g_1 , to be δ with δ assumed as normal(0, σ_1^2). The number of parameters is $K_1 = 3$. The normal assumption is true, but $\sigma_1^2 \neq \sigma^2$. In model g_r is

$$x = \beta_o + \sum_{j=1}^r \beta_j z_j + \delta,$$

with δ assumed as normal(0, σ_r^2); the number of parameters is K_r .

We can consider the sequence of nested models for $r = 1, \dots, m$. The ordering of regressors is arbitrary, but is used in the following. The nested models are considered sequentially (as per the theory). It is then convenient, but not required, to specify the regressors satisfy $|\beta_j| > |\beta_{j+1}|$. Doing so yields insights into AIC model selection easily (such a structured situation is also considered in Section 7.7.1).

From the point of view of insights to be gained, the assumption that the regressors is not restrictive if $R = m$, because any regression can be transformed into the case of orthogonal regressors, for example, by regression on principal components. However, if we assume $R < m$, the orthogonality assumption is restrictive, because our observed regressors are not orthogonal to the regressors not observed (realistic of real data). For this reason there is no advantage in assuming $R < m$.

The regressors are random variables. Therefore, to generalize, we take an additional expectation over certain matrices to be random. Let $\underline{z}_r = (1, z_1, \dots, z_r)$. For a sample of n observations, the matrix notation is $\underline{x} = \underline{Z}\underline{\beta} + \underline{\delta}$, so matrices such as $\underline{Z}'\underline{Z}$ and

by their expectations with respect to the random \underline{z}_r (for large n , $E(Z'Z)^{-1} \approx I$). However, the nonlinearities in some theoretical formulas below are only large-sample approximations.

Under this scenario of regressor independence we determine the noncentrality parameters of Section 7.7.1 that apply, for large n , to the nested sequence of models defined here (g_1 to g_R):

$$\lambda_i = n \log \left[1 + \frac{(\beta_{i+1})^2}{(\beta_{i+2})^2 + \cdots + (\beta_m)^2 + \sigma^2} \right],$$

$$\lambda_{m-1} = n \log \left[1 + \frac{(\beta_m)^2}{\sigma^2} \right].$$

Given this context and the theory in Section 7.7.1 we can approximate (expected Δ_i values for AIC and AIC_c, and (approximate) theoretical expected K-L optimal model. We compare these to the actual sets of Δ_i values. The needed partial sums of noncentrality parameters are

$$\lambda_{i+} = n \log \left[1 + \frac{\sum_{j=i+1}^m (\beta_j)^2}{\sigma^2} \right], \quad i = 1, \dots, m-1.$$

The value of exploring this situation is that we can also compare theoretical quantities that can be related, or compared, to actual values. Under any of the models the regression coefficient estimates are unbiased for the true parameters (because all regressors are orthogonal to the error term, expected in general). Under model g_r the value of σ_r^2 is

$$\sigma_r^2 = (\beta_{r+1})^2 + \cdots + (\beta_m)^2 + \sigma^2.$$

In notation used elsewhere in Chapter 7, the above σ_r^2 is the error variance. The usual (conditional on the model) parameter cv's under model g_r are

$$cv(\hat{\beta}_r) = \frac{\beta_r}{se(\hat{\beta}_r | g_r)} = \frac{\sqrt{n}\beta_r}{\sigma_r}, \quad r = 1, \dots, m.$$

However, rather than compute the above cv's it is informative to compute the cv's under the global model, $\sqrt{n}\beta_r/\sigma$, and observe the difference between the cv's versus what parameters are included or excluded from the optimal model.

We can also determine the overall mean square error of the model. Minimum MSE is generally accepted as a good criterion for model selection; here we have

$$MSE(Z) = \sum_{j=1}^n E_f \left[\hat{E}(x_j | \underline{z}_r) - E(x_j | \underline{z}_m) \right]^2$$

here over Z to get the unconditional result

$$\text{MSE} = \sigma_r^2(r+1) + \sum_{j=r+1}^m n(\beta_j)^2.$$

Sometimes model selection is based on minimum MSE response variable for a single additional (independent of g_r) regressors. Then this average mean square error of prediction g_r is $\text{MSEP} = \text{MSE} + n\sigma^2$. Thus, it suffices to consider C_p implements minimum MSE model selection for regression (see, e.g., Akaike, 1995).

The K-L-based target criterion T (7.20) can be determined. We express that result as $-2T$ for direct comparison to MSE-based results, and we will label it here as KL. Hence the model g_r is exactly

$$\begin{aligned} \text{KL} = & nE[\log(\chi_{n-r-1}^2(0))] + n \log \left[\frac{\sigma_r^2}{n} \right] \\ & + \frac{n}{n-r-3} \left[(r+1) + \left(\frac{n\sigma^2}{\sigma_r^2} \right) + \sum_{j=r+1}^m \right] \end{aligned}$$

In MSE (7.88) the term $\sum_{j=r+1}^m n(\beta_j)^2$ is bias squared for excluded regressors. Note that in MSE these components of bias are due to theoretical precision of the excluded $\hat{\beta}_j$, whereas in K-L the bias of excluded regressors is “judged” relative to the theoretical precision of $\hat{\beta}_j$ have. That is, the comparable term reflecting bias is $n(\beta_j/\sigma_r)^2$, not just $(\beta_j)^2$. This feature of the KL criterion selection is more than just optimizing on pure bias versus variance as in MSE.

Most insights based on all of these results will need to come from examples and simulation. Because of the possible volume of such considerations (all the variables here), we leave such computing to the reader. We have done a lot of computing and simulation to verify these formulas and Section 7.7.1 results. One result (known about C_p versus AIC) is that K-L-based model selection is about the same as selection based on minimum theoretical variance. We terminate this by computing the theoretical criteria KL and AIC, and is convenient to rescale the R values of those criteria to range from zero. Table 7.6 gives such results for one case: $R = 0.6$, $\beta_i = (0.6)^{i-1}$, and σ taking several values in the range 0.1 to 1.0.

We have not undertaken a detailed analysis of these criteria for regression, let alone for this case of orthogonal regressors. The results are representative of cases we have looked at in the past. The K-L best model has either the same number of parameters

0.0, results for both KL and MSE have been rounded so that they clearly indicating the theoretical optimal model under these criteria.

r	$\sigma = 1.0$		$\sigma = 0.5$		$\sigma = 0.25$		$\sigma = 0.1$	
	KL	MSE	KL	MSE	KL	MSE	KL	MSE
1	8.01	14.08	27.17	16.66	57.72	17.58	142.67	17.58
2	1.70	3.76	11.16	5.59	33.52	6.32	113.80	6.32
3	0.00	0.56	2.75	1.64	15.09	2.18	85.50	2.18
4	0.59	0.00	0.00	0.33	4.34	0.68	58.46	0.68
5	2.26	0.42	0.28	0.00	0.20	0.17	34.20	0.17
6	4.52	1.21	2.01	0.04	0.00	0.02	15.33	0.02
7	7.20	2.13	4.50	0.21	1.74	0.00	4.10	0.00
8	10.29	3.10	7.52	0.43	4.48	0.03	0.00	0.00
9	13.83	4.09	11.03	0.67	7.90	0.09	0.46	0.00
10	17.90	5.08	15.09	0.91	11.92	0.14	3.33	0.00

parsimonious than the theoretically best model under minimum expected models selected under operational criteria such as AIC, CAIC differ slightly from these results (but less so as n gets large).

7.7.6 AIC_c and Models for Multivariate Data

The derivation of large-sample AIC in Section 7.2 *does* apply to independent multivariate observations, each with p nonindependent components. The small-sample improvement of AIC that applies to $p = 1$ linear models with homogeneous normal residuals also apply in the corresponding multivariate case. This problem was studied by Fujikoshi and Satoh (1997). They focused on selection of the best model, i.e., inclusion or exclusion of the same set of possible regressors among the p regressions. They assume that a general $p \times p$ covariance matrix Σ applies for the residual vector of each observation. With k regressors (this may include an intercept) there are k parameters. Each model also includes $p(p + 1)/2$ unknown parameters. Thus, $K = (k \cdot p) + p(p + 1)/2$.

For their data analysis context Fujikoshi and Satoh (1997) give AIC, CAIC in their notation, analogous to the univariate case. Their result (their formula 7) can be expressed as follows:

$$CAIC = AIC + 2 \frac{K(k + 1 + p)}{n - k - 1 - p}.$$

The univariate case corresponds to $p = 1$, and then $K = k + 1$. The general result in (7.90) reduces to the univariate AIC_c .

other multivariate modeling problems we do not know what to do. Our key point here is that the univariate result for AIC_c does not generalize to the multivariate setting.

The form of (7.90), by virtue of including variable selection, is a restricted context considered by Fujikoshi and Satoh (1984). In (7.90) we hypothesize a generalization of the univariate result to corresponding multivariate applications:

$$AIC_c = AIC + 2 \frac{K(K+v)}{np - K - v}.$$

In (7.91) v is the number of distinct parameters used in the model. $1 \leq v \leq p(p+1)/2$. Note that the count K includes the error variance σ^2 . This count is correct for the univariate case wherein $v = p = 1$. It seems reasonable until a derivation is published for the multivariate case of AIC_c to multivariate applications.

The emphasis of this section is the generalization of AIC_c to multivariate applications. Many multivariate analysis methods, such as principal component analysis, analysis of variance or covariance, are done in a least squares fashion while assuming a multivariate normal model such as $MVN(\mu, \Sigma)$ (where Σ is assumed to be full rank). Given that least squares is used, the software packages typically report the value of the maximized log-likelihood. However, the residuals and cross products matrix (SSCP) is nearly always provided. The log-likelihood of the software packages and the MLE of Σ is $\hat{\Sigma} = SSCP/n$. The maximized log-likelihood is proportional to $-(n/2)\log(|\hat{\Sigma}|)$, where $|\hat{\Sigma}|$ is the determinant of $\hat{\Sigma}$. Hence, as long as all models considered are multivariate normal residuals, we may use

$$AIC = n \times \log(|\hat{\Sigma}|) + 2K$$

and

$$AIC_c = n \times \log(|\hat{\Sigma}|) + 2K + \frac{2K(K+v)}{n - K - v}.$$

In the univariate case SSCP is just the residual sum of squares and the MLE of σ^2 is RSS/n . Thus the multivariate case is just a generalization of the univariate case. Theory for the multivariate case is due to Fujikoshi and Satoh (1984:61).

If the determinant $|SSCP|$ is directly available, one could use $AIC = n \times \log(|SSCP|) + 2K$. However, because $\log(|\hat{\Sigma}|) = -p \times \log(n) + \log(|SSCP|)$ and p and n are constants we can just as well take, for example,

$$AIC = n \times \log(|SSCP|) + 2K.$$

7.7.7 *There Is No True TIC_c*

Fujikoshi and Satoh (1997) also consider a small sample size. In essence they want to extend AIC_c to TIC_c . However, AIC_c (Burnham 1989) arises by computing the exact value of the target model (7.20) for a linear model with constant normally distributed error condition (assumption) that this model is the true data-generating process. That derivation is derived without any assumption that truth, f , is the true model. That derivation can be justified only for large sample sizes. For a small sample, an exact version of TIC would require us to specify the distribution f (i.e., specify truth). Even if we could do so, in general, or at all, the result would depend upon assuming the truth. Thus no defensible, general small-sample analytical version of a TIC_c seems possible.

The issue did not escape the attention of Fujikoshi and Satoh. What they did (and they knew it) was to assume that the linear, global model defined by using all available regressors contained the true model as an unknown submodel. Thus the true model is, by assumption, one of models considered and is a special case of the global model. Under these assumptions Fujikoshi and Satoh (1997) derived an analytical formula for sample size, for the target criterion of (7.20). Their formula is beyond $-2 \log(\mathcal{L})$, that must be estimated from the data, and involves elements that are simple functions of known n and K , and that are not estimable only by virtue of the strong assumptions made. Applying the same small-sample-size adjustments with TIC as we do with AIC_c may be a good idea; but we cannot find truly general adjustments for TIC.

7.7.8 *Kullback–Leibler Information Relationship and Information Matrix*

The Fisher information matrix is defined by (7.6) for any

$$\mathcal{I}(\underline{\theta}) = E_g \left[- \frac{\partial^2 \log(g(\underline{x} | \underline{\theta}))}{\partial \theta_i \partial \theta_j} \right].$$

In taking this expectation it is assumed that the true data-generating process is $g(\underline{x} | \underline{\theta})$ (hence the underlying integration is with respect to \underline{x}) when this one particular member (i.e., g at $\underline{\theta}_o$) of the set $\{g(\underline{x} | \underline{\theta})\}$ for fixed structure, and any $\underline{\theta} \in \Theta$, is the generating model.

$$I(\underline{\theta}_o) = E_f \left[- \frac{\partial^2 \log p(\mathbf{x} | \underline{\theta}_o)}{\partial \theta_i \partial \theta_j} \right],$$

in general, $I(\underline{\theta}_o) \neq \mathcal{I}(\underline{\theta}_o)$. Moreover, $\mathcal{I}(\underline{\theta}_o)^{-1}$ is guaranteed to be the sample variance–covariance matrix of the MLE $\hat{\underline{\theta}}$ only if $\hat{\underline{\theta}} = \underline{\theta}_o$. For a value of $\underline{\theta} \in \Theta$ that is near the K-L minimizing distribution, the quadratic approximation to the K-L difference is

$$I(f, g(\cdot | \underline{\theta})) - I(f, g(\cdot | \underline{\theta}_o)) \approx \frac{1}{2}(\underline{\theta} - \underline{\theta}_o)' \mathcal{I}(\underline{\theta}_o)(\underline{\theta} - \underline{\theta}_o).$$

For the case of $f = g$ we get the result

$$I(g(\cdot | \underline{\theta}_o), g(\cdot | \underline{\theta})) \approx \frac{1}{2}(\underline{\theta} - \underline{\theta}_o)' \mathcal{I}(\underline{\theta}_o)(\underline{\theta} - \underline{\theta}_o).$$

Thus if one member of the set of models $g(\underline{x} | \underline{\theta})$, $\underline{\theta} \in \Theta$, is the true generating distribution, then the approximate K-L information is obtained from using a nearby distribution as the approximating distribution. This is a quadratic form in the Fisher information matrix. This is not surprising given the definitions of both $\mathcal{I}(\underline{\theta})$ and $I(g(\cdot | \underline{\theta}_o), g(\cdot | \underline{\theta}))$. We will show that the two underlying concepts of “information” are quite different concepts. It was in the 1920s that Fisher chose to use the term “information matrix” of second mixed partials of a probability distribution. There is no relationship to information theory, which is a concept that arose mostly since Shannon’s pioneering work in the late 1940s. Information theory deals with logs of probabilities. The Fisher information matrix relates to the precision of ML estimators. The Kullback–Leibler divergence (Kullback 1951) was a result of their attempt to understand and explain the relationship by “information” in relation to sufficiency (personal communication with Leibler).

7.7.9 Entropy and Jaynes Maxent Principle

In Section 2.9 we noted that the Akaike weights w_i can be reinterpreted in a Bayesian approach based on prior probabilities τ_i . To choose the most probable model in a manner philosophically consistent with the data, we suggest resorting to the use of the Jaynes maximal entropy principle (Jaynes 1957, 1982, Jessop 1995). This principle arises from information theory. The maxent principle says that if we must completely specify a probability distribution with only partial knowledge about moments of that distribution, then we should choose the distribution that is most uninformative with regard to missing information. This means choosing the distribution that has maximal entropy subject to any information constraints we can justify, such as constraints based on data. Mathematically,

justifying (for example, constraints about the mean and/or distribution). The result is a distribution that conveys no information about what we explicitly build into it. If the only constraint was that the prior probabilities sum to 1, then the maximum entropy distribution is the uniform distribution.

That we can justify this uninformative prior for the maximum entropy theory is yet another example of how deeply information theory underlies statistical model selection. We will not discuss the merits of exploring information-theoretic data-based model selection, but we recommend that interested readers pursue some of the references here on the subject. An introductory reference that ties together the theory of statistics and information theory, including the Jaynes connection, is Jessop (1995). A nontechnical reference is Lucky (1992), from a practical perspective see Yockey (1992), while Cover and Thomas (1991) is a very thorough overview of information theory. Short, historical treatments are given by Wehrl (1978) and Ullah (1996).

As a general comment we emphasize the extensive use of information and entropy theory, and how these ideas underlie many scientific and technical areas (from Boltzmann to Einstein, and Kullback–Leibler, for example). There is thus a deep foundation for the Kullback–Leibler information measure and a firm basis for its use in many other aspects of statistics. K-L is not just another (of many) measures of discrepancies between probability distributions; it is a fundamental data-based model selection in science when truth is very “noisy,” and models can be only approximations to truth.

In saying that this theoretical foundation for use of K-L information we would liken it to the theoretical basis for the importance of the golden ratio (≈ 2.7183) in mathematics. It is not at all obvious why such a number should universally be the basis for logarithms in most of mathematics and science. But just as with K-L information, there is a compelling, deep reason, not easily perceived, for the importance of the golden ratio.

7.7.10 Akaike Weights w_i Versus Selection Probabilities

The model selection probabilities can be expressed as expected values of random variables that are a function of the sample data:

$$M_i(\underline{x}) = \begin{cases} 1 & \text{if model } i \text{ is selected by the data} \\ 0 & \text{otherwise.} \end{cases}$$

By definition, $E(M_i(\underline{x})) = \pi_i$. We assume no ties for the maximum likelihood estimate.

The Akaike weights (see Section 2.9) defined by

$$w_i = \frac{\exp(-\frac{1}{2}\Delta_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta_r)}$$

$0 < \gamma < \infty$,

$$w_k(\gamma) = \frac{1}{1 + \sum_{r \neq k} \exp(-\gamma \Delta_r)}$$

and

$$w_i(\gamma) = \frac{\exp(-\gamma \Delta_i)}{1 + \sum_{r \neq k} \exp(-\gamma \Delta_r)}, \quad i \neq k$$

In the limit as γ goes to infinity we have the result

$$\lim_{\gamma \rightarrow \infty} w_i(\gamma) = M_i(\underline{x}),$$

whence

$$\lim_{\gamma \rightarrow \infty} E(w_i(\gamma)) = \pi_i$$

(the implied interchange of limits will be valid here). There may or may not be the case that $E(w_i(0.5)) \equiv E(w_i) \neq \pi_i$; also, $E(w_i)$ may be unrelated.

This result does not rule out $E(w_i) \approx \pi_i$, which is sometimes a useful approximation. Moreover, use of the w_i as an estimator for the *set* of selection probabilities seems reasonable where such $\hat{\pi}_i$ are needed. (Research could be done to find out how well on the Akaike weights).

7.8 Kullback–Leibler Information Is Always Nonnegative

It is not obvious that the Kullback–Leibler discrepancy,

$$I(f, g) = \int f(x) \log \left(\frac{f(x)}{g(x)} \right) dx,$$

is strictly nonnegative for any possible $g(x)$. Here we reduce the model g to just $g(x)$ rather than $g(\underline{x} | \underline{\theta})$. Also, the possible nature of f and g is not emphasized in the proofs in this section.

Rigorous proofs exist that $I(f, g) \geq 0$ and that $I(f, g) = 0$ if $g(x) \equiv f(x)$ for all x . Here we give a valid, but not rigorous, proof that $I(f, g) \geq 0$. We do so for both the case of continuous and discrete distributions such as the Poisson, binomial, etc., wherein

$$I(f, g) = \sum_{i=1}^k p_i \log \left(\frac{p_i}{q_i} \right).$$

the q_1, \dots, q_k constitute the approximating probability model). Hence, here f and g correspond to the p_i and q_i .

In the first case, both $f(x)$ and $g(x)$ must be valid probability distributions, hence satisfy $f(x) \geq 0$, $g(x) \geq 0$ and both integrate to 1:

$$\int f(x)dx = 1, \quad \int g(x)dx = 1.$$

The exact limits of integration need not be specified here, the same for both f and g . Moreover, without loss of generality, $f(x) > 0$, $g(x) > 0$; hence the ratio $f(x)/g(x)$ is never undefined (but may be taken as ∞). For the discrete case we have $0 < q_i < 1$, and

$$\sum_{i=1}^k p_i = 1, \quad \sum_{i=1}^k q_i = 1.$$

We consider first the case of continuous probability distributions. One line of proof is to define a new function

$$h(x) = \frac{g(x) - f(x)}{f(x)};$$

thus,

$$\frac{g(x)}{f(x)} = 1 + h(x).$$

The lower bound on $h(x)$ is -1 , because for any x over which the integral is performed, $g(x)$ can be arbitrarily close to 0. The upper bound is thus $-1 < h(x) < \infty$. Note also that $\log(a) = -\log(1/a)$

$$\begin{aligned} I(f, g) &= \int f(x) \log \left(\frac{f(x)}{g(x)} \right) dx \\ &= - \int f(x) \log \left(\frac{g(x)}{f(x)} \right) dx \\ &= 0 - \int f(x) \log \left(\frac{g(x)}{f(x)} \right) dx \\ &= \int f(x) h(x) dx - \int f(x) \log \left(\frac{g(x)}{f(x)} \right) dx \end{aligned}$$

The last step above uses the fact that

$$\begin{aligned} 0 &= \int f(x) h(x) dx = \int f(x) \frac{g(x) - f(x)}{f(x)} dx \\ &= \int (g(x) - f(x)) dx \end{aligned}$$

Returning now to the main proof, we have

$$\begin{aligned}
 I(f, g) &= \int f(x)h(x)dx - \int f(x)\log\left(\frac{g(x)}{f(x)}\right)dx \\
 &= \int f(x)h(x)dx - \int f(x)\log(1 + h(x))dx \\
 &= \int f(x)[h(x) - \log(1 + h(x))]dx \\
 &= \int f(x)t(h(x))dx,
 \end{aligned}$$

where $t(h(x)) = h(x) - \log(1 + h(x))$. We do not need to know the values of $t(h(x))$. Nor do we need to consider $t(\cdot)$ as a function of x also, x may be univariate or multivariate. It suffices to consider $t(h)$ as a function of h , hence $t(h)$, over the full range of h , $-1 < h < \infty$, possible by varying x . All we care about is some basic aspects of $t(h)$, namely that it is strictly nonnegative. It is.

Calculus can be used to show that $t(h) \geq 0$, and that $t(h)$ has a minimum, and for any $h \neq 0$, then $t(h) > 0$. A simple heuristic is to plot $t(h)$ over, say, $-1 < h \leq 5$, and check $t(h)$ at a few values (Figure 7.1, and $t(10) = 7.6021$, $t(20) = 16.9555$, $t(1000) = 993.0913$). Given that $t(x) \geq 0$ for all x ,

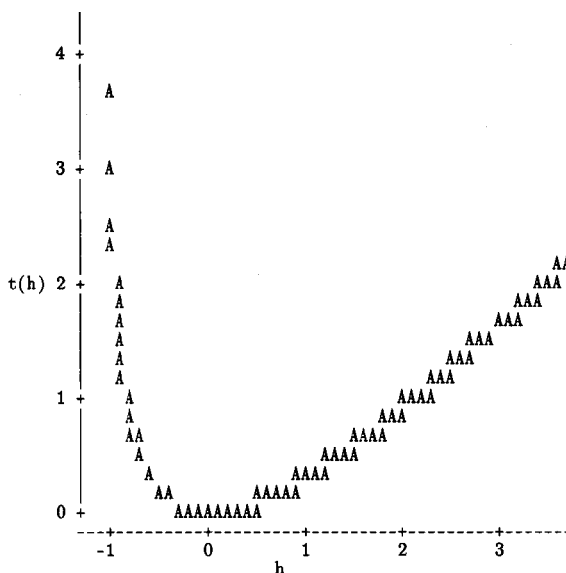


FIGURE 7.1. Plot of the function $t(h)$ near 0

The calculus proof that $t(x) \geq 0$ makes use of the first and second derivatives of the function $t(h) = h - \log(1 + h)$:

$$t'(h) = \frac{h}{1+h},$$

$$t''(h) = \frac{1}{(1+h)^2}.$$

The set of critical points (which includes minima, maxima, and saddle points) of $t(h)$ consists of the solutions to $t'(h) = 0$ plus the endpoints $h = -1$ and infinity. In this case the unique solution is $h = 0$. The fact that $h(x) = 0$ could occur for more than one value of x is not an issue. An extremum is deduced from $t''(0) = 1$, which, because it is positive, shows that $h = 0$ is a minimum of the function $t(h)$ (and by uniqueness it is the only minimum). Therefore, for all h (and hence for all x) $t(h(x)) \geq 0$. Also, from these results, $t(h(x))$ is a convex function.

Deeper mathematical theory is required to prove that $f(x) = g(x)$ for all x (in the relevant range of integration). If $f(x) = g(x)$, then $I(f, g) = 0$. Part of the “deeper” mathematics says that when $f(x)$ is a continuous probability density function, $t(x) \geq 0$, then

$$I(f, g) = \int f(x)t(x)dx = 0$$

if and only if $t(x) = 0$ for all x in the range of integration. This result seems reasonably intuitive, so we will not belabor the point.

$$h(x) - \log(1 + h(x)) \equiv 0, \quad \text{for all } x$$

or

$$h(x) = \log(1 + h(x)),$$

and finally,

$$e^{h(x)} = 1 + h(x).$$

The standard series expansion for e^h can be used here, with h replaced by $h(x)$:

$$1 + h(x) + \sum_{i=2}^{\infty} \frac{1}{i!} [h(x)]^i = 1 + h(x)$$

or

$$\sum_{i=2}^{\infty} \frac{1}{i!} [h(x)]^i = 0.$$

set of x values \mathcal{N} for which

$$\int_{\mathcal{N}} f(x) dx > 0,$$

then we would have

$$\int g(x) dx < \int f(x) dx,$$

which cannot be true. Thus, because both $f(x)$ and $g(x)$ are functions, $h(x) \leq 0$ for all x , which implies that we must have $h(x) = 0$ for all x .

Now we consider (in less detail) the discrete case

$$I(f, g) = \sum_{i=1}^k p_i \log \left(\frac{p_i}{q_i} \right),$$

such that $0 < p_i < 1$, $0 < q_i < 1$ for all i , and $\sum_{i=1}^k p_i = \sum_{i=1}^k q_i = 1$. For fixed k , this $I(f, g)$ is a function of $k - 1$ variables, with q_1, q_2, \dots, q_{k-1} (the p_i distribution is considered fixed here).

$$h_i = \frac{q_i - p_i}{p_i}, \quad i = 1, \dots, k,$$

whence

$$\frac{q_i}{p_i} = 1 + h_i, \quad -1 < h_i < \infty, \quad i = 1, \dots, k.$$

As in the continuous case,

$$\sum_{i=1}^k p_i h_i = \sum_{i=1}^k (q_i - p_i) = 0,$$

so we can derive

$$I(f, g) = \sum_{i=1}^k p_i (h_i - \log(1 + h_i)) = \sum_{i=1}^k p_i t(h_i).$$

It was proved above that $t(h) \geq 0$; thus it must be that even $I(f, g) \geq 0$.

It is clear that if $p_i = q_i$ for all i , then $I(f, g) = 0$. Also, if $h_i = 0$ for all i , then $I(f, g) = 0$. Then it must be that $t(h_i) = 0$ for all i (otherwise, $I(f, g) > 0$). Therefore, we must have

$$e^{h_i} = 1 + h_i, \quad i = 1, \dots, k.$$

The set of indices $\{1, \dots, k\}$ can be partitioned into two sets \mathcal{N} and \mathcal{P} . For i in \mathcal{N} , $h_i < 0$ and for i in \mathcal{P} , $h_i \geq 0$. For i in \mathcal{P} ,

Next, note that

$$\begin{aligned}
 0 &= \sum_{i=1}^k p_i h_i \\
 &= \sum_{i=1}^k (q_i - p_i) \\
 &= \sum_{i \in \mathcal{N}} (q_i - p_i) + \sum_{i \in \mathcal{P}} (q_i - p_i) \\
 &= \sum_{i \in \mathcal{N}} (q_i - p_i).
 \end{aligned}$$

But i in \mathcal{N} means $h_i < 0$, or $q_i - p_i < 0$, which would sum would be strictly < 0 , which is a contradiction. This that the set \mathcal{N} is empty: There cannot be any $h_i < 0$ if $I(f, g) = 0$, then $f \equiv g$ (i.e., $p_i = q_i$, for all i in the disc

7.9 Summary

Most of this chapter is quite technical; we will try to provide a summary of key points or results. Sections 7.1 through 7.6 provide the mathematical theory for K-L information-theoretic model selection. The general theory is given (Sections 7.1–7.3) along with several special cases (Sections 7.4–7.5) and some specific explorations (Section 7.6). In particular, a very detailed derivation of the K-L divergence (Section 7.2, along with the relationship to AIC. Then Section 7.7 provides a (not exhaustive) examination of the issue of whether we can use AIC_c rather than TIC; the results strongly support use of AIC_c as acceptable, relative to TIC, but actually preferable. Section 7.8 provides simple (though not general) methods to explore model selection that are operationally based on the log-likelihood. Theoretical results are also given. Sections 7.7.6–7.7.10 give specific results and considerations that do not fit elsewhere. Section 7.9 is less technical than Sections 7.1–7.6, and we urge you to read it for the general insights therein. Section 7.8 gives a proof of the Kullback–Leibler inequality.

There are several rigorous derivations from Kullback–Leibler divergence leading to various information-theoretic criteria: The most common derivation leads to Takeuchi’s (1976) information criterion. The derivation of AIC_c is given in detail in Section 7.4. It is also the unique small-sample version of AIC, but AIC_c is recommended (that could change in the future, especially for discrete distributions). Section 7.5 gives a derivation of AIC for the exponential family of distributions.

from K-L information is much more exact, relies on fewer assumptions, and is easier to understand.

The fact that such derivations exist is important to know. The detailed explanations are very detailed because the theory underlying model selection based on K-L information is important to have clearly stated and understood. Such understanding of the theory puts one in a better position to accept use of the information-theoretic criteria and understand their strengths and weaknesses.

While Kullback–Leibler information is the logical basis for model selection, it turns out we must use *expected* (over the true model) K-L information as the quantity of interest when model parameters are estimated. This, of course, is the reality of actual data analysis. The bias–variance tradeoff and the associated model parsimony achieved by AIC for model selection is an important byproduct of the approach. The detailed derivations make it clear that K-L–based model selection does not have an explicit objective of meeting the principle of parsimony. Rather, it is a natural consequence of data-based K-L model selection. The bias–variance tradeoff happens. In fact, it is because the model parameters, must minimize *expected* K-L information. The cross-validation property of AIC.

The detailed derivations make it clear that use of information-theoretic criteria in the analysis of real data is not based on the existence of a true model or the notion that such a true model is in the set of candidate models. Literature contrary to this point is mistaken.

Model selection attempts to establish some rigorous basis for choosing the most parsimonious model(s) used for inference. The relationship between AIC and TIC is made clear, and investigations were undertaken to show that AIC is a good proxy for TIC. It seems poetic that AIC can be thought of as a simplified implementation of the more general TIC. The trace term in AIC is about equal to K for “good” models and does not depend on sample size if sample size is large. Some insights are provided to help understand the relationship between $\text{tr}[J(\theta_o)I(\theta_o)^{-1}]$ and K in a variety of situations. Evaluations were conducted for logistic regression, multinomial regression, count data, and normal regression models. In all cases we found the trace term of TIC is very close to being K as long as the model structure and error distribution are not drastically different from truth. If the model is not the true data-generating model, the trace term was not equal to K or $< K$. Rather, the matter is unpredictable; the model could be more or less general than the true model and still the trace term can be any of $= K$, $> K$, or $< K$. When examined, however, if the model was less general than the true model (the worst case), we predominantly found $\text{tr}[J(\theta_o)I(\theta_o)^{-1}] < K$. Thus, models that are less general than the true model then often lead to slightly more parsimonious models.

If the set of models contains one or more good models, will tend not to be selected (or even ranked high), because the value for that model, hence rendering the issue of use of TIC moot. As noted above, for good models use of AIC is acceptable for small samples and even for large samples if values of K are small relative to sample size. More research on such second-order improvements, especially for discrete random variables.

Monte Carlo methods seem to be the only tool to assess the relative merits of aspects of model selection and methods to incorporate model uncertainty. In some cases, asymptotic results can be obtained, but they are of little interest or practical use. We present some quick results on model selection using Monte Carlo simulation in the case of a large number of models. For that same context, the theoretical expectations of model selection can be easily found, and this is explored in some detail in the case of models with normal errors. The issue of AIC overfitting is also explored. Extreme overfitting can occur, but the probability of this occurring and one way to minimize the problem is to keep the set of models small. Doing searches over “all possible models” (e.g., all subsets of K variables) increases the risk of overfitting. In linear regression it seems that model selection is very similar to model selection based on minimum variance unbiased (course, in the analysis of real data we cannot do selection based on theoretical MSE).

Summary

This book covers some philosophy about data analysis, the interface between mathematical statistics and information theory, practical statistical methodology useful in the applied sciences, and we present a general strategy for modeling and data analysis. We give challenging examples from our fields of interest, provide a list of things not to do, and suggest some areas needing further theoretical work. We side with the fast-growing ranks that see limited utility in frequentist hypothesis testing. Finally, we provide references for those wishing to study further.

Conceptually, there is information in the observed data. We try to express this information in a compact form via a “model.” The model represents a scientific hypothesis and is then a basis for making predictions about the process or system that generated the data. One can view the information in data as a change in “coding” like a change in the way a thought or emotion expressed in one language (e.g., French) looks when expressed in another language (e.g., Russian). A good model has only a finite, fixed amount of information. The (unachievable) goal is to attain a perfect 1-to-1 translation such that no information is lost in going from the data to a model of the information. Models are only approximations, and we cannot hope to perfectly achieve the idealized goal. However, we can attempt to find a model of the data in the sense that the model loses as little information as possible. This leads directly to Kullback–Leibler information $I(f, g)$: the amount of information lost when model g is used to approximate full reality f . We choose the model that minimizes K-L information loss. Because we

information-theoretic criteria (e.g., AIC , AIC_c , $QAIC$, $QAIC_c$), the model allows the efficient and objective separation or filtering from *noise*. In an important sense, we are not really trying to model the *information* in the data; instead, we are trying to model the *information* in the data.

While we use the notation f to represent truth or the existence of a “true model” in the life sciences, be the process (truth) that generates the sample data we to make inferences about truth, while realizing that full truth be beyond us when we have only sample data. Data analysis is thought of as an attempt to identify f ; instead, we must use good approximations to truth and from which therefore make inferences concerning truth. We do not want merely to develop a model that has a very large number of parameters; instead, we use data to aid in the selection of a parsimonious model that allows to be made about the system or process under study. A process representing a well-defended scientific hypothesis, aids in the of the system of interest.

Relatively few statistics books provide a summary of the yet fewer provide an effective, unified strategy for data analysis where there is substantial complexity. The breadth of the covered here makes a summary difficult to write. Undergraduate occasionally ask the professor, “What is important for me to examination?” The professor is typically irritated by such the student should realize that it is *all* important! Indeed, Akaike’s pioneering work is that it *is* all important. The information paradigm is a *package*; each of the package’s contents is but it is the integration of the contents that makes for an effective consistent strategy, and a practical and powerful methodology. package that has been so frequently left out is the critical generation, and modeling *before* examination of the data; thinking should occur prior even to data collection. This the science of the issue formally enters the overall “analysis” (Burnham 1999a).

The information-theoretic methods we present can be used to best model that can be used in making inferences from error often portrayed in the literature in this simple manner. This is much richer than this simplistic portrayal of model selection. **In fact, an emphasis of this second edition is multimodel selection. MMI has several advantages; all relate to the broad selection uncertainty.** One can easily rank alternative models from best to worst using the convenient differences Δ_i for each model, given the data [i.e., $\mathcal{L}(g_i \mid \text{data})$], can be easily

a subset of good models. Evidence ratios are useful for support of one model versus another, given the data; but irrespective of other models in the set.

Model selection uncertainty can be easily quantified using the bootstrap (the bootstrap is an alternative). Estimates of this component can be incorporated into unconditional estimates of precision and bias. For many problems (e.g., prediction) model-averaging methods are available. We treat this important issue in Chapters 4–5. Thus, we formal inference from all models in the set.

For those who have scanned through the pages of this book, it may be a surprise at the general lack of mathematics and formulas (Chapters 6 and 7 being the exceptions). That has been our intent. The underlying information-theoretic methods is relatively simple. They are easy to understand and use (“low tech”), while the underlying theory is complex (Chapter 7). As we wrote the book and tried to understand the underlying papers (see Parzen et al. 1998) we found the need to delve into the philosophical aspects that are generally philosophical. The science of the problem is to get into modeling *before* one begins to rummage through the data. In some critical respects, applied statistics courses are far from treating statistics as an integral part of scientific discovery, with little emphasis on model selection methods or their importance, while successful teaching null hypothesis testing methods and data analysis often rely on the assumption that the model is both true and given. See Anderson et al. 2001b, note, “The standard approach in teaching—stressing the use of a p value while warning against its misinterpretation—has been an abysmal failure.” It seems necessary to greatly reduce the reliance on p values (Anderson et al. 2001b and d).

8.1 The Scientific Question and the Collection of Data

The formulation of the research question is crucial in investigating complex systems and processes in the life sciences. A good research question is a mistake all too often seen in the published literature. It is better to have a poor answer to a poor question. Investigators should readdress the importance and quality of the question to be answered. Scientific hypotheses, represented by models, must have a clear place on the table.

A careful program of data collection must follow from the formulation of the research question. Particular attention should be placed on the variables that are of interest and interesting covariates. Observational studies, done well,

tors are controlled or balanced, experimental units are randomized to treatment and control groups with adequate replication), al. (1980), Gail (1996), Beyers (1998), and Glymour (1998) are different philosophies. Valid inference must assume that these basic principles have been carefully planned and conducted. Before one can answer general questions must be answered in the affirmative:

*Are the study objectives sound, relevant, and achievable?
Has there been proper attention to study design and laboratory methods?*

8.2 Actual Thinking and A Priori Modeling

Fitting models, each representing a scientific hypothesis, is an important activity in many biological, ecological, and medical studies. Statistical inferences about the system of interest are made from a parsimonious model of the observational or experimental data. We see this activity increase as more complicated scientific and medical problems are addressed. In particular, a priori modeling becomes important as several data sets are collected on the same issue by different investigators at widely differing field sites over several years.

We recommend much more emphasis on thinking! Leave time for a while, giving time to think hard about the overall problem. What information is contained in the published literature, even on issues not related to the issue at hand? What nonlinearities are likely to be predicted? What interactions are hypothesized to be important? Two or more variables be combined to give a more meaningful analysis? Should some variables be dropped from consideration? Should some variables be encouraged with the people in the field or laboratory? What parameters might be similar across different data sets? Model building should be driven by the underlying problem combined with a good understanding of mathematical modeling. An important conceptual phase might take several days or even weeks. This seems far more time than is often spent under current conditions.

Biologists generally subscribe to the philosophy of "a priori hypotheses" (Chamberlain 1890, Platt 1964, Mayr 1996). This philosophy forms the basis for the set of candidate models to be considered. Model building can begin during the time that the a priori hypotheses are being sorted out. Modeling must carefully quantify the scientific problem of interest. Often it is effective to begin with the global model and then consider some lower-dimensional models. Others may favor a bottom-up approach. The critical matter here is that one arrives, eventually, at a model that is both

however, every attempt should be made to keep the number small.

Critical Thinking

Our science culture does not do enough to regularly encourage critical thinking. This failure has slowed the scientific progress.

We fail to fault the trivial content of the typical ecology paper.

There is a need for more *careful thinking* (than is usual) and a *better balance* between scientific hypotheses, data, and analysis.

Chamberlin's concept of *multiple working hypotheses*, introduced 100 years ago, has a deep level of support among scientists. It is a thought the method led to "certain distinctive habits of mind" of value in education." Why has this principle not become more common than the rare exception, in so many fields of applied science?

Platt (1964) noted that years and decades can be wasted in science unless one thinks carefully in advance about what the results of conclusive experiments would be.

With the information-theoretic approach, there is no hypothesis, or a statistical hypothesis test, or an arbitrary significance level, or testable power, or the multiple testing problem, or the fact that a hypothesis is nearly always *obviously* false in the first place. The application of statistical hypothesis testing arbitrarily classifies results into meaningless categories of "significant" and "nonsignificant." This has little to contribute to the advancement of science (Arnold 2008). We recommend that researchers stop using the term "significant," which is overused, uninformative, and misleading. The results of model selection on estimates of expected (relative) Kullback–Leibler information are different from the results of some form of statistical hypothesis testing. See the simulated starling data, Section 3.4, or the sage grouse example.

So, investigators may proceed with inferential or confirmatory statistics if they feel satisfied that they can objectively address two questions:

Was the set of candidate models derived a priori?

What justifies this set?

The justification should include a rationale for models excluded from the set. A carefully defined set of models is required. If information-theoretic methods are used to select the single best model, the entire set of models is used to reach defensible inferences. If no known about the system under study that a large number of models are included in the candidate set, then the analysis should proceed only exploratory (if models are developed as data analysis is exploratory and risky). **One should check the fit or adequacy of the global model using standard methods. If the global model is**

tem under study and the data collected. There is no hard thinking at this point (Platt 1964).

8.3 The Basis for Objective Model Selection

Statistical inference from a data set, *given a model*, is supported by a very large amount of theory. Theorists and practitioners employ this theory, either likelihood or least squares, to solve a wide variety of problems in the applied sciences. The most compelling question is “*Which model to use?*” Valid inference must usually be based on a model, but which one?

Akaike chose the celebrated Kullback–Leibler distance as a basis for model selection. This is a fundamental quantity that has earlier roots in Boltzmann’s concept of *entropy*, a cornerstone of nineteenth-century science. The K-L distance between two models f and g is defined for continuous functions as the integral

$$I(f, g) = \int f(x) \log \left(\frac{f(x)}{g(x | \theta)} \right) dx,$$

where \log denotes the natural logarithm and f and g are probability distributions. Kullback and Leibler (1951) developed “information theory,” thus the notation $I(f, g)$ as it relates to “information” lost when model g is used to approximate truth f . An approximating model that loses as little information as possible is equivalent to minimizing $I(f, g)$ over the models in the set considered to be fixed. An interpretation equivalent to minimizing $I(f, g)$ is that we seek an approximating model that is the “shortest distance” from the truth. Both interpretations seem useful and compelling.

The K-L distance can be written equivalently as

$$I(f, g) = \int f(x) \log(f(x)) dx - \int f(x) \log(g(x | \theta)) dx$$

The two terms on the right in the above expression are expectations with respect to f (truth). Thus, the K-L distance (above) is a difference between two expectations,

$$I(f, g) = E_f[\log(f(x))] - E_f[\log(g(x | \theta))]$$

each with respect to the true distribution f . The first expectation is a constant that depends only on the unknown true distribution. Treating this unknown term as a constant, only a measure of the distance between the two models is needed.

or

$$I(f, g) - \text{constant} = -E_f[\log(g(x | \theta))]$$

Thus, the term $(I(f, g) - \text{constant})$ is a *relative* distance between model f and model g . This provides a deep theoretical basis for model selection: compute or estimate $E_f[\log(g(x | \theta))]$.

Akaike (1973, 1974, 1985, 1994) showed that the relative K-L information was

$$E_y E_x [\log(g(x | \hat{\theta}(y)))],$$

where y and x are independent random samples from the target distribution, and both statistical expectations are taken with respect to the target distribution. f , is the target distribution, both with respect to truth f , is the target distribution. approaches based on K-L information.

8.4 The Principle of Parsimony

Parsimony is the concept that a model should be as simple as possible with respect to the included variables, model structure, and number of parameters. Parsimony is a desired characteristic of a model used for model selection, usually visualized as a suitable tradeoff between square error and the number of parameter estimators (Figure 1.3). Parsimony lies between underfitting and overfitting (Forster and Sober 1994, Forster 1999). K-L information is a fundamental basis for achieving parsimony in modeling.

The concept of parsimony has a long history in the scientific method, expressed as “Occam’s razor”: shave away all that is unnecessary. It is to make things “as simple or small as possible.” Parsimony represents a tradeoff between bias and variance as a function of the model (K). A good model is a proper balance between underfitting and overfitting, given a particular sample size (n). Most model selection methods are based on the concept of a squared bias versus variance tradeoff. The selection of a model from a set of approximating models must emphasize parsimony. These philosophical issues are stressed in the scientific method, some experience and reconsideration to reach a full understanding of importance.

Roots of Theory

As deLeeuw (1992) noted, Akaike found a form (between Boltzmann's entropy and Kullback–Leibler information paradigms in information and coding theory) and maximization (the dominant paradigm in statistics).

This finding makes it possible to combine estimation (maximum likelihood estimation) and model selection under a single theoretical optimization.

Akaike's (1973) breakthrough was the finding of an estimator of relative K-L information, based on a bias-corrected maximum likelihood value. His estimator was an approximation and, under certain conditions, asymptotically unbiased. He found that

$$\text{estimated expected (relative) K-L information} \approx 1$$

where $\log(\mathcal{L}(\hat{\theta}))$ is the maximized log-likelihood value and θ is the vector of estimable parameters in the approximating model (this is the log-likelihood term). Akaike multiplied through by -2 and provided the Akaike information criterion (AIC)

$$\text{AIC} = -2 \log(\mathcal{L}(\hat{\theta})) + 2K.$$

Akaike considered his information-theoretic criterion an extension of maximum likelihood theory. Conceptually, the principle of parsimony was added "penalty" (i.e., $2K$) while minimizing AIC.

Assuming that a set of a priori candidate models has been specified, then AIC is computed for each of the approximating models. The model where AIC is minimized is selected as best for the empirical data. This is a simple, compelling concept, based on deep theoretical foundations (i.e., K-L information). Given a focus on a priori issues, model selection, scientific hypotheses, and model selection, *the inference is in the model*. In a sense, parameter estimates are almost byproducts of model selection. This inference relates to the estimated best approximation of the true model. Information seems to be contained in the data.

Important refinements followed shortly after the pioneer work. Most relevant was Takeuchi's (1976) information criterion which provided an asymptotically unbiased estimate of relative K-L information. TIC is little used, since it requires the estimation of first and second partial derivatives of the log-likelihood function. Its practical use hinges on the availability of a relatively large sample size. AIC can be viewed as a parsimonious version of TIC. Akaike's work was motivated by Sugiura's (1978) work, and resulted in the work of Hurvich and Tsai (1989, 1990b, 1991, 1994, 1995a and

$$\text{AIC}_c = -2 \log(\mathcal{L}(\hat{\theta})) + 2K + \frac{2K(K+1)}{(n-K-1)}$$

where n is sample size. The final bias-correction term vanishes as n increases, so that AIC_c approaches AIC with respect to K (and AIC_c becomes AIC), but the addition of the bias-correction term is important if n is not large relative to K (we suggest using AIC_c as a default criterion, alternatively, always using AIC_c).

A third extension was a simple modification to AIC and AIC_c for overdispersed count data (Lebreton et al. 1992). A variance-to-mean ratio is computed from the goodness-of-fit statistic, divided by the degrees of freedom, $\hat{c} = \chi^2 / \text{df}$. The value of the maximized log-likelihood is then divided by the estimate of overdispersion to provide a profile log-likelihood. These criteria are denoted by QAIC and QAIC_c, and derived from quasi-likelihood theory (Wedderburn 1974):

$$\text{QAIC} = -[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c}] + 2K,$$

and

$$\begin{aligned} \text{QAIC}_c &= -[2 \log(\mathcal{L}(\hat{\theta})) / \hat{c}] + 2K + \frac{2K(K+1)}{n-K-1} \\ &= \text{QAIC} + \frac{2K(K+1)}{n-K-1}. \end{aligned}$$

When no overdispersion exists, $c = 1$, and the formulas for QAIC and QAIC_c reduce to AIC and AIC_c, respectively. There are other, more sophisticated methods to account for overdispersion in count data, but this simple method is quite satisfactory. Methods are given in Chapter 6 to allow for the possibility of the data to have partition-specific estimates of overdispersion. The number of estimable parameters (K) must include the number of estimates of c . Thus, if males and females have different degrees of overdispersion, these are to be estimated from the data, then K must include the number of these estimates.

AIC is often presented in the scientific literature in an oversimplified way, the bias-correction term K (the so-called penalty term) is omitted, yet, perhaps, is that AIC is often given without reference to its theoretical basis with Kullback–Leibler information. Such shallow presentations have had very negative effects, and have misled many into thinking that AIC is a whole class of selection criteria that are “information-theoretic” (Burnham 1996). Criteria such as AIC, AIC_c, QAIC, and TIC are essentially different estimates of (relative) Kullback–Leibler distance and are useful in the “noisy” sciences.

Because only relative K-L information can be estimated, for information criteria, it is convenient to rescale these values so that the model with the minimum AIC (or AIC_c or TIC) has a value of 0. The criterion values can be rescaled as simple differences,

$$\begin{aligned}\Delta_i &= AIC_i - AIC_{min} \\ &= \hat{E}_{\hat{\theta}}[\hat{I}(f, g_i)] - \min \hat{E}_{\hat{\theta}}[\hat{I}(f, g_i)]\end{aligned}$$

While the value of minimum $\hat{E}_{\hat{\theta}}[\hat{I}(f, g_i)]$ is not known (only if we have an estimate of the size of the increments of information for the various models compared to the estimated best model, then the minimum $E_{\hat{\theta}}[\hat{I}(f, g_i)]$). The Δ_i values are easy to compute and allow a quick comparison and ranking of candidate models. Models with Δ_i in computing Akaike weights. As a rough rule of thumb, models with Δ_i within 1–2 of the best model have substantial support and deserve consideration in making inferences. Models having Δ_i between 4–7 of the best model have considerably less support, and models with $\Delta_i > 10$ have either essentially no support and might not deserve further consideration or at least fail to explain some systematic variation in the data. If the observations are not independent (e.g., as such) or if the sample size is quite small, or if there is a large number of models, then the simple guidelines above cannot be taken too literally.

There are cases where a model with $\Delta_i > 10$ might still be of interest, particularly if the sample size is very large (e.g., see Section 3.5.5). For example, let model *A*, with year-specific structure on one of the parameters, be the best model in the set ($\Delta_A = 0$) and model *B*, with less structure, be the second best model. If models *A* and *B* have year-specific parameters, have $\Delta_B = 11.4$. Assume that all candidate models were derived prior to data analysis (i.e., no data-driven model selection). Model *A* is able to identify important variation in a parameter that is important. However, in terms of understanding and generalizing the results based on the data, it might sometimes be justified to use model *B*, because it may seem to “capture” the important fixed effects. Models *A* and *B* should both be detailed in any resulting publication. The understanding and interpretation might be enhanced using model *B* if the information in the data would be (intentionally) lost. Such information could be partially recovered by, for example, using a random effect model (see Section 3.5.5) to estimate the mean of the time-effect distribution and the variance of its distribution.

The principle of parsimony provides a philosophical basis for model selection; Kullback–Leibler information provides an objective measure of information; a fundamental theory; and the information criteria (particularly AIC) provide a practical, general methodology for use in data analysis.

in the set of candidates (although this is sometimes stated in the technical literature). [We note that several “dimension-consistent” estimates of the dimension (K) of the “true model” have been published that attempt to provide asymptotically unbiased estimates of the dimension (K) of the “true model” but only estimates of K-L information in a strained way, are based on different assumption sets, and often perform poorly (even toward the true model) unless a very large sample size is available (or when the sample size is small, such as in many problems in the physical sciences). We recommend these dimension-consistent criteria for the analysis of life sciences.]

8.7 Scaling Alternative Models

The information-theoretic approach does more than merely select the model that is best for making inference, given the set of a priori models and the data. The Δ_i allow a ranking of the models from most plausible to the worst; the larger the Δ_i , the less plausible is model i . **not reasonable to expect to be able to make inferences from a single model; biology is not simple; why should we hope for a single model? The information-theoretic paradigm is a paradigm for examination of alternative models and, where appropriate, for inference from more than one model (MMI).**

The simple transformation $\exp(-\frac{1}{2}\Delta_i)$ results in the (diagnostic) likelihood of model i , given the data $\mathcal{L}(g_i|x)$. These are functions in θ ; $\mathcal{L}(\theta|x, g_i)$ is the likelihood of the parameters θ , given the data x and model (g_i). These likelihoods are very useful; for example, the ratio of the likelihood for model i versus model j is merely

$$\mathcal{L}(g_i|x)/\mathcal{L}(g_j|x).$$

It is convenient to normalize these likelihoods such that the

$$w_i = \frac{\exp(-\frac{1}{2}\Delta_i)}{\sum_{r=1}^R \exp(-\frac{1}{2}\Delta_r)},$$

and interpret these as a weight of evidence. Akaike (e.g., 1973, 1974, 1979, 1980, and 1981b; also see Kishino 1991 and Buckland et al. 1999) has used these values, and we have found them to be simple and useful. The weight of evidence ratio of model i versus model j is then just w_i/w_j , the ratio of the likelihood $\mathcal{L}(g_i|x)/\mathcal{L}(g_j|x)$. Drawing on Burnham and Anderson (1998), we interpret w_i as the estimated probability that model i is the best model, given the data at hand, given the set of models considered (see

through here was to consider priors on models that are π_i and K (we call this class of model priors “savvy,” i.e., shrinkage), then AIC and AIC_c fall out as a strictly Bayesian result. In the Bayesian derivative, it is compelling to interpret the Akaike prior model probabilities. While many (objective) Bayesians object with the use of a diffuse or noninformative prior on model parameters (a uniform prior on a model parameter), use of such diffuse priors (such as $1/R$) may have poor properties or unintended consequences. Some priors on models may be uninformative, but not in the Bayesian derivation of AIC (or AIC_c) and BIC differ on models. However, these criteria are fundamentally different in substantive ways. In this book we place an emphasis on the use of AIC_c as bias-corrected estimates of Kullback–Leibler information. This seems so much more objective and fundamental.

The w_i are useful as the “weight of evidence” in favor of model i as being the actual K–L best model in the set. The bigger the Δ_i , the smaller the weight, the less plausible is model i as being the best approximation to the true model, conditional on both the data and the set of a priori models.

Alternatively, one could draw B bootstrap samples (B should be at least 1,000), use the appropriate information criterion to select a model for each of the B samples, and tally the proportion of the i th model was selected. Denote such bootstrap-selected model $\hat{\pi}_i$. While w_i and $\hat{\pi}_i$ are not estimates of exactly the same thing, they are closely related and provide information concerning the uncertainty of model for use. The Akaike weights are simple to compute, while bootstrap weights are computer-intensive and not practical to compute for large B (e.g., the simulated starling experiment, Section 3.4), but bootstrap repetitions must be drawn and analyzed.

Under the hypothesis-testing approach, nothing can be said about ranking or scaling models, particularly if the models were nested. For least squares problems one could turn to adjusted R^2 values to compare models, but other kinds of models cannot be scaled up or down (a very poor) approach (see the analogy in Section 2.5).

8.8 MMI: Inference Based on Model Averaging

Rather than base inferences on a single selected best model from a set of models, we can base our inferences on the entire set of models, averaging. The key to this inference methodology is the Akaike weights if a parameter θ is in common over all models (as θ_i in model i). If we are predicting, by using the weighted average we are basing

or

$$\hat{\theta} = \sum_{i=1}^R \hat{\pi}_i \hat{\theta}_i.$$

This approach has both practical and philosophical advantages. If an averaged estimator can be used, it appears to have better properties than $\hat{\theta}$ from the selected best model.

If one has a large number of closely related models, such as those from stepwise or based variable selection (all-subsets selection), designating a single best model is unsatisfactory, because that estimated “best” model can change from data set to data set. In this situation model-averaging seems to provide much more stabilized inference. The concept of inference averaging over the models can be used to reduce model selection bias in estimating the coefficient estimates in all-subsets selection. For the regression coefficient associated with predictor x_j we use the estimate $\hat{\beta}_j$, where $\hat{\beta}_j$ is the regression coefficient β_j averaged over all models in which x_j appears.

$$\hat{\beta}_j = \frac{\sum_{i=1}^R w_i I_j(g_i) \hat{\beta}_{j,i}}{w_+(j)},$$

$$w_+(j) = \sum_{i=1}^R w_i I_j(g_i),$$

where i is for model $i = 1, \dots, R$, j is for predictor variable x_j .

$$I_j(g_i) = \begin{cases} 1 & \text{if predictor } x_j \text{ is in model } g_i \\ 0 & \text{otherwise.} \end{cases}$$

Conditional on model g_i being selected, model selection biasing $\hat{\beta}_{j,i}$ away from zero. Thus a new estimator, denoted $\tilde{\beta}_j$, is

$$\tilde{\beta}_j = \sum_{i=1}^R w_i \hat{\beta}_{j,i}.$$

Investigation of this idea, and extensions of it, is an open research topic. One point here is that while $\hat{\beta}_j$ can be computed ignoring models in which x_j does not appear, $\tilde{\beta}_j$ does require fitting all R of the models.

8.9 MMI: Model Selection Uncertainty

At first, one might think that one could use an information criterion to select an approximating model that was “close” to truth (remember that we are assuming that

been specified a priori as the only model considered. This would not be terrible, since at least one would have a selected objectively, based on a valid theory and a priori approach would often be superior to much of current practice where the best model has an Akaike weight > 0 . considering only this model, and the usual measures of precision on this selected model, is that this tends to overestimate (1992) calls the failure to acknowledge model selection scandal.” [We might suggest that the widespread use of selection testing and blatant data dredging in model selection represents In fact, there is a variance component due to model selection should be incorporated into estimates of precision such as conditional (on the selected model). While this is a research development, several useful methods are suggested in the will surely appear in the technical literature in the next few additional Bayesian approaches.

The Akaike (w_i) or bootstrap (π_i) weights that are used for models can also be used to estimate unconditional precision in the parameter θ over R models (model g_i , for $i = 1, \dots, R$).

$$\widehat{\text{var}}(\hat{\theta}_i) = \left[\sum_{i=1}^R w_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]$$

$$\widehat{\text{var}}(\hat{\theta}_i) = \left[\sum_{i=1}^R \pi_i \sqrt{\widehat{\text{var}}(\hat{\theta}_i | g_i) + (\hat{\theta}_i - \hat{\theta})^2} \right]$$

These estimators, from Buckland et al. (1997), include an additional sampling variance, given model g_i (denoted by $\widehat{\text{var}}(\hat{\theta}_i | g_i)$) incorporate a variance component for model selection uncertainty. These estimators of unconditional variance are also appropriate when one wants a model-averaged estimate of the parameter over R models.

Chapter 4 gives some procedures for setting confidence intervals for model selection uncertainty, and it is noted that achieved coverage is then a useful measure of the utility of methods for incorporating selection uncertainty into inference. Only a limited aspect of this can be currently handled. *Given* a set of candidate models and a selection method, we can assess selection uncertainty. The uncertainty in the set of models cannot be addressed; we lack a theory for this. We lack good, general guidelines for defining the a priori set of candidate papers to appear on these scientific and philosophical issues.

inference on the importance of a variable is similarly important on all the models. If one selects the best model and says that are the important ones and the other variables are not important, naive, unreliable inference. We suggest that the relative importance x_j be measured by the sum of the Akaike weights over all models that that variable appears:

$$w_+(j) = \sum_{i=1}^R w_i I_j(g_i).$$

Thus again, proper inference requires fitting all the models of a given type of model-averaging. A certain balance in the number of models j , must be achieved. When possible, one should average on all the models, via model-averaging and selection bias reduction, rather than risk making inference based only on the model estimated, and, often, ignoring other models that are also quite good.

8.11 More on Inferences

Information-theoretic methods do not offer a mechanical, recipe for science. While these methods can certainly be misused, thinking as models are developed to represent the multiple hypotheses that must be the focus of the entire study. A central theme of this book is attention to the need to ask better scientific questions in the first place (Platt 1964). Rather than test trivial null hypotheses, it is better to ask questions relating to well-defined alternative hypotheses. When achieved, a great deal more hard thinking will be required.

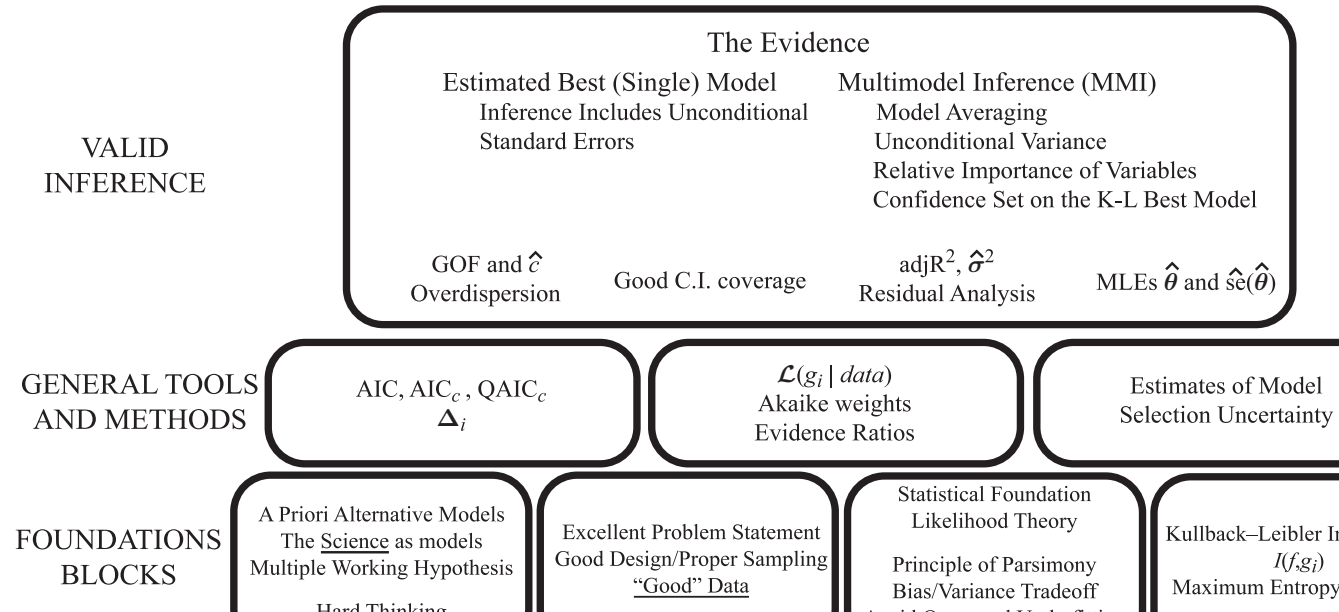
There needs to be increased attention to separating those hypotheses on a priori considerations from those resulting from some other consideration. White (2000:1097) comments, “Data snooping is a disease that should be avoided, but in fact is endemic.”

Essentially no justifiable theory exists to estimate precise probabilities (for those still so inclined) when data dredging has taken place (mis)used is for a priori analyses, assuming that the model selected is the best (to the data). A major concern here is the finding of effects and hypotheses that are actually spurious where inferences are made post hoc (Anderson et al. 2001b). This glaring fact is either not understood by reviewers and journal editors or is simply ignored. Two typical approaches include (1) an iterative approach, in which patterns and hypotheses after initial analysis are “chased” by repeatedly building new models with effects included and (2) analysis of “all possible models.” The first is a poor approach to making inferences about the sampled population.

The information theoretic paradigm avoids statistical testing concepts and focuses on relationships of variables (vs. the estimation of effect size and measures of its precision) primarily in the context of making inferences from a single model or making robust inference from many models (e.g., using model averaging based on Akaike weights). Data analysis is a process of making inferences supported by the data and the degree of complexity of the model set. Often, models other than just the estimated best model are reported to provide information. Evidence ratios and confidence sets on model selection are inferences on all, or several of the best, models in the set. In general, these approaches should not be used unthinkingly; a good set of models is essential, and this involves professional judgment and the integration of scientific hypotheses into the model set.

When the analysis of data has been completed using the information theoretic approach, one should gather and report on the totality of evidence at hand. The primary evidence might be the selected model, the model estimates and appropriate measures of precision (including evidence ratios and Akaike weights) (including a component for model selection uncertainty.) The ranks of each model, the Akaike weights should be reported and interpreted. The parameter estimates are often important, particularly for parameter estimates, ratios, confidence sets on the K-L best model, and a ranking of the importance of predictor variables are often useful evidence. Quantities such as adjusted R^2 and $\hat{\sigma}^2$ should be reported for the best model. The results from an analysis of residuals for the selected model should also be important to report and interpret. Every effort should be made to report and objectively report on all the evidence available. If some post hoc activities, this should be clearly stated in publication. Figure 8.1 provides a simplistic graphical representation of the information theoretic approach. The point of Figure 8.1 is to reinforce some of the foundational (bottom building blocks) and the practical tools and methods (top building blocks) that rest on these foundations. If these are used thoughtfully, one can hope to provide compelling evidence allowing for a better understanding of the science problem and the alternative hypotheses.

It seems worth noting that K-L information and MMI can be used in a variety of types of conflict resolution where data exist that are central to the resolution of the conflict (Anderson et al. 1999, 2001c). While this may take us too far afield; however, as Hoeting et al. (1999) note in a similar context, "Model averaging also allows users to incorporate uncertainty from many models in the estimation process; thus model averaging may provide a better estimation method than the traditional method of selecting a best model of scientists a better estimation method than the traditional method of selecting a best model to get the committee to agree on a best model."



At a conceptual level, reasonable data and a good model of “information” from “noise.” Here, information relationships, estimates of model parameters, and computation. Noise then refers to the residuals; variation left unexplained information extracted from the data to make proper inference.

Summary

We want an approximating model that minimizes information and properly separates noise (noninformation, or entropy) from information. The philosophy for this separation is the principle of maximum entropy. The conceptual target for such partitioning is Kullback–Leibler information, and the tactic for selection of a best model is an information-theoretic approach (AIC, AIC_c , $QAIC_c$, or TIC). The notion of data-based model selection resulting inference is a very difficult subject, but we do know that uncertainty about the selected model can often be expected and incorporated into estimates of precision.

Still, model selection (in the sense of parsimony) is the central part of the analysis. In using the more advanced methods presented here, model selection can be thought of as a way to compute Akaike weights. The more models in the set as a way to make robust inferences from the data. More research is needed on the quantification of model uncertainty, the plausibility of alternative models, ways to reduce model uncertainty, and ways to provide effective measures of precision (with confidence intervals on a given model). Confidence intervals with good coverage are a goal of inference following data-based model selection.

Information-theoretic methods are relatively simple to use and are applicable to employ across a very wide class of empirical situations in many disciplines. The information-theoretic approach unifies parameter estimation and model selection under an optimization framework, based on the Kullback–Leibler information and likelihood theory. With the exception of the methods are easy to compute by hand if necessary (as with the MLEs, maximized log-likelihood values, and $\widehat{\text{var}}(\hat{\theta}_i | \text{data})$ for the models). Researchers can easily understand the information-theoretic methods presented here; we believe that it is *very* important that researchers use the methods they employ.

References

- Agresti, A. (1990). *Categorical data analysis*. John Wiley and Sons.
- Aitkin, M. (1991). Posterior Bayes factors (with discussion). *Journal of the Royal Statistical Society, Series B*, **53**, 111–143.
- Akaike, H. (1973). Information theory as an extension of the maximum likelihood principle. Pages 267–281 in B. N. Petrov, and F. Csaki, (eds.) *Second International Symposium on Information Theory*. Akademiai Kiado, Budapest.
- Akaike, H. (1974). A new look at the statistical model identification. *Automatic Control AC* **19**, 716–723.
- Akaike, H. (1976). Canonical correlation analysis of time series and the Akaike information criterion. Pages 27–96 in R.K. Mehra, and D.G. Lainiotis (eds.) *Advances and Case Studies*. Academic Press, New York, NY.
- Akaike, H. (1977). On entropy maximization principle. Pages 27–36 in B. N. Petrov, and F. Csaki, (eds.) *Applications of statistics*. North-Holland, Amsterdam, The Netherlands.
- Akaike, H. (1978a). A new look at the Bayes procedure. *Biometrika* **65**, 237–242.
- Akaike, H. (1978b). A Bayesian analysis of the minimum AIC procedure. *Institute of Statistical Mathematics* **30**, 9–14.
- Akaike, H. (1978c). On the likelihood of a time series model. *The Annals of Statistics* **6**, 463–474.
- Akaike, H. (1979). A Bayesian extension of the minimum AIC procedure for model fitting. *Biometrika* **66**, 237–242.
- Akaike, H. (1980). Likelihood and the Bayes procedure (with discussion). J.M. Bernardo, M.H. De Groot, D.V. Lindley, and A.F.M. Smith (eds.) *Bayesian Statistics*. University Press, Valencia, Spain.

- Akaike, H. (1981b). Modern development of statistical methods. Eykhoff (ed.) *Trends and progress in system identification*. Pergamon Press, Oxford.
- Akaike, H. (1983a). Statistical inference and measurement of entropy. G.E.P. Box, T. Leonard, and C-F. Wu (eds.) *Scientific inference and robustness*. Academic Press, London.
- Akaike, H. (1983b). Information measures and model selection. *IEEE Transactions on Systems, Man, and Cybernetics* **13**, 267–276.
- Akaike, H. (1983c). On minimum information prior distributions. *Journal of Statistical Mathematics* **35A**, 139–149.
- Akaike, H. (1985). Prediction and entropy. Pages 1–24 in A.C. Atkinson and A.E. Hoaglin (eds.) *A celebration of Statistics*. Springer, New York, NY.
- Akaike, H. (1987). Factor analysis and AIC. *Psychometrika* **52**, 311–332.
- Akaike, H. (1992). Information theory and an extension of the maximum likelihood principle. Pages 610–624 in S. Kotz, and N.L. Johnson (eds.) *Breakthroughs in Statistics*. Springer-Verlag, London.
- Akaike, H. (1994). Implications of the informational point of view for statistical science. Pages 27–38 in H. Bozdogan (ed.) *Engineering in Statistics*. Vol. 3, Proceedings of the First US/Japan Conference on the Application of Statistics: Modeling: An Informational Approach. Kluwer Academic Publishers, Dordrecht, Netherlands.
- Akaike, H., and Nakagawa, T. (1988). *Statistical analysis and control*. KTK Scientific Publishers, Tokyo. (English translation by H. Akaike).
- Allen, D.M. (1970). Mean square error of prediction as a criterion for selection of models. *Technometrics* **13**, 469–475.
- Amari, S. (1993). Mathematical methods of neurocomputing. Pages 1–10 in J.L. Nelson, J.L. Jensen, and W.S. Kendall (eds.) *Networks and probabilistic aspects*. Chapman and Hall, New York, NY.
- Anonymous. (1997). *The Kullback Memorial Research Conference*. The George Washington University, Washington, D.C.
- Anderson, D. R. (2001). The need to get the basics right in wildlife management. *Society Bulletin* **29**, 1294–1297.
- Anderson, S., Auquier, A., Hauck, W.W., Oakes, D., Vandaele, P. (1980). *Statistical methods for comparative studies*. John Wiley & Sons, New York.
- Anderson, D.R., and Burnham, K.P. (1976). *Population ecology and the effect of exploitation on survival*. U.S. Fish and Wildlife Service, Washington, D.C. No. 128.
- Anderson, D.R., and Burnham, K.P. (1999a). General strategies for the analysis of ringing data. *Bird Study* **46** (suppl.), S261–270.

- Anderson, D.R., Burnham, K.P., and White, G.C. (1994). A model for overdispersed capture–recapture data. *Ecology* **75**, 1780–1793.
- Anderson, D.R., Burnham, K.P., and White, G.C. (1998). Comparison for model selection and statistical inference from capture–recapture data. *Applied Statistics* **25**, 263–282.
- Anderson, D.R., Burnham, K.P., Franklin, A.B., Gutierrez, R.J., Farnsworth, R.G., White, G.C., and Shenk, T.M. (1999). A protocol for conflict resolution of empirical data related to natural resource controversies. *Wildlife Society Bulletin* **27**, 1050–1058.
- Anderson, D.R., Burnham, K.P., and Thompson, W.L. (2000). Model selection problems, prevalence, and an alternative. *Journal of Wildlife Management* **64**, 1–13.
- Anderson, D.R., and K.P. Burnham. 2001a. Commentary on model selection. *The Ecological Society of America* **82**, 160–161.
- Anderson, D.R., Burnham, K.P., Gould, W.R., and Cherry, S. (2001b). Finding effects that are actually spurious. *Wildlife Society Bulletin* **29**, 1260–1270.
- Anderson, D.R., Link, W.A., Johnson, D.H. and Burnham, K.P. (2001c). Presenting results of data analyses. *Journal of Wildlife Management* **65**, 1–13.
- Apostol, T.M. (1957). *Mathematical analysis: a modern approach*. Addison-Wesley Publishing Co., Inc. Reading, MA.
- Armitage, P. (1957). Studies in the variability of pock counts. *Journal of the Royal Statistical Society* **20**, 564–581.
- Atilgan, T. (1996). Selection of dimension and basis for density estimation of dimension, basis and error distribution for regression. *Communications in Statistics Theory and Methods* **25**, 1–28.
- Atkinson, A.C. (1978). Posterior probabilities for choosing a regression model. *Biometrika* **65**, 39–48.
- Atkinson, A. C. (1980). A note on the generalized information criterion model. *Biometrika* **67**, 413–18.
- Atmar, W. (2001). A profoundly repeated pattern. *Bulletin of the Ecological Society of America* **82**, 208–211.
- Augustin, N.H., Muggleston, M.A., and Buckland S.T. (1996). Analysis of the spatial distribution of wildlife. *Journal of Applied Ecology* **33**, 1–13.
- Azzalini, A. (1996). *Statistical inference based on the likelihood*. London.

- Barron, A., Rissanen, J., and Yu, B. (1998). The minimum description length principle for coding and modeling. *IEEE Transactions on Information Theory* **44**, 2244–2260.
- Bartlett, M.S. (1936). Some notes on insecticide tests in the laboratory. *Journal of the Royal Statistical Society, Supplement* **1**, 185–194.
- Bedrick, E.J., and Tsai, C-L. (1994). Model selection for multivariate data with correlated samples. *Biometrics* **50**, 226–231.
- Berger, J.O., and Pericchi, L.R. (1996). The intrinsic Bayes factor for model selection and prediction. *American Statistical Association* **91**, 109–122.
- Berger, J.O., and Wolpert, R.L. (1984). The likelihood principle. *International Statistical Review* **52**, 359–382. *Statistics Monograph* 6.
- Berk, R.H. (1966). The limiting behavior of posterior distributions for normal distributions with incorrect priors. *Annals of Mathematical Statistics* **37**, 51–58.
- Bernardo, J.M., and Smith, A.F.M. (1994). *Bayesian theory*. John Wiley & Sons, Chichester, UK.
- Berry, D.A. (1988). Multiple comparisons, multiple tests and data dredging: a Bayesian perspective. *Bayesian Statistics 3*. Clarendon, Oxford, UK.
- Berryman, A.A., Gutierrez, A.P., and Arditi, R. (1995). Credible, precise estimates for predator–prey models—a reply to Abrams, Gleeson, and Sarnau. *Ecology* **76**, 1985.
- Beyers, D.W. (1998). Causal inference in environmental impact assessment. *North American Benthological Society* **17**, 367–373.
- Bliss, C.I. (1935). The calculation of the dosage-mortality curve. *Biometrics* **22**, 134–167.
- Bickel, P., and Zhang, P. (1992). Variable selection in nonparametric regression with categorical covariates. *Journal of the American Statistical Association* **87**, 1123–1134.
- Bishop, Y.M.M., Fienberg, S.E., and Holland, P.W. (1975). *Discrete Bayesian models: theory and practice*. The MIT Press, Cambridge, MA.
- Blansali, R.J. (1993). Order selection for linear time series models. In T. S. Rao (ed.) *Developments in time series analysis*, Chapman & Hall, London.
- Blau, G.E., and Neely, W.B. (1975). Mathematical model building to determine the distribution of DURSBAN® insecticide added to water. Pages 133–163 in A. Macfadyen (ed.) *Advances in Ecological Research*. Academic Press, London.
- Bollen, K.A., and Long, J.S. (1993). *Testing structural equations*. Sage Publications, Thousand Oaks, CA.
- Boltzmann, L. (1877). Über die Beziehung zwischen dem Hauptsatze der mechanischen Wärmetheorie und der Wahrscheinlichkeitsrechnung resp. der Theorie des Wärmegleichgewicht. *Wiener Berichte* **76**, 373–435.

- Box, G.E.P. (1967). Discrimination among mechanistic models. *Technometrics* **9**, 51–64.
- Box, G.E.P. (1976). Science and statistics. *Journal of the American Statistical Association* **71**, 791–799.
- Box, G.E.P., and Jenkins, G.M. (1970). *Time series analysis: forecasting and control*. Holden-Day, London.
- Box, G.E.P., Hunter, W.G., and Hunter, J.S. (1978). *Statistics for experimenters*. Wiley and Sons, New York, NY, USA.
- Box, G.E.P., Leonard, T., and Wu, C-F. (eds.) (1981). *Scientific inference and robustness*. Academic Press, London.
- Box, J.F. (1978). *R. A. Fisher: the life of a scientist*. John Wiley and Sons, New York, NY, USA.
- Boyce, M.S. (1992). Population viability analysis. *Annual Review of Ecology and Systematics* **23**, 481–506.
- Bozdogan, H. (1987). Model selection and Akaike's information criterion: a review of theory and its analytical extensions. *Psychometrika* **52**, 345–370.
- Bozdogan, H. (1988). A new model-selection criterion. Pages 599–610. In *Classification and related methods of data analysis*. North-Holland Publishing Co., Amsterdam, The Netherlands.
- Bozdogan, H. (1994). Editor's general preface. Pages ix–xii in H. Bozdogan, *Classification and Scientific Applications*. Vol. 3, Proceedings of the First U.S. Conference on the Frontiers of Statistical Modeling: An Informational Approach. Kluwer Academic Publishers, Dordrecht, the Netherlands.
- Breiman L. (1992). The little bootstrap and other methods for dimension reduction: regression: X-fixed prediction error. *Journal of the American Statistical Association* **87**, 738–754.
- Breiman, L. (1995). Better subset regression using the nonnegative garrote. *Journal of the American Statistical Association* **90**, 373–384.
- Breiman L. (1996). Heuristics of instability and stabilization in model selection. *Journal of Statistics* **24**, 2350–2383.
- Breiman, L. (2001). Statistical modeling: the two cultures (with commentary). *Science* **16**, 199–231.
- Breiman, L., and Freedman, D.F. (1983). How many variables should be put on the regression equation? *Journal of the American Statistical Association* **78**, 1115–1126.
- Brisbin I.L., Jr., Collins, C.T., White, G.C., and McCallum, D.A. (1984). Guidelines for the analysis and interpretation of growth data: the shape of the growth curve. *Journal of Wildlife Management* **48**, 552–554.
- Brockwell, P.J., and Davis, R.A. (1987). *Time series: theory and methods*. Chapman and Hall, New York, NY.

- Broda, E. (1983). *Ludwig Boltzmann: man, physicist, philosopher* (C. D. Broda, trans.). Oxford: Oxford University Press, Oxford, Connecticut.
- Brown, D. (1992). A graphical analysis of deviance. *Applied Statistics* **41**, 1–12.
- Brown, D., and Rothery, P. (1993). *Models in biology: mathematical computing*. John Wiley and Sons. New York, NY.
- Brown, L.D., Cai, T.T., and DasGupta, A. (2001). Interval estimation of a proportion. *Statistical Science* **16**, 101–135.
- Brown, P.J. (1993). *Measurement, regression, and calibration*. Chapman and Hall, UK.
- Brownie, C., Anderson, D.R., Burnham, K.P., and Robson, D.S. (1985). *Design for band recovery data—a handbook*. 2nd ed. U.S. Fish and Wildlife Service, Publication 156.
- Brownie, C., Hines, J.E., Nichols, J.D., Pollock, K.H., and Hestbeck, J.B. (1993). Capture-recapture studies for multiple strata including non-Markovian recapture. *Biometrics* **49**, 1173–1187.
- Brush, S.G. (1965). *Kinetic theory*. Vol. 1 Pergamon Press, Oxford.
- Brush, S.G. (1966). *Kinetic theory*. Vol. 2 Pergamon Press, Oxford.
- Bryant, P.G., and Cordero-Braña, O.I. (2000). Model selection and the description length principle. *The American Statistician* **54**, 257–267.
- Buckland, S.T. (1982). A note on the Fourier series model for analysis of variance. *Biometrics* **38**, 469–477.
- Buckland, S.T. (1984). Monte Carlo confidence intervals. *Biometrika* **71**, 1–11.
- Buckland, S.T., Anderson, D.R., Burnham, K.P., and Laake, J.L. (1993). *Estimating abundance of biological populations*. Chapman and Hall, London.
- Buckland, S.T., Anderson, D.R., Burnham, K.P., Laake, J.L., Borchers, D.L. (2001). *An introduction to distance sampling*. Oxford University Press, Oxford.
- Buckland, S.T., Burnham, K.P., and Augustin, N.H. (1997). Model selection: an introduction to the part of inference. *Biometrics* **53**, 603–618.
- Buckland, S.T., and Elston, D.A. (1993). Empirical models for the distribution of wildlife. *Journal of Applied Ecology* **30**, 478–495.
- Burman, P. (1989). A comparative study of ordinary cross-validation, the bootstrap and repeated learning-testing methods. *Biometrika* **76**, 239–244.
- Burman, P., and Nolan, D. (1995). A general Akaike-type criterion for model selection in robust regression. *Biometrika* **82**, 877–886.
- Burnham, K.P. (1988). A comment on maximum likelihood estimation of distributions. *Biometrical Journal* **30**, 379–384.

- Burnham, K.P. (in review). Basic random effects models in ringing data. *Ecological and Environmental Statistics*.
- Burnham, K.P., and Anderson, D.R. (1992). Data-based selection of statistical model: the key to modern data analysis. Pages 16–30 in D.R. Barrett (eds.) *Wildlife 2001: Populations*. Elsevier Scientific Publishing Co., Amsterdam.
- Burnham, K.P., and Anderson, D.R. (2001). Kullback-Leibler information criterion: strong inference in ecological studies. *Wildlife Research* **28**, 111–114.
- Burnham, K.P., Anderson, D.R., White, G.C., Brownie, C., and Pollock, K.H. (1999). *Model selection and analysis methods for fish survival experiments based on release-recapture data*. Fisheries Society, Monograph **5**.
- Burnham, K.P., Anderson, D.R., and White, G.C. (1994). Evaluation of the Akaike and Schwarz criteria with empirical results. *Biometrical Journal* **36**, 299–315.
- Burnham, K.P., White, G.C., and Anderson, D.R. (1995a). Model selection of capture–recapture data. *Biometrics* **51**, 888–898.
- Burnham, K.P., Anderson, D.R., and White, G.C. (1995b). Selection of capture–recapture models when capture probabilities are heterogeneous. *Applied Statistics* **22**, 611–624.
- Burnham, K.P., Anderson, D.R., and White, G.C. (1996). Meta-analysis of the Northern Spotted Owl. *Studies in Avian Biology* **17**, 92–101.
- Burnham, K.P., and White, G.C. (2002). Evaluation of some random effects models applicable to bird ringing data. *Journal of Applied Statistics* **29**, 111–124.
- Bystrak, D. (1981). The North American breeding bird survey, 1966–1979. Ralph, and J.M. Scott (eds.) *Estimating numbers of terrestrial birds*. *Biological Conservation* **6**, 1–10.
- Carlin, B., and Louis, T. (2000). *Bayes and empirical Bayes methods for data analysis*. Chapman and Hall, London.
- Carlin, B.P., and Chib, S. (1995). Bayesian model choice via Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society, Series B* **57**, 459–474.
- Carlin, B.P., and Louis, T.A. (1996). *Bayes and empirical Bayes methods for data analysis*. Chapman and Hall, London.
- Carpenter, S.R. (1990). Large-scale perturbations: opportunities for research. *Conservation Biology* **71**, 2038–2043.
- Carroll, R.J., and Ruppert, D. (1988). *Transformation and weighting in regression*. Chapman and Hall, New York, NY.
- Carroll, R., Ruppert, D., and Stefanski, L. (1995). *Measurement error in regression models*. Chapman and Hall, London.

- Caswell, H. (2001). *Matrix population models: construction, analysis, and programming*. Sunderland, MA: Sinauer Associates, Inc., Publishers, Sunderland, MA.
- Cavanaugh, J.E., and Neath, A.A. (1999). Generalizing the decision-theoretic information criterion. *Communication in Statistics—Theory and Methods* **27**, 1131–1144.
- Cavanaugh, J.E., and Shumway, R.H. (1997). A bootstrap variant of the Schwarz model selection. *Statistica Sinica* **7**, 473–496.
- Chamberlain, T.C. (1890). The method of multiple working hypotheses. (Reprinted 1965, *Science* **148**, 754–759).
- Chamberlin, T. C. (1965). (1890) The method of multiple working hypotheses. 754–759. (reprint of 1890 paper in *Science*).
- Chatfield, C. (1991). Avoiding statistical pitfalls (with discussion). *Journal of the Royal Statistical Society, Series A* **154**, 240–268.
- Chatfield, C. (1995a). *Problem solving: a statistician's guide*. Second edition. John Wiley and Hall, London.
- Chatfield, C. (1995b). Model uncertainty, data mining and statistical model selection. *the Royal Statistical Society, Series A* **158**, 419–466.
- Chatfield, C. (1996). Model uncertainty and forecast accuracy. *Journal of Forecasting* **15**, 495–508.
- Chen, M.-H., Shao, Q.-M., and Ibrahim, J.G. (2000). *Monte Carlo methods in finance and computation*. Springer, New York, NY.
- Cherry, S. (1998). Statistical tests in publications of The Wildlife Society. *Bulletin* **26**, 947–953.
- Chib, S., and Jeliazkov, I. (2001). Marginal likelihood from the Metropolis-Hastings algorithm. *Journal of the American Statistical Association* **96**, 270–281.
- Chow, G.C. (1981). A comparison of the information and posterior probability criteria for model selection. *Journal of Econometrics* **16**, 21–33.
- Chung, H-Y., Lee, K-W., and Koo, J-A. (1996). A note on bootstrap model selection criterion. *Statistics and Probability Letters* **26**, 35–41.
- Clayton, D., and Hills, M. (1993). *Statistical models in epidemiology*. Oxford University Press, Oxford, UK.
- Clayton, M.K., Geisser, S., and Jennings, D. (1986). A comparison of model selection procedures. Pages 425–439 in P. Goel, and A. Zellner (eds.) *Bayesian decision analysis*. Elsevier, New York, NY.
- Cochran, W.G. (1963). *Sampling techniques*. 2nd ed., John Wiley and Sons, New York, NY.
- Cohen, E.G.D., and Thirring, W. (eds.) (1973). *The Boltzmann equation and its applications*. Springer-Verlag, New York, NY.

- Collopy, F., Adya, M., and Armstrong, J.S. (1994). Principles for validity: the case of information systems spending forecasts. *Information Systems Research* **5**, 170–179.
- Conner, M.M., McCarty, C.W., and Miller, M.W. (2000). Detection of chronic wasting disease prevalence in mule deer. *Journal of Wildlife Management* **36**, 691–699.
- Conner, M.M., White, G.C., and Freddy, D.J. (2001). Elk movement and season hunting in northwest Colorado. *Journal of Wildlife Management* **65**, 1001–1010.
- Cook, T.D., and Campbell, D.T. (1979). *Quasi-experimentation: design and inference for field settings*. Houghton Mifflin Company, Boston, MA.
- Cook, R., Cook, J.G., Murray, D.L., Zager, P., Johnson, B.K., and Johnson, J.L. (1998). Development of predictive models of nutritional condition for mule deer. *Journal of Wildlife Management* **65**, 973–987.
- Copas, J.B. (1983). Regression, prediction and shrinkage (with discussion). *Royal Statistical Society, Series B*, **45**, 311–354.
- Cover, T.M., and Thomas, J.A. (1991). *Elements of information theory*. John Wiley & Sons, New York, NY.
- Cox, D.R. (1990). Role of models in statistical analysis. *Statistical Science* **5**, 1–47.
- Cox, D.R. (1995). The relation between theory and application in statistics. *Test* **4**, 207–261.
- Cox, D.R., and Reid, N. (2000). *The theory of the design of experiments*. John Wiley & Sons/Hall/CRC, Boca Raton, FL.
- Cox, D.R., and Snell, E.J. (1989). *Analysis of binary data*. 2nd ed., Chapman & Hall, New York, NY.
- Craven, P., and Wahba, G. (1979). Smoothing noisy data with spline functions: the correct degree of smoothing by the method of generalized cross-validation. *Mathematics* **31**, 377–403.
- Cressie, N.A.C. (1991). *Statistics for spatial data*. John Wiley and Sons, New York, NY.
- Cutler, A., and Windham, M.P. (1994). Information-based validity analysis. Pages 149–170 in H. Bozdogan (ed.) *Engineering and Artificial Intelligence*, vol. 2, Proceedings of the First US/Japan Conference on the Engineering and Artificial Intelligence Modeling: An Informational Approach. Kluwer Academic Publishers, Dordrecht, Netherlands.
- Daniel, C., and Wood, F.S. (1971). *Fitting equations to data*. Wiley, New York, NY.
- de Gooijer, J.G., Abraham, B., Gould, A., and Robinson, L. (1985). On the order of an autoregressive-moving average process: a survey. *Journal of the Royal Statistical Society, Series B* **53**, 301–329.

- deLeeuw, J. (1992). Introduction to Akaike (1973) information theory and the maximum likelihood principle. Pages 599–609 in S. Kotz, and J. R. van der Meer (eds.), *Breakthroughs in statistics*. Vol. 1. Springer-Verlag, London.
- Dempster, A.P. (1971). Model searching and estimation in the logic of probability. *on the Foundations of Statistical Inference*. University of Waterloo, Ontario, Canada.
- Dempster, A.P. (1997). The direct use of likelihood for significance testing. *Computing* **7**, 247–252.
- Desu, M.M., and Roghavarao, D. (1991). *Sample size methodology*. Marcel Dekker, New York, NY.
- Dijkstra, T.K. (ed). (1988). *On model uncertainty and its statistical implications*. Notes in Economics and Mathematical Systems, Springer-Verlag, New York, NY.
- Dijkstra, T.K., and Veldkamp, J.H. (1988). Data-driven selection of models using bootstrap. Pages 17–38 in T.K. Dijkstra (ed.) *On model uncertainty and its implications*. Lecture Notes in Economics and Mathematical Systems, Springer-Verlag, New York, NY.
- Draper, D. (1995). Assessment and propagation of model uncertainty. *Journal of the Royal Statistical Society, Series B* **57**, 45–97.
- Draper, N.R., and Smith, H. (1981). *Applied regression analysis*. Second edition. Wiley and Sons, New York, NY.
- Eberhardt, L.L. (1978). Appraising variability in population studies. *Management Science* **42**, 207–238.
- Eberhardt, L.L., and Thomas, J.M. (1991). Designing environmental monitoring systems. *Ecological Monographs* **61**, 53–73.
- Edwards, A.W.F. (1976). *Likelihood: an account of the statistical concept and its application to scientific inference*. Cambridge University Press, Cambridge, UK.
- Edwards, A.W.F. (1992). *Likelihood: expanded edition*. The Johns Hopkins University Press, Baltimore, MD.
- Edwards, A.W.F. (2001). Occam's bonus. In A. Zellner, H.A. Keuzer, and J. Zellner (eds.), *Simplicity, inference and modelling*, pgs. 128–132. Cambridge University Press, Cambridge, UK.
- Edwards, D. (1998). Issues and themes for natural resources trend analysis. *Ecological Applications* **8**, 323–325.
- Efron, B. (1979) Bootstrap methods: another look at the jackknife. *Biometrika* **7**, 1–26.
- Efron, B. (1983). Estimating the error rate of a prediction rule: empirical Bayes validation. *Journal of the American Statistical Association* **78**, 311–319.

- Efron, B. (1986). How biased is the apparent error rate of a prediction? *American Statistical Association* **81**, 461–470.
- Efron, B., and Morris, C. (1975). Data analysis using Stein's paradox: some generalizations. *Journal of the American Statistical Association* **70**, 119–127.
- Efron, R., and Gong, G. (1983). A leisurely look at the bootstrap, the jackknife, and cross-validation. *The American Statistician* **37**, 36–48.
- Efron, B., and Tibshirani, R.J. (1993). *An introduction to the bootstrap*. London: Chapman and Hall.
- Ellison, A.M. (1996). An introduction of Bayesian inference for environmental decision-making. *Ecological Applications* **6**, 103–114.
- Feder, M., Merhav, N., and Gutman, M. (1992). Universal prediction for ergodic sequences. *IEEE Transactions on Information Theory* **38**, 1258–1266.
- Fienberg, S.E. (1970). The analysis of multidimensional contingency tables. *Biometrika* **57**, 445–454.
- Fildes, R., and Makridakis, S. (1995). The impact of empirical autocorrelation on time series analysis and forecasting. *International Statistics Review* **63**, 1–12.
- Findley, D.F. (1985). On the unbiasedness property of AIC for exact autoregressive stochastic time series models. *Journal of Time Series Analysis* **6**, 1–12.
- Findley, D.F. (1991). Counterexamples to parsimony and BIC. *American Statistical Mathematics* **43**, 505–514.
- Findley, D.F., and Parzen, E. (1995). A conversation with Hirotugu Akaike. *Science* **10**, 104–117.
- Finney, D.J. (1971). *Probit analysis*. 3rd. ed. Cambridge University Press.
- Fisher, R.A. (1922). On the mathematical foundations of theoretical statistics. London: Cambridge University Press.
- Fisher, R.A. (1936). Uncertain inference. *Proceedings of the American Philosophical Association and Sciences* **71**, 245–58.
- Fisher, R.A. (1949). A biological assay of tuberculins. *Biometrics* **5**, 1–14.
- Flack, V.F., and Chang, P.C. (1987). Frequency of selecting non-significant regression analysis: a simulation study. *The American Statistician* **41**, 1–5.
- Flather, C.H. (1992). Patterns of avian species-accumulation rates across different landscapes. Ph.D. dissertation. Colorado State University. Fort Collins, CO.
- Flather, C.H. (1996). Fitting species-accumulation functions and assessing the use impacts on avian diversity. *Journal of Biogeography* **23**, 155–168.
- Ford, E.D. (2000). *Scientific method for ecological research*. Cambridge University Press, Cambridge, UK.

- Forster, M.R. (2000). Key concepts in model selection: performance. *Journal of Mathematical Psychology* **44**, 205–231.
- Forster, M.R. (2001). The new science of simplicity. In A. Zellner and M. McAleer (eds.), *Simplicity, inference and modelling*, pp. University Press, Cambridge, UK.
- Forster, M.R., and Sober, E. (1994). How to tell when simpler, *ad hoc* theories will provide more accurate predictions. *British Journal for the Philosophy of Science* **45**, 399–424.
- Franklin, A.B., Anderson, D.R., and Burnham, K.P. (2002). Estimation and variation in avian survival probabilities using random effects. *Applied Statistics* **29**, 267–287.
- Franklin, A.B., Shenk, T.M., Anderson, D.R., and Burnham, K.P. (2002). Model selection: the alternative to null hypothesis testing. Pages 75–90. In A.B. Franklin (eds.) *Modeling in Natural Resource Management*. Island Press, Washington, D. C.
- Freedman, D.A. (1983). A note on screening regression equations. *Technometrics* **37**, 152–155.
- Freedman, D. (1999). From association to causation: some remarks on causal statistics. *Statistical Science* **14**, 243–258.
- Freedman, D.A., Navidi, W., and Peters, S.C. (1988). On the implications of overfitting in fitting regression equations. Pages 1–16 in T.K. Dijkstra (ed.) *Overfitting and its statistical implications*. Lecture Notes in Economics and Mathematical Systems, Springer-Verlag, New York, NY.
- Fujikoshi, Y. and Satoh, K. (1997). Modified AIC and C_p in multivariate analysis. *Biometrika* **84**, 707–716.
- Gail, M.H. (1996). Statistics in action. *Journal of the American Statistical Association* **91**, 1–13.
- Gallant, A.R. (1987). *Nonlinear statistical models*. John Wiley and Sons, New York.
- Gamerman, D. (1997). *Markov Chain Monte Carlo*. Chapman and Hall, London.
- Garthwaite, P.H., Jolliffe, I.T., and Jones, B. (1995). *Statistical inference*. John Wiley, London.
- Gause, G.F. (1934). *The struggle for existence*. Williams and Wilkins, Baltimore.
- Geisser, S. (1975). The predictive sample reuse method with applications. *American Statistical Association* **70**, 320–328.
- Gelfand, A.E., and Smith, A.F.M. (1990). Sampling-based approaches to calculating marginal densities. *Journal of the American Statistical Association* **85**, 398–409.
- Gelfand, A., and Dey, D.K. (1994). Bayesian model choice: an empirical Bayes calculation. *Journal of the Royal Statistical Society, Series B*, **56**, 501–514.

- Gelman, A., Carlin, J.B., Stern, H.S. and Rubin, D.B. (1995). *B*. Chapman and Hall, London.
- George, E.I., and Foster, D.P. (2000). Calibration and empirical B. *Biometrika* **87**, 731–748.
- George, E.I., and McCulloch, R.E. (1993). Variable selection via C. *of the American Statistical Association* **88**, 881–889.
- Gerard, P.D., Smith, D.R., and Weerakkody, G. (1998). Limits analysis. *Journal of Wildlife Management* **62**, 801–807.
- Gilchrist, W. (1984). *Statistical modelling*. Chichester, Wiley and S.
- Gilks, W. R., Richardson, S., and Spiegelhalter, D.J. (1996). *Marko practice*. Chapman and Hall, London.
- Glymour, C. (1998). Causation. Pages 97–109 in S. Kotz (ed.), *En sceinces*. John Wiley and Sons, New York, NY.
- Gochfeld, M. (1987). On paradigm vs. methods in the study of gro
- Gokhale, D.V., and Kullback, S. (1978). *The information in con*. Dekker, New York, NY.
- Golub, G.H., Health, M., and Wahba, G. (1979). Generalized cross for choosing a good ridge parameter. *Technometrics* **21**, 215–22
- Goodman, S.N. (1993). *p* values, hypothesis tests, and likelihood demiology of a neglected historical debate (with discussion) *Epidemiology* **137**, 485–501.
- Goodman, S.N., and Berlin, J.A. (1994). The use of predicted cor planning experiments and the misuse of power when interpre *Internal Medicine* **121**, 200–206.
- Goutis, C., and Casella, G. (1995). Frequentist post-data inference. *Review* **63**, 325–344.
- Granger, C.W.J., King, M.L., and White, H. (1995). Comments on te and the use of model selection criteria. *Journal of Econometrics*
- Graybill, F.A., and Iyer, H.K. (1994). *Regression analysis: con*. Duxbury Press, Belmont, CA.
- Greenhouse, S.W. (1994). Solomon Kullback: 1907–1994. *Ins Statistics Bulletin* **23**, 640–642.
- Guiasu, S. (1977). *Information theory with applications*. McGraw-
- Guisan, A., and Zimmermann, N.E. (2000). Predictive habitat ecology. *Ecological Modelling* **135**, 147–186.
- Ha, I.D., Lee, Y., and Song, J-K. (2001). Hierarchical likelihood app *Biometrika* **88**, 233–243.

- Hald, A. (1952). *Statistical theory with engineering applications*. John Wiley, New York, NY.
- Hald, A. (1998). *A history of mathematical statistics*. John Wiley, New York, NY.
- Hand, D.J. (1994). Statistical strategy: step 1. Pages 1–9 in P. Cheesbrough and D.J. Hand (eds.) *Selecting models from data*. Springer-Verlag, New York, NY.
- Hand, D.J. (1998). Data mining: statistics and more? *The American Statistician* **52**, 111–118.
- Hand, D.J., Blunt, G., Kelly, M.G., and Adams, N.M. (2000). Data mining: a new statistical science? *Statistical Science* **15**, 111–131.
- Hannan, E.J., and Quinn, B.G. (1979). The determination of the order of the model. *Journal of the Royal Statistical Society, Series B* **41**, 190–195.
- Hansen, M.H., and Yu, B. (2001). Model selection and the principle of minimum length. *Journal of the American Statistical Association* **96**, 746–757.
- Harlow, L.L., Mulaik, S.A., and Steiger, J.H. (eds.) (1997). *What if the results don't work? What if the tests don't work?* Lawrence Erlbaum Associates, Publishers, Mahwah, NJ.
- Hasenöhr, F. (ed.) (1909). *Wissenschaftliche Abhandlungen*. 3 Volumes. Leipzig, Germany.
- Hastie, T.J., and Tibshirani, R.J. (1990). *Generalized additive models*. Chapman and Hall, London.
- Haughton, D. (1989). Size of the error in the choice of a model to fit a family. *Sankhya, Series A* **51**, 45–58.
- Hayne, D. (1978). Experimental designs and statistical analyses. Pages 1–10 in D. Hayne (ed.) *Populations of small mammals under natural conditions*. Pennsylvania State University Press, University of Pittsburgh, Vol 5.
- Henderson, H., and Velleman, P. (1981). Building multiple regression models. *Biometrics* **37**, 391–411.
- Heyde, C.C. (1997). *Quasi-likelihood and its application: a general approach to parameter estimation*. Springer-Verlag, New York, NY.
- Hilborn, R., and Mangel, M. (1997). *The ecological detective: computer intensive statistical methods for ecology*. Princeton University Press, Princeton, NJ.
- Hjorth, J.S.U. (1994). *Computer intensive statistical methods: validation, bootstrapping and bootstrap*. Chapman and Hall, London.
- Hobson, A., and Cheng, B-K. (1973). A comparison of the entropy and information measures. *Journal of Statistical Physics* **7**, 301–310.
- Hocking, R.R. (1976). The analysis and selection of variable models. *Biometrics* **32**, 1–49.
- Hoenig, J.M., and Heisey, D.M. (2001). The abuse of power: the probability of power calculations for data analysis. *The American Statistician* **55**, 19–24.

- Hoeting, J.A., Madigan, D., Raftery, A.E., and Volinsky, C.T. (1999). Model averaging: a tutorial (with discussion). *Statistical Science* **14**, 382–411.
- Hosmer, D.W., and Lemeshow, S. (1989). *Applied logistic regression*. John Wiley and Sons, New York, NY.
- Howard, R.A. (1971). *Dynamic probabilistic systems*. John Wiley and Sons, New York, NY.
- Hurvich, C.M., Simonoff, J.S., and Tsai, C-L. (1998). Smoothing nonparametric regression using an improved Akaike information criterion. *Royal Statistical Society, Series B*, **60**, 271–293.
- Hurvich, C.M., and Tsai, C-L. (1989). Regression and time series models in small samples. *Biometrika* **76**, 297–307.
- Hurvich, C.M., and Tsai, C-L. (1990a). Model selection for linear regression in small samples. *Statistics and Probability Letters* **9**, 1–7.
- Hurvich, C.M., and Tsai, C-L. (1990b). The impact of model selection on regression. *The American Statistician* **44**, 214–217.
- Hurvich, C.M., and Tsai, C-L. (1991). Bias of the corrected AIC criterion in regression and time series models. *Biometrika* **78**, 499–509.
- Hurvich, C.M., and Tsai, C-L. (1994). Autoregressive model selection using a bias-corrected version of AIC. Pages 137–157 in H. Bozdogan, *Statistical and Scientific Applications*. Vol. 1, Proceedings of the First U.S.-Turkey Symposium on the Frontiers of Statistical Modeling: An Informational Approach. Kluwer Publishers, Dordrecht, the Netherlands.
- Hurvich, C.M., and Tsai, C-L. (1995a). Relative rates of convergence of model selection criteria in linear regression. *Biometrika* **82**, 418–425.
- Hurvich, C.M., and Tsai, C-L. (1995b). Model selection for extreme values in small samples. *Biometrics* **51**, 1077–1084.
- Hurvich, C.M., and Tsai, C-L. (1996). The impact of unsuspected model selection in linear regression. *Statistics and Probability Letters* **28**, 1–10.
- Hurvich, C.M., Shumway, R., and Tsai, C-L. (1990). Improved cross-entropy and Leibler information for autoregressive model selection in small samples. *Biometrika* **77**, 709–719.
- Ibrahim, J.G., and Chen, M-H. (1997). Predictive variable selection in a linear model. *Biometrics* **53**, 465–478.
- Inman, H.F. (1994). Karl Pearson and R. A. Fisher on statistical tests. *Nature. The American Statistician* **48**, 2–11.
- Irizarry, R.A. (2001). Information and posterior probability criteria for local likelihood estimation. *Journal of the American Statistical Association* **96**, 111–121.

- Jaffe, A.J., and Spierer, H.F. (1987). *Misused statistics: straight talk*. Marcel Dekker, Inc., New York, NY.
- James, F.C., and McCulloch, C.E. (1990). Multivariate analysis in ecology: a panacea or Pandora's box? *Annual Reviews of Ecology and Systematics* **21**, 129–166.
- Jaynes, E.T. (1957). Information theory and statistical mechanics. *Physical Review* **106**, 620–630.
- Jaynes, E.T. (1982). On the rationale of maximum-entropy methods. *IEEE* **70**, 939–952.
- Jaynes, E.T. (in prep.). *Probability theory: the logic of science*. Cambridge University Press, Cambridge, UK.
- Jeffreys, H. (1948). *Theory of probability*. Oxford University Press, Oxford, UK.
- Jeffreys, H. (1973). *Scientific inference*. 3rd ed. Cambridge University Press, Cambridge, UK.
- Jessop, A. (1995). *Informed assessments: an introduction to informed statistics*. Ellis Horwood, London.
- Jevons, W. S. (1874, 1877). *The principles of science*. MacMillan, London, Kingdom.
- Jiménez, J.A., Hughes, K.A., Alaks, G., Graham, L., and Lacy, R.C. (1999). A study of inbreeding depression in a natural habitat. *Science* **266**, 1558–1562.
- Johnson, D.H. (1995). Statistical sirens: the allure of nonparametric tests. *Wildlife Management* **63**, 763–772.
- Johnson, D.H. (1999). The insignificance of statistical significance testing. *Wildlife Management* **63**, 763–772.
- Johnson, J.W. (1996). Fitting percentage of body fat to simple body weight. *Journal of Statistics Education* **4** (e-journal).
- Johnson, N.L., and Kotz, S. (1970). *Continuous univariate distributions*. John Wiley & Sons, New York, NY.
- Johnson, N.L., and Kotz, S. (1992). *Univariate discrete distributions*. Interscience Publication, New York, NY.
- Jones, D., and Matloff, N. (1986). Statistical hypothesis testing in ecology: a review in terms. *Journal of Economic Entomology* **79**, 1156–1160.
- Judge, G.C., and Yancey, T. (1986). *Improved methods of inference*. North-Holland, Amsterdam, the Netherlands.
- Kabaila, P. (1995). The effect of model selection on confidence regions. *Econometric Theory* **11**, 537–549.
- Kapur, J.N., and Kesavan, H.K. (1992). *Entropy optimization principles*. Academic Press, London.

- Kass, R.E., and Raftery, A.E. (1995). Bayes factors. *Journal of the American Statistical Association* **90**, 773–795.
- Kass, R.E., and Raftery, A.E. (1995). “Bayes factors.” *Journal of the American Statistical Association* **90**, 773–795.
- Kass, R.E., and Wasserman, L. (1995). A reference Bayesian test for its relationship to the Schwarz criterion. *Journal of the American Statistical Association* **90**, 928–934.
- Kishino, H., Kato, H., Kasamatsu, F., and Fujise, Y. (1991). Detection and estimation of population characteristics from field survey data: a feasibility study of the Southern Hemisphere minke whales. *Annals of the Institute of Statistical Mathematics* **43**, 435–453.
- Kiso, K., Akamine, T., Ohnishi, S., and Matsumiya, Y. (1992). Maturation of the growth of sea-run and fluviatile forms of the female masu salmon (*Salmo masou*) in rivers of the southern Sanriku district, Honshu, Japan. *Marine Biology* **58**, 1779–1784.
- Kittrell, J.R. (1970). Mathematical modelling of chemical reactors. *Chemical Engineering* **8**, 97–183.
- Knopf, F.L., Sedgwick, J.A., and Cannon, R.W. (1988). Guild structure relative to seasonal cattle grazing. *Journal of Wildlife Management* **52**, 103–112.
- Konishi, S., and Kitagawa, G. (1996). Generalized information criterion. *Biometrika* **83**, 875–890.
- Kooperberg, C., Bose, S., and Stone, C.J. (1997). Polychotomous regression. *American Statistical Association* **92**, 117–127.
- Kreft, I., and deLeeuw, J. (1998). *Introducing multilevel modeling*. Thousand Oaks, CA, USA.
- Kuhn, T.S. (1970). *The structure of scientific revolutions*. 2nd ed. Chicago Press, Chicago, IL.
- Kullback, S. (1959). *Information theory and statistics*. John Wiley & Sons, NY.
- Kullback, S. (1987). The Kullback-Leibler distance. *The American Statistician* **41**, 77–80.
- Kullback, S., and Leibler, R.A. (1951). On information and entropy. *Mathematical Statistics* **22**, 79–86.
- Laake, J.L., Buckland, S.T., Anderson, D.R., and Burnham, K.P. (1998). *Distance 2.1: a guide*. Version 2.1. Colorado Cooperative Fish and Wildlife Research Unit, Colorado State University, Fort Collins, CO.
- Lahiri, P. (ed.) (2001). Model selection. Institute of Mathematical Statistics, No. 38.

- Larimore, W.E., and Mehra, R.K. (1985). The problem of overfitting. *Statistical Society, Series B* **57**, 247–262.
- Laud, P.W., and Ibrahim, J.G. (1995). Predictive model selection. *Statistical Society, Series B* **57**, 247–262.
- Laud, P.W., and Ibrahim, J.G. (1996). Predictive specification of p in variable selection. *Biometrika* **83**, 267–274.
- Leamer, E.E. (1978). *Specification searches: ad hoc inference with*. John Wiley and Sons, New York, NY.
- Lebreton, J-D., Burnham, K.P., Clobert, J., and Anderson, D.R. (1992). and testing biological hypotheses using marked animals: a unified framework. *Ecological Monograph* **62**, 67–118.
- Lee, Y., and Nelder, J.A. (1996). Hierarchical generalized linear models. *Royal Statistical Society, Series B* **58**, 619–678.
- Lehmann, E.L. (1983). *Theory of point estimation*. John Wiley and Sons, New York, NY.
- Lehmann, E.L. (1990). Model specification: the views of Fisher and Neyman. *Statistical Science* **5**, 160–168.
- Leirs, H., Stenseth, N.C., Nichols, J.D., Hines, J.E., Verhagen, R., and Schmid, B. (1998). Stochastic seasonality and nonlinear density-dependent factors in an African rodent. *Nature* **389**, 176–180.
- Leonard, T., and Hsu, J.S.J. (1999). *Bayesian methods: an analytical approach for interdisciplinary researchers*. Cambridge University Press, Cambridge, UK.
- Leroux, B.G. (1992). Consistent estimation of a mixing distribution. *Journal of the Royal Statistical Society, Series B* **20**, 1350–1360.
- Lewis, P. (1998). Maximum likelihood as an alternative to parsimony in molecular systematics. Pages 132–163 in D. Soltis, P. Soltis, and M. Donnell-Smith (eds.) *Molecular systematics of plants II*. Kluwer Publishing, Boston, MA.
- Liang, K-Y, and McCullagh, P. (1993). Case studies in binary distribution. *Journal of the Royal Statistical Society, Series B* **55**, 623–630.
- Lindley, D.V. (1986). The relationship between the number of factors in a mixture. Pages 459–470 in P.K. Goel, and A. Zellner (eds.) *Bayesian techniques*. Elsevier Science Publishers, New York, NY.
- Lindsey, J.K. (1995). *Modeling frequency and count data*. Oxford University Press, Oxford, UK.
- Lindsey, J.K. (1996). *Parametric statistical inference*. Oxford Science Publications, Oxford, UK.
- Lindsey, J.K. (1999a). On the use of corrections for overdispersion in count data. *Journal of the Royal Statistical Society, Series B* **61**, 553–561.
- Lindsey, J.K. (1999b). Some statistical heresies. *The Statistician* **48**, 1–10.

- Lindsey, J.K., and Jones, B. (1998). Choosing among generalized linear models using medical data. *Statistics in Medicine* **17**, 59–68.
- Linhart, H. (1988). A test whether two AIC's differ significantly. *Statistica* **48**, 153–161.
- Linhart, H., and Zucchini, W. (1986). *Model selection*. John Wiley & Sons, New York, NY.
- Longford, N.T. (1993). *Random coefficient models*. Oxford University Press, New York, NY.
- Longford, N.T. and Nelder, J.A. (1999). Statistics versus statistical process control. *Statistics in Medicine* **18**, 2311–2320.
- Lucky, R.W. (1991). *Silicon dreams: information, man, and machine*. John Wiley & Sons, New York, NY.
- Ludwig, D. (1989). Small models are beautiful: efficient estimators and model selection. Pages 274–284 in C. Castillo-Chavez, S.A. Levin, and C.A. Shoemaker (eds.) *Statistical approaches to problems in resource management and epidemiology*. Chapman and Hall, London.
- Lunneborg, C.E. (1994). *Modeling experimental and observational data*. Sage Publications, Belmont, CA.
- Lytle, D.A. (2002). Flash floods and aquatic insect life-history traits: comparing multiple models. *Ecology* **83**, 370–385.
- Madigan, D., and Raftery, A.E. (1994). Model selection and account selection in graphical models using Occam's window. *Journal of the Royal Statistical Society B* **56**, 1133–1153.
- Madigan, D.M., Raftery, A.E., York, J.C., Bradshaw, J.M., and Almond, D. (1994). Graphical model selection. Pages 91–100 in P. Chessa, D. Madigan, and A. Raftery (eds.) *Selecting models from data: AI and statistics IV*. Springer-Verlag, New York, in Statistics **89**.
- Mallows, C.L. (1973). Some comments on C_p . *Technometrics* **12**, 135–139.
- Mallows, C.L. (1995). More comments on C_p . *Technometrics* **37**, 117–119.
- Manly, B.F.J. (1991). *Randomization and Monte Carlo methods in biology*. Chapman and Hall, New York, NY.
- Manly, B.F.J., McDonald, L.L., and Thomas, D.L. (1993). *Resource selection by animals: statistical design and analysis for field studies*. Chapman and Hall, New York, NY.
- Manly, B.F.J. (1992). *The design and analysis of research studies in wildlife management*. Cambridge University Press, Cambridge, UK.
- Manly, B.F.J. (2001). *Statistics for environmental science and management*. Chapman and Hall, London.
- Marshall, J.R. (1990). Data dredging and noteworthiness. *Epidemiology* **1**, 10–12.

- MATLAB® (1994) *High-performance numerical computations and graphics*. The MathWorks, Inc., Natick, MA.
- Maurer, B.A. (1998). Ecological science and statistical paradigms: a review. *Ecology* **79**, 502–503.
- Mayr, E. (1997). *This is biology: the science of the living world*. Harvard University Press. Cambridge, MA.
- McBride, G.B., Loftis, J.C., and Adkins, N.C. (1993). What do signs tell us about the environment? *Environmental Management* **17**, 423–430.
- McCullagh, P., and Nelder, J.A. (1989). *Generalized linear models*. Chapman and Hall, New York, NY.
- McCullagh, P., and Pregibon, D. (1985). Discussion comments on the work of Efron. *Annals of Statistics* **13**, 898–900.
- McLachlan, G.J., and Peel, D. (2000). *Finite Mixture Models*. John Wiley & Sons, New York, NY.
- McQuarrie, A.D. (1999). A small-sample correction of the Schwarz criterion. *Statistics and Probability Letters* **44**, 79–86.
- McQuarrie, A.D.R., and Tsai, C-L. (1998). *Regression and time series analysis*. World Scientific Publishing Company, Singapore.
- Mead, R. (1988). *The design of experiments: statistical principles for experimental design*. Cambridge University Press, New York, NY.
- Miller, A.J. (1990). *Subset selection in regression*. Chapman and Hall, New York, NY.
- Mooney, C.Z., and Duval, R.D. (1993). *Bootstrapping: a nonparametric approach to statistical inference*. Sage Publications, London.
- Moore, D.F. (1987). Modelling the extraneous variance in the presence of overdispersion. *Journal of the Royal Statistical Society* **36**, 8–14.
- Morgan, B.J.T. (1992). *Analysis of quantal response data*. Chapman and Hall, New York, NY.
- Morgan, B.J.T. (2000). *Applied stochastic modelling*. Arnold Press, London.
- Morris, C.N. (1983). Parametric empirical Bayes inference: theory and practice. *Journal of the American Statistical Association* **78**, 47–65.
- Mosteller, F., and Tukey, J.W. (1968). Data analysis, including statistical computing. In E. Aronson (eds.) *Handbook of Social Psychology*, Vol. 2. Academic Press, New York, MA.
- Myers, R.A., Barrowman N.J., Hutchings, J.A., and Rosenberg, A. (1998). The dynamics of exploited fish stocks at low populations levels. *Science* **280**, 1480–1483.
- Naik, P.A., and Tsai, C-L. (2001). Single-index model selections. *Biometrika* **88**, 585–595.
- Nester, M. (1996). An applied statistician's creed. *Applied Statistics* **45**, 1–10.

- Nichols, J.D., and Kendall, W.L. (1995). The use of multi-strata capture-recapture models to address questions in evolutionary ecology. *Journal of Applied Ecology* **32**, 401–412.
- Nishii, R. (1988). Maximum likelihood principle and model selection when the true model is unspecified. *Journal of Multivariate Analysis* **27**, 392–403.
- Noda, K., Miyaoka, E., and Itoh, M. (1996). On bias correction of Akaike's information criterion in linear models. *Communications in Statistics—Theory and Methods* **25**, 1857–1871.
- Norris, J.L., and Pollock, K.H. (1995). A capture–recapture model for estimating survival in the presence of behavioural response. *Environmental and Ecological Statistics* **2**, 139–154.
- Norris, J.L., and Pollock, K.H. (1997). Including model uncertainty in population estimates from multiple capture studies. *Environmental and Ecological Statistics* **4**, 127–141.
- O'Connor, M.P., and Spotila, J.R. (1992). Consider a spherical lizard model for estimating survival approximations. *American Zoologist* **32**, 179–193.
- O'Connor, R.J. (2000). Why ecology lags behind biology. *The Science of Nature* **147**, 101–109.
- O'Hagan, A. (1995). Fractional Bayes factors for model comparison. *Journal of the Royal Statistical Society, Series B* **57**, 99–138.
- Olden, J.D., and Jackson, D.A. (2000). Torturing data for the sake of model selection: are our regression models? *Ecoscience* **7**, 501–510.
- Otis, D.L., Burnham, K.P., White, G.C., and Anderson, D.R. (1978). Estimating survival from capture data on closed animal populations. *Wildlife Monographs* **68**, 3–64.
- Pan, W. (1999). Bootstrapping likelihood for model selection with applications. *Journal of Computational and Graphical Statistics* **8**, 687–698.
- Pan, W. (2001a). Akaike's information criterion in generalized linear models. *Biometrics* **57**, 120–125.
- Pan, W. (2001b). Model selection in estimating equations. *Biometrics* **57**, 126–134.
- Parzen, E. (1994). Hirotugu Akaike, statistical scientist. Pages 25–34. In *Engineering and Scientific Applications*. Vol. 1, Proceedings of the International Conference on the Frontiers of Statistical Modeling: An Information-Theoretic Perspective. Academic Publishers, Dordrecht, the Netherlands.
- Parzen, E., Tanabe, K., and Kitagawa, G. (eds.) (1998). *Selected papers of Hirotugu Akaike*. Springer-Verlag Inc., New York, NY.
- Pascual, M.A., and Kareiva, P. (1996). Predicting the outcome of experiments from experimental data: maximum likelihood and Bayesian approaches. *Ecology* **77**, 1000–1010.
- Peirce, C.S. (1955). Abduction and induction. Pages 150–159. In *Philosophical writings of Peirce*. Dover, New York, NY.

- Peters, R.H. (1991). *A critique for ecology*. Cambridge University USA.
- Peterson, T.S. (1960). *Elements of calculus* (2nd ed.). Harper Brothers.
- Platt, J.R. (1964). Strong inference. *Science* **146**, 347–353.
- Pollock, K.H., Nichols, J.D., Brownie, C., and Hines, J.E. (1990). *capture–recapture experiments*. *Wildlife Monographs*. **107**, 1–97.
- Pope, S.E., Fahrig, L., and Merriam, H.G. (2000). Landscape metapopulation effects on leopard frog populations. *Ecology* **81**.
- Posada, D., and Crandall, K. (1998). MODELTEST: testing the model. *Bioinformatics* **14**, 817–818.
- Posada, D., and Crandall, K. (2001). Selecting models of nucleotide substitution of human immunodeficiency virus 1 (HIV-1). *Molecular Biology and Evolution* **18**, 897–906.
- Poskitt, D.S., and Tremayne A.R. (1987). Determining a portfolio of models. *Biometrika* **74**, 125–137.
- Pötscher, B.M. (1991). Effects of model selection on inference. *Biometrika* **78**, 163–185.
- Qian, G., Gabor, G., and Gupta, R.P. (1996). Generalized linear predictive least quasi-deviance criterion. *Biometrika* **83**, 41–54.
- Qin, J., and Lawless, G. (1994). Empirical likelihood and generalized likelihood. *Annals of Statistics* **22**, 300–325.
- Quinn, J.F., and Dunham, A.E. (1983). On hypothesis testing in the field. *American Naturalist* **122**, 22–37.
- Raftery, A.E. (1995). Bayesian model selection in social research. *Sociological Methodology* **25**, 111–195.
- Raftery, A. (1996a). Approximate Bayes factors and accounting for model uncertainty in generalized linear regression models. *Biometrika* **83**, 251–266.
- Raftery, A. (1996b). Hypothesis testing and model selection. Pages 165–190. In S. Richardson, and D.J. Spiegelhalter (eds.), *Markov chain Monte Carlo: practical issues*. Chapman and Hall, London.
- Raftery, A., Madigan, D.M., and Hoeting, J. (1993). *Model selection and model uncertainty in linear regression models*. Technical Report 154, Department of Statistics, University of Washington, Seattle, WA.
- Raftery, A.E., Madigan, D., and Hoeting, J.A. (1997). Bayesian model selection in generalized linear regression models. *Journal of the American Statistical Association* **92**, 133–147.
- Rawlings, J.O. (1988). *Applied regression analysis: a research synthesis*. Belmont, CA.

- Renshaw, E. (1991). *Modelling biological populations in space*. University Press, Cambridge, UK.
- Reschenhofer, E. (1996). Prediction with vague prior knowledge. *Statistics—Theory and Methods* **25**, 601–608.
- Reschenhoffer, E. (1999). Improved estimation of the expected discrepancy in case of misspecification. *Econometric Theory* **15**.
- Rexstad, E. (2001). Back cover of T.M. Shenk and A.B. Franklin, *in natural resource management*. Island Press, Washington, D.C.
- Rexstad, E.A., Miller, D.D., Flather, C.H., Anderson, E.M., Hupp, J. (1988). Questionable multivariate statistical inference in wildlife studies. *Journal of Wildlife Management* **52**, 794–798.
- Rexstad, E.A., Miller, D.D., Flather, C.H., Anderson, E.M., Hupp, J. (1990). Questionable multivariate statistical inference in wildlife studies: a reply. *Journal of Wildlife Management* **54**, 189–193.
- Ripley, B.D. (1996). *Pattern recognition and neural networks*. Cambridge University Press, Cambridge, UK.
- Rissanen, J. (1989). *Stochastic complexity in statistical inquiry*. World Scientific, Computer Science, Vol 15. Singapore.
- Rissanen, J. (1996). Fisher information and stochastic complexity. *Information Theory* **42**, 40–47.
- Robert, C.P., and Casella, G. (1999). *Monte Carlo statistical methods*. John Wiley, New York, NY.
- Roecker, E.B. (1991). Prediction error and its estimation for stochastic processes. *Technometrics* **33**, 459–468.
- Ronchetti, E., and Staudte, R.G. (1994). A robust version of Mallows' test. *American Statistical Association* **89**, 550–559.
- Rosenblum, E.P. (1994). A simulation study of information theoretic hypothesis tests in one factor ANOVA. Pages 319–346. *Engineering and Scientific Applications*. Vol. 2, Proceedings of the Conference on the Frontiers of Statistical Modeling: An Informatic Perspective. Academic Publishers, Dordrecht, the Netherlands.
- Roughgarden, J. (1979). *Theory of population genetics and evolutionary dynamics: an introduction*. Macmillan Publishing Company, New York, NY.
- Royall, R.M. (1997). *Statistical evidence: a likelihood paradigm*. Chapman and Hall, London.
- Royle, J.A., and Link, W.A. (2002). Random effects and shrinkage in capture-recapture methods. *Journal of Applied Statistics* **29**, 329–351.

- Sakamoto, Y. (1991). *Categorical data analysis by AIC*. KTK Scientific Publishers, Tokyo, Japan.
- Sakamoto, Y., and Akaike, H. (1978). Analysis of cross classified data. *Journal of the Institute of Statistical Mathematics Part B* **30**, 185–197.
- Sakamoto, Y., Ishiguro, M., and Kitagawa, G. (1986). *Akaike information criterion and its application*. KTK Scientific Publishers, Tokyo, Japan.
- Santer, T.J., and Duffy, D.E. (1989). *The statistical analysis of data*. Springer-Verlag, New York, NY.
- SAS Institute Inc. (1985). SAS[®] language guide for personal computers. SAS Institute Inc., Cary, NC.
- SAS Institute. (1988). SAS/STAT[®] user's guide. Edition 6.03. SAS Institute Inc., Cary, NC.
- Sauerbrei, W., and Schumacher, M. (1992). A bootstrap resampling method for variable selection and model building: application to the Cox regression model. *Statistics in Medicine* **11**, 1375–1402.
- Sawa, T. (1978). Information criteria for discriminating among competing models. *Econometrica* **46**, 1273–1291.
- Scheiner, S.M., and Gurevitch, J. (eds.) (1993). *Design and analysis of ecological experiments*. Chapman and Hall, London.
- Schmidt, B.R., and Anholt, B.R. (1999). Analysis of survival data for common toads. *Amphibia-Reptilia* **20**, 97–108.
- Schoener, T.W. (1970). Nonsynchronous spatial overlap of lizards on three desert shrubs. *Ecology* **51**, 408–418.
- Schreuder, H.T., Gregoire, T.G., and Wood, G.B. (1993). *Statistical inference for multisource forest inventory*. John Wiley and Sons, New York, NY.
- Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics* **6**, 461–464.
- Sclove, S.L. (1987). Application of some model-selection criteria to some problems in multivariate analysis. *Psychometrika* **52**, 333–343.
- Sclove, S.L. (1994a). Small-sample and large-sample statistical models. Pages 31–39 in P. Cheeseman, and R.W. Oldford (eds.) *Selected topics in statistical modeling*. Springer-Verlag, New York, NY.
- Sclove, S.L. (1994b). Some aspects of model-selection criteria. Page 105 in S. S. Gupta (ed.) *Engineering and Scientific Applications*. Vol. 2. Proceedings of the 1994 Conference on the Frontiers of Statistical Modeling: An Information Systems Perspective. Academic Publishers, Dordrecht, the Netherlands.
- Seber, G.A.F. (1977). *Linear regression analysis*. John Wiley and Sons, New York, NY.
- Seber, G.A.F. (1984). *Multivariate observations*. John Wiley and Sons, New York, NY.
- Seber, G.A.F., and Wild, C.J. (1989). *Nonlinear regression*. John Wiley and Sons, New York, NY.

- Severini, T.A. (2000). *Likelihood methods in statistics*. Oxford University Press, Oxford, UK.
- Shannon, C.E. (1948). A mathematical theory of communication. *Journal of the American Statistical Association* **27**, 379–423 and 623–656.
- Shao, J. (1993). Linear model selection by cross-validation. *Journal of the American Statistical Association* **88**, 486–494.
- Shao, J. (1996). Bootstrap model selection. *Journal of the American Statistical Association* **91**, 655–665.
- Shao, J. (1997). An asymptotic theory for linear model selection. *Journal of the American Statistical Association* **92**, 264.
- Shao, J. and Tu, D. (1995). *The jackknife and bootstrap*. Springer-Verlag, New York.
- Shefferson, R.P., Sandercock, B.K., Proper, J., and Beissinger, S. (1991). Seed bank dynamics, dormancy and survival of a rare herbaceous perennial using mark-recapture. *Ecology* **82**, 145–156.
- Shenk, T.M., and Franklin, A.B., (eds.) (2001). *Modeling in natural resource management*. Island Press, Washington, D. C.
- Shi, P., and Tsai, C-L. (1998). A note on the unification of the Akaike information criterion. *Journal of the Royal Statistical Society, Series B*, **60**, 551–558.
- Shi, P., and Tsai, C-L. (1999). Semiparametric regression and model selection. *Statistical Planning and Inference* **77**, 341–349.
- Shibata, R. (1976). Selection of the order of an autoregressive moving average process by the information criterion. *Biometrika* **63**, 117–26.
- Shibata, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *Annals of Statistics* **8**, 147–164.
- Shibata, R. (1981). An optimal selection of regression variables. *Annals of Statistics* **9**, 872–882. Correction (1982). **69**, 492.
- Shibata, R. (1983). A theoretical view of the use of AIC. Pages 237–246. In R. A. Fisher (ed.) *Time series analysis: theory and practice*. Elsevier Science Publishers, Amsterdam, Holland, the Netherlands.
- Shibata, R. (1986). Consistency of model selection and parameter estimation. Pages 141–151 in J. Gani, and M.B. Priestly (eds.) *Essays in time series and related topics*. *Journal of Applied Probability*, Special Volume 23A.
- Shibata, R. (1989). Statistical aspects of model selection. Pages 207–227. In J. Gani, and M.B. Priestly (eds.) *From data to model*. Springer-Verlag, London.
- Shibata, R. (1997a). Bootstrap estimate of Kullback–Leibler divergence. *Statistica Sinica* **7**, 375–394.

- Shimizu, R. (1978). Entropy maximization principle and selection of the best autoregressive Gaussian process. *Annals of the Institute of Statistical Mathematics* **26**, 263–270.
- Shono, H. (2000). Efficiency of the finite correction of Akaike's information criterion. *Fisheries Science* **66**, 608–610.
- Silverman, B.W. (1982). Algorithm AS 176: kernel density estimation using a fast Fourier transform. *Applied Statistics* **31**, 93–99.
- Silvey, S.D. (1975). *Statistical inference*. Chapman and Hall, London.
- Simonoff, J.S., and Tsai, C-L. (1999). Semiparametric and additive models with an improved AIC criterion. *Journal of Computational and Graphical Statistics* **9**, 1–27.
- Skalski, J.R., Hoffman, A., and Smith, S.G. (1993). Testing the significance of individual and cohort-level covariates in animal survival studies. Pages 9–27. In P.M. North (eds.) *Marked individuals in the study of bird populations*. Birkhäuser, Basel, Switzerland.
- Skalski, J.R., and Robson, D.S. (1992). *Techniques for wildlife inventory: the analysis of capture data*. Academic Press, New York, NY.
- Smith, G.N. (1966). Basic studies on DURSBAN® insecticide. *Do not cite*.
- Smith, S.C., Skalski, J.R., Schlechte, J.W., Hoffman, A., and Casse, B. (1993). *Statistical Survival Analysis of Fish and Wildlife Tagging Studies*. Springer-Verlag, Sciences, University of Washington, Seattle, WA.
- Sober, E. (1999). Instrumentalism revisited. *Crítica* **31**, 3–39.
- Sober, E. (2001). Instrumentalism, parsimony and the Akaike framework. *the Philosophy of Science Association* (in press).
- Sommer, S., and Huggins, R.M. (1996). Variables selection using the Akaike information criterion. *Applied Statistics* **45**, 15–29.
- Soofi, E.S. (1994). Capturing the intangible concept of information. *Statistical Association* **89**, 1243–1254.
- Southwell, C. (1994). Evaluation of walked line transect counts for estimating density. *Journal of Wildlife Management* **58**, 348–356.
- Speed, T.P., and Yu, B. (1993). Model selection and prediction: nonparametric methods. *Annals of the Institute of Statistical Mathematics* **1**, 35–54.
- Spiegelhalter, D.J., Best, N.G., and Carlin, B.P., and van der Linde, A. (1998). Measures of model complexity and fit. *Journal of the Royal Statistical Society* **64**, 1–34.
- Sprott, D.A. (2000). *Statistical inference in science*. Springer Series in Statistics.
- Steel, M., and Penny, D. (2001). Parsimony, likelihood and the role of information in phylogenetics. *Molecular Biology and Evolution* **17**, 839–850.

- Sterling, T.D., Rosenbaum, W.L., and Weinkam, J.J. (1995). Publication bias: the effect of the outcome of statistical tests on the decision to publish. *American Statistician* **49**, 108–112.
- Stewart-Oaten, A. (1995). Rules and judgments in statistics: three decades of experience, 2001–2009.
- Stigler, S.M. (1986). *The history of statistics*. Harvard University Press, Cambridge, MA.
- Stoica, P., Eykhoff, P., Janssen, P., and Söderström, T. (1986). Model selection by cross-validation. *International Journal of Control* **43**, 1841–1873.
- Stone, M. (1974). Cross-validatory choice and assessment of statistical models (with discussion). *Journal of the Royal Statistical Society, Series B* **39**, 446–462.
- Stone, M. (1977). An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. *Journal of the Royal Statistical Society, Series B* **39**, 44–47.
- Stone, C.J. (1982). Local asymptotic admissibility of a generalized likelihood selection rule. *Annals of the Institute of Statistical Mathematics* **30**, 293–307.
- Stone, M., and Brooks, R.J. (1990). Continuum regression: cross-validated constructed prediction embracing ordinary least squares, partial least squares, principal components regression (with discussion). *Journal of the Royal Statistical Society, Series B* **52**, 237–269.
- Stromborg, K.L., Grue, C.E., Nichols, J.D., Hepp, G.R., Hines, R.B., et al. (1988). Postfledging survival of European starlings exposed to a sublethal dose of insecticide. *Ecology* **69**, 590–601.
- Sugiura, N. (1978). Further analysis of the data by Akaike's information criterion: finite corrections. *Communications in Statistics, Theory and Methods* **7**, 137–147.
- Takeuchi, K. (1976). Distribution of informational statistics and a central limit theorem. *Suri-Kagaku* (Mathematic Sciences) **153**, 12–18. (In Japanese).
- Taub, F.B. (1993). Book review: Estimating ecological risks. *Ecology* **74**, 1000–1001.
- Taubes, G. (1995). Epidemiology faces its limits. *Science* **269**, 164–165.
- Thabane, L., and Haq, M.S. (1999). On Bayesian selection of the best model using the Kullback–Leibler divergence measure. *Statistica Neerlandica* **53**, 1–14.
- Thompson, S.K. (1992). *Sampling*. Wiley, New York, NY.
- Thompson, M.E. (1997). *Theory of sample surveys*. Chapman and Hall, New York, NY.
- Thompson, W.L., and Lee, D.C. (2000). Modeling relationships between physical attributes and snorkel counts of chinook salmon and steelhead pike. *Journal of Fisheries and Aquatic Sciences* **57**, 1834–1842.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* **58**, 267–288.

- Tong, H. (1994). Akaike's approach can yield consistent order determination. In H. Bozdogan (ed.) *Engineering and Scientific Applications*. Vol. 1. First US/Japan Conference on the Frontiers of Statistical Modeling and Data Analysis. Kluwer Academic Publishers, Dordrecht, the Netherlands.
- Tukey, J.W. (1980). We need both exploratory and confirmatory. *The American Statistician* **34**, 23–25.
- Turchin, P., and Batzli, G.O. (2001). Availability of food and the population dynamics of arvicoline rodents. *Ecology* **82**, 1521–1534.
- Ullah, A. (1996). Entropy, divergence and distance measures with ecological applications. *Journal of Statistical Planning and Inference* **49**, 137–162.
- Umbach, D.M., and Wilcox, A.J. (1996). A technique for measuring the relative importance of useful features of birthweight distributions. *Statistics in Medicine* **15**, 1097–1107.
- Venter, J.H., and Snyman, J.L.J. (1995). A note on the generalized cross-validation criterion in linear model selection. *Biometrika* **82**, 215–219.
- Ver Hoef, J.M. (1996). Parametric empirical Bayes methods for ecological data. *Ecological Applications* **6**, 1047–1055.
- Wade, P.R. (2000). Bayesian methods in conservation biology. *Conservation Biology* **14**, 1308–1316.
- Walters, C.J. (1996). Computers and the future of fisheries. Pages 22–33. In J. G. Garman and E. Mckness (eds.) *Computers in fisheries research*. Chapman and Hall, New York.
- Wang, C. (1993). *Sense and nonsense of statistical inference: comments on the literature*. Marcel Dekker, Inc., New York, NY.
- Wang, P., Puterman, M. L., Cockburn, I., and Le, N. (1996). Mixed-effects models with covariate dependent rates. *Biometrics* **52**, 381–400.
- Wasserman, L. (2000). Bayesian model selection and model averaging. *Mathematical Psychology* **44**, 92–107.
- Wedderburn, R.W.M. (1974). Quasi-likelihood functions, generalized linear models and the Gauss–Newton method. *Biometrika* **61**, 439–447.
- Wehrl, A. (1978). General properties of entropy. *Reviews of Modern Physics* **50**, 221–238.
- Weiner, J. (1995). On the practice of ecology. *Journal of Ecology* **83**, 103–113.
- Weisberg, S. (1985). *Applied linear regression*. 2nd ed. Wiley, New York.
- Wel, J. (1975). Least squares fitting of an elephant. *Chemtech* Feb. 1975, 10–11.
- Westfall, P.H., and Young, S.S. (1993). *Resampling-based multiple testing: methods for p-value adjustment*. John Wiley and Sons, New York.
- White, G.C. (1983). Numerical estimation of survival rates from capture-recapture and biotelemetry data. *Journal of Wildlife Management* **47**, 716–728.

- White, G.C., Anderson, D.R., Burnham, K.P., and Otis, D.L. (1982). *Capture-recapture and removal methods for sampling closed populations*. Los Alamos Laboratory Report LA-8787-NERP, Los Alamos, NM.
- White, G.C., and Burnham, K.P. (1999). Program MARK-surveys: Estimating parameters of populations of marked animals. *Bird Study* **46**.
- White, G.C., Burnham, K.P., and Anderson, D.R. (2001). Advances in capture-recapture MARK. in R. Fields (ed.) *Integrating People and Wildlife for Sustainable Management*. Proceedings of the Second International Wildlife Management Conference, Society, Bethesda, MD.
- White, H. (1994). *Estimation, inference and specification analysis*. Cambridge University Press, Cambridge, UK.
- White, H. (2000). A reality check for data snooping. *Econometrica* **68**, 1097–1119.
- Williams, B.K., Nichols, J.D., and Conroy, M.J. (2002). *Analysis and design of capture-recapture populations: modeling, estimation, and decision making*. Academic Press, San Diego, CA.
- Williams, D.A. (1982). Extra-binomial variation in logistic linear models. *Biometrics* **31**, 144–148.
- Wood, S.N. and Thomas, M.B. (1999). Super-sensitivity to structural model misspecification. *Proceedings of the Royal Society* **266**, 565–570.
- Woods, H., Steinour, H.H., and Starke, H.R. (1932). Effect of curing temperature and cement on heat evolved during hardening. *Industrial and Engineering Chemistry* **24**, 1207–1214.
- Ye, J. (1998). On measuring and correcting the effects of data mining on statistical inference. *Journal of the American Statistical Association* **93**, 120–131.
- Yockey, H.P. (1992). *Information theory and molecular biology*. Cambridge University Press.
- Yoccoz, N.G. (1991). Use, overuse, and misuse of significance tests in ecology. *Bulletin of the Ecological Society of America* **72**, 1–10.
- Young, L.J., and Young, J.H. (1998). *Statistical ecology*. Kluwer Academic Publishers, London, UK.
- Yu, B. (1996). Minimum description length principle: a review. In *Proceedings of the Conference on Information Sciences and Systems*, Princeton University.
- Yu, B. (1999). Coding and model selection: a brief tutorial on the minimum description length principle. *Statistical Computing and Graphics* **9**, 1, 27–41.
- Zablan, M.A. (1993). Evaluation of sage grouse banding program in Colorado. M.S. thesis, Colorado State University, Fort Collins, CO.

Zellner, A., Keuzenkamp, H. A., and McAleer, M., (eds.) (2001). *Simplicity in modelling: keeping it sophisticatedly simple*. Cambridge University Press, Cambridge, UK.

Zhang, P. (1992a). Inferences after variable selection in linear regression. *Journal of the Royal Statistical Society* **79**, 741–746.

Zhang, P. (1992b). On the distributional properties of model selection criteria. *The American Statistical Association* **87**, 732–737.

Zhang, P. (1993a). Model selection via multifold cross-validation. *Journal of the Royal Statistical Society* **56**, 299–313.

Zhang, P. (1993b). On the convergence rate of model selection criteria. *Journal of the Royal Statistical Society, Part A—Theory and Methods* **22**, 2765–2775.

Zhang, P. (1994). On the choice of penalty term in generalized FPE. *Journal of the Royal Statistical Society* **57**, 1–13.
in P. Cheeseman, and R.W. Oldford (eds.) *Selecting models from data*. John Wiley, New York, NY.

Zucchini, W. (2000). An introduction to model selection. *Journal of the Royal Statistical Society* **44**, 41–61.

Index

Acute effect, 113
adjusted coefficient of multiple
determination, 37
AIC, 60, 61, 96, 385, 444
AIC differences, 70, 97
AIC_c, 66, 326, 374, 445
Akaike weights, 75, 303, 439
Akaike's information criterion, 60, 96
all possible models, 89, 147
all-subsets selection, 225, 250
analysis paradigms, 203
ANODEV, 329
anolis lizards, 137, 146
approximate model probabilities, 75
approximating model, 23
averaging across model parameters, 151

Bayes factor, 170, 302, 305
Bayesian, 153
Bayesian approaches, 77, 98
Bayesian model selection, 183, 301
best approximating model, 32, 287
bias–variance tradeoff, 31, 116, 130
bias-corrected log-likelihood, 64
BIC, 270, 286, 290, 293, 303
Boltzmann's entropy, 52, 76

bootstrap, 90, 100
276, 372
bootstrap sample
bootstrap-select

CAICF, 287
candidate model
cement data, 100
cement hardening
chain binomial m
chronic effects, 1
confidence interv
confidence set fo
169
confidence set of
critical thinking,
cross entropy, 54
cross-validation,
curse of dimensi
cutoff value, 171

Data analysis, 47
data dredging, 2,
140, 147
data snooping, 4
 Δ_i , 71
density estimatio

- dimension-consistent criteria, 447
- discrimination information, 52
- Durban storm data, 187
- Elephant**, 29
- empirical Bayes, 313
- entropy, 86, 427, 442
- entropy maximization principle, 49, 86
- estimated models, 25
- estimated selection probabilities, 172
- European starlings, 112
- evidence, 146
- evidence ratios, 77, 97, 99, 111, 198, 447
- exploratory data analysis, 39, 84
- exponential family, 380
- extrabinomial variation, 67
- Field protocol**, 440
- Fisher information matrix, 358, 359, 426
- fixed effects, 310
- Flather's models, 94
- flour beetle mortality, 195
- Fourier series model, 188
- Freedman's paradox, 17, 38, 72
- full reality, 22, 438
- G-statistic**, 336
- gamma distribution, 327
- generalized Akaike weights, 76
- generalized cross-validation, 62, 316
- generalized entropy, 51
- generalized estimating equations, 316
- generalized linear models, 116
- generalized smoothing, 316
- generating model, 20, 26, 111, 117, 155
- global model, 17, 26, 88, 100, 104, 125, 146, 202
- goodness-of-fit, 68, 262, 305, 309, 336
- H-likelihood**, 317
- higher-order interactions, 140
- historic emphasis on hypothesis testing, 41
- hypothesis testing, 121, 131, 337

- importance of m
- importance of th
- inference philos
- inferential uncer
- infinite-dimensi
- information, 9, 2
- information theo
- interval scale, 59
- iterative process

Judgment, 333

K-L model prior
 Kullback–Leible
 Kullback–Leible
 96, 352, 4

Least squares, 6,
 least-absolute de
 likelihood functi
 likelihood of a m
 likelihood ratio t
 likelihood theory
 linear regression
 log-likelihood, 9
 logistic regressio
 loglinear model
 138

Markov chain M
 masu salmon, 14
 maxent principle
 maximized log-l
 97, 144
 maximum entrop
 maximum likelih
 mean square erro
 meaningful ques
 misspecification,
 mixture models,
 model-averaged
 149
 model-averaged
 model averaging
 model interpreta
 model probabil

- model selection methods, 35
- model selection uncertainty, 3, 45, 77, 99, 124, 149, 153, 176, 203, 271, 439, 449
- models of fish growth, 142
- Monte Carlo data, 113
- Monte Carlo studies, 89
- multidimensional contingency tables, 140
- multimodel inference, 4, 47, 80, 99, 149, 438
- multinomial distribution, 397
- multiple comparison, 83
- multiple regression, 268
- multiple testing problem, 133
- multiple working hypotheses, 96, 201, 440
- multivariate data, 424

- Nestling starling, 111
- nonlinear least squares, 108
- nonnested models, 88
- Northern Spotted Owl, 146
- null hypothesis testing, 83, 84, 98, 439
- number of parameters, 12

- Occam's razor**, 29, 443
- overdispersed count data, 67
- overdispersion, 67, 200, 305
- overdispersion parameter, 70
- overfitting, 32, 33, 35, 250, 417

- Parameters**, 330
- parsimonious model, 143
- philosophical issues, 96, 219
- philosophy, 220, 248, 437
- philosophy of science, 145
- pine wood data, 183
- pivotal quantity, 170
- post hoc power, 42
- prediction, 35, 150
- principle of parsimony, 29, 31, 48, 62, 97, 443
- probability model, 7
- profile likelihood, 165
- profile likelihood interval, 11, 166

- QAIC_c, 70, 308,
- quantile regression
- Quasi-likelihood
- quasi-true model

- Random coefficient**
- random effects,
- ranking alternati
- recommendation
- regression mode
- relative expected
- relative importan
- variable,
- relative importan
- 167, 168,
- residual sum of s
- robust inferences

- Sage grouse**, 146
- sage grouse surv
- sample size, 332
- Sampling variati
- scaling, 149
- scaling alternativ
- science culture,
- scientific hypoth
- scientific questio
- selection bias, 24
- selection target,
- Shannon's conce
- shrinkage, 152, 2
- shrinkage estima
- simulated data, 1
- simulated ecosys
- simulation, 229
- small sample AI
- spurious correlat
- spurious effects,
- statistical eviden
- statistical pitfalls
- statistical science
- stepwise selectio
- strategy, 2, 18, 1
- strength of eviden
- strict experiment
- study objective,

tapering effects, 21, 44, 89, 115, 287
thought experiment, 246
TIC, 65, 96, 188, 368, 371, 372, 385,
444
trace function, 361
transformations of the response variable,
81
treatment effects, 119
trivial null hypotheses, 42
true model, 20, 30, 32, 89, 97, 143, 438,
447
truth or full reality, 23, 438

Unconditional confidence intervals, 164
unconditional estimates of precision, 178

unconditional va
underfitting, 32,

Valid inference,
variable selection
244, 267
variance compon
variance–covaria
variance inflation
variance inflation

 w_i , 75
weight of eviden
447
weighting predic