Maximum Likelihood Estimation and Inference

With Examples in R, SAS and ADMB

Russell B. Millar

Plug in

Pseudo Bayes Bootstrap Bayes, prior 1 Bayes, prior 2



STATISTICS IN PRACTICE

Maximum Likelihood Estimation and Inference

Statistics in Practice

Series Advisors

Human and Biological Sciences Stephen Senn

University of Glasgow, UK

Earth and Environmental Sciences

Marian Scott University of Glasgow, UK

Industry, Commerce and Finance

Wolfgang Jank University of Maryland, USA

Statistics in Practice is an important international series of texts which provide detailed coverage of statistical concepts, methods and worked case studies in specific fields of investigation and study.

With sound motivation and many worked practical examples, the books show in down-to-earth terms how to select and use an appropriate range of statistical techniques in a particular practical field within each title's special topic area.

The books provide statistical support for professionals and research workers across a range of employment fields and research environments. Subject areas covered include medicine and pharmaceutics; industry, finance and commerce; public services; the earth and environmental sciences, and so on.

The books also provide support to students studying statistical courses applied to the above areas. The demand for graduates to be equipped for the work environment has led to such courses becoming increasingly prevalent at universities and colleges.

It is our aim to present judiciously chosen and well-written workbooks to meet everyday practical needs. Feedback of views from readers will be most valuable to monitor the success of this aim.

A complete list of titles in this series appears at the end of the volume.

Maximum Likelihood Estimation and Inference

With Examples in R, SAS and ADMB

Russell B. Millar

Department of Statistics, University of Auckland, New Zealand



This edition first published 2011 © 2011 John Wiley & Sons, Ltd

Registered office

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, United Kingdom

For details of our global editorial offices, for customer services and for information about how to apply for permission to reuse the copyright material in this book please see our website at www.wiley.com.

The right of the author to be identified as the author of this work has been asserted in accordance with the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by the UK Copyright, Designs and Patents Act 1988, without the prior permission of the publisher.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Designations used by companies to distinguish their products are often claimed as trademarks. All brand names and product names used in this book are trade names, service marks, trademarks or registered trademarks of their respective owners. The publisher is not associated with any product or vendor mentioned in this book. This publication is designed to provide accurate and authoritative information in regard to the subject matter covered. It is sold on the understanding that the publisher is not engaged in rendering professional services. If professional advice or other expert assistance is required, the services of a competent professional should be sought.

Library of Congress Cataloging-in-Publication Data

```
Millar, R. B. (Russell B.)
Maximum likelihood estimation and inference : with examples in R, SAS, and ADMB /
Russell B. Millar.
p. cm.
Includes bibliographical references and index.
ISBN 978-0-470-09482-2 (hardback)
1. Estimation theory. 2. Chance–Mathematical models. I. Title.
QA276.8.M55 2011
519.5'44–dc22
```

2011013225

A catalogue record for this book is available from the British Library.

Print ISBN: 978-0-470-09482-2 ePDF ISBN: 978-0-470-09483-9 oBook ISBN: 978-0-470-09484-6 ePub ISBN: 978-1-119-97771-1 Mobi ISBN: 978-1-119-97772-8

Set in 10.25/12pt Times by Thomson Digital, Noida, India

Contents

Pr	eface		xiii
Pa	rt I 🛛	PRELIMINARIES	1
1	A ta	ste of likelihood	3
	1.1	Introduction	3
	1.2	Motivating example	4
		1.2.1 ML estimation and inference for the binomial	4
		1.2.2 Approximate normality versus likelihood ratio	8
	1.3	Using SAS, R and ADMB	9
		1.3.1 Software resources	11
	1.4	Implementation of the motivating example	11
		1.4.1 Binomial example in SAS	11
		1.4.2 Binomial example in R	14
		1.4.3 Binomial example in ADMB	15
	1.5	Exercises	17
2	Esse	ntial concepts and iid examples	18
	2.1	Introduction	18
	2.2	Some necessary notation	19
		2.2.1 MLEs of functions of the parameters	22
	2.3	Interpretation of likelihood	23
	2.4	IID examples	25
		2.4.1 IID Bernoulli (i.e. binomial)	25
		2.4.2 IID normal	26
		2.4.3 IID uniform	28
		2.4.4 IID Cauchy	29
		2.4.5 IID binormal mixture model	31
	2.5	Exercises	33
Pa	rt II	PRAGMATICS	37
3	Нур	othesis tests and confidence intervals or regions	39
	3.1	Introduction	39

vi	CC	ONTENTS	
	3.2	Approximate normality of MLEs	40
		3.2.1 Estimating the variance of $\hat{\theta}$	41
	3.3	Wald tests, confidence intervals and regions	43
		3.3.1 Test for a single parameter	43
		3.3.2 Test of a function of the parameters	43
		3.3.3 Joint test of two or more parameters	44
		3.3.4 In R and SAS: Old Faithful revisited	45
	3.4	Likelihood ratio tests, confidence intervals and regions	49
		3.4.1 Using R and SAS: Another visit to Old Faithful	51
	3.5	Likelihood ratio examples	54
		3.5.1 LR inference from a two-dimensional contour plot	54
		3.5.2 The <i>G</i> -test for contingency tables	56
	3.6	Profile likelihood	57
		3.6.1 Profile likelihood for Old Faithful	58
	3.7	Exercises	59
4	Wha	at you really need to know	64
	4.1	Introduction	64
	4.2	Inference about $g(\boldsymbol{\theta})$	65
		4.2.1 The delta method	65
		4.2.2 The delta method applied to MLEs	67
		4.2.3 The delta method using R, SAS and ADMB	70
		4.2.4 Delta method examples	72
	4.3	Wald statistics – quick and dirty?	75
		4.3.1 Wald versus likelihood ratio revisited	77
		4.3.2 Pragmatic considerations	78
	4.4	Model selection	79
		4.4.1 AIC	79
	4.5	Bootstrapping	81
		4.5.1 Bootstrap simulation	82
		4.5.2 Bootstrap confidence intervals	83
		4.5.3 Bootstrap estimate of variance	85
		4.5.4 Bootstrapping test statistics	85
		4.5.5 Bootstrap pragmatics	86
		4.5.6 Bootstrapping Old Faithful	86
		4.5.7 How many bootstrap simulations is enough?	90
	4.6	Prediction	91
		4.6.1 The plug-in approach	92
		4.6.2 Predictive likelihood	93
		4.6.3 Bayesian prediction	93
		4.6.4 Pseudo-Bayesian prediction	94
		4.6.5 Bootstrap prediction	94
	4.7	Things that can mess you up	95
		4.7.1 Multiple maxima of the likelihood	95
		4.7.2 Lack of convergence	96

CONTENTS	vii

		4.7.3 Parameters on the boundary of the parameter space	96
		4.7.4 Insufficient sample size	98
	4.8	Exercises	98
5	Max	imizing the likelihood	101
	5.1	Introduction	101
	5.2	The Newton-Raphson algorithm	103
	5.3	The EM (Expectation–Maximization) algorithm	104
		5.3.1 The simple form of the EM algorithm	105
		5.3.2 Properties of the EM algorithm	107
		5.3.3 Accelerating the EM algorithm	111
		5.3.4 Inference	112
	5.4	Multi-stage maximization	113
		5.4.1 Efficient maximization via profile likelihood	114
		5.4.2 Multi-stage optimization	116
	5.5	Exercises	118
6	Som	e widely used applications of maximum likelihood	121
	6.1	Introduction	121
	6.2	Box-Cox transformations	122
		6.2.1 Example: the Box and Cox poison data	123
	6.3	Models for survival-time data	125
		6.3.1 Notation	126
		6.3.2 Accelerated failure-time model	127
		6.3.3 Parametric proportional hazards model	128
		6.3.4 Cox's proportional hazards model	130
		6.3.5 Example in R and SAS: Leukaemia data	132
	6.4	Mark–recapture models	134
		6.4.1 Hypergeometric likelihood for integer valued N	136
		6.4.2 Hypergeometric likelihood for $N \in \mathbb{R}^+$	137
		6.4.3 Multinomial likelihood	138
		6.4.4 Closing remarks	140
	6.5	Exercises	141
7	Gen	eralized linear models and extensions	143
	7.1	Introduction	143
	7.2	Specification of a GLM	144
		7.2.1 Exponential family distribution	144
		7.2.2 GLM formulation	146
	7.3	Likelihood calculations	148
	7.4	Model evaluation	149
		7.4.1 Deviance	149
		7.4.2 Model selection	151
		7.4.3 Residuals	151
		7.4.4 Goodness of fit	152

viii	CONTENTS

	7.5	Case study 1: Logistic regression and inverse prediction in R \ldots	154
		7.5.1 Size-selectivity modelling in R	155
	7.6	Beyond binomial and Poisson models	161
		7.6.1 Quasi-likelihood and quasi-AIC	162
		7.6.2 Zero inflation and the negative binomial	164
	7.7	Case study 2: Multiplicative vs additive models	
		of over-dispersed counts in SAS	167
		7.7.1 Background	167
		7.7.2 Poisson and quasi-Poisson fits	169
		7.7.3 Negative binomial fits	171
	7.8	Exercises	173
8	Qua	si-likelihood and generalized estimating equations	175
	8.1	Introduction	175
	8.2	Wedderburn's quasi-likelihood	177
		8.2.1 Quasi-likelihood analysis of barley blotch data in R	177
	8.3	Generalized estimating equations	181
		8.3.1 GEE analysis of multi-centre data in SAS	183
	8.4	Exercises	187
9	ML	inference in the presence of incidental parameters	188
	9.1	Introduction	188
		9.1.1 Analysis of paired data: an intuitive use of conditional	
		likelihood	190
	9.2	Conditional likelihood	192
		9.2.1 Restricted maximum likelihood	197
	9.3	Integrated likelihood	198
		9.3.1 Justification	199
		9.3.2 Uses of integrated likelihood	200
	9.4	Exercises	201
10	Late	nt variable models	202
	10.1	Introduction	202
	10.2	Developing the likelihood	203
	10.3	Software	204
		10.3.1 Background	204
		10.3.2 The Laplace approximation and Gauss-Hermite quadrature	206
		10.3.3 Importance sampling	208
		10.3.4 Separability	209
		10.3.5 Overview of examples	210
	10.4	One-way linear random-effects model	210
		10.4.1 SAS	212
		10.4.2 R	215
		10.4.3 ADMB	216
	10.5	Nonlinear mixed-effects model	217

viii

CONTENTS ix

233

	10.5.1 SAS	219
	10.5.2 ADMB	220
10.6	Generalized linear mixed-effects model	221
	10.6.1 R	222
	10.6.2 SAS	224
	10.6.3 ADMB	224
	10.6.4 GLMM vs GEE	225
10.7	State-space model for count data	227
10.8	ADMB template files	228
	10.8.1 One-way linear random-effects model using REML	228
	10.8.2 Nonlinear crossed mixed-effects model	229
	10.8.3 Generalized linear mixed model using GREML	230
	10.8.4 State-space model for count data	231
10.9	Exercises	232

Part III THEORETICAL FOUNDATIONS

11	Cramér-Rao inequality and Fisher information	235
	11.1 Introduction	235
	11.1.1 Notation	236
	11.2 The Cramér-Rao inequality for $\theta \in \mathbb{R}$	236
	11.3 Cramér-Rao inequality for functions of θ	239
	11.4 Alternative formulae for $\mathcal{I}(\theta)$	241
	11.5 The iid data case	243
	11.6 The multi-dimensional case, $\theta \in \Theta \subset \mathbb{R}^s$	243
	11.6.1 Parameter orthogonality	244
	11.6.2 Alternative formulae for $\mathcal{I}(\theta)$	245
	11.6.3 Fisher information for re-parameterized models	246
	11.7 Examples of Fisher information calculation	247
	11.7.1 Normal (μ, σ^2)	247
	11.7.2 Exponential family distributions	249
	11.7.3 Linear regression model	251
	11.7.4 Nonlinear regression model	252
	11.7.5 Generalized linear model with canonical link function	252
	11.7.6 $\operatorname{Gamma}(\alpha, \beta)$	252
	11.8 Exercises	253
12	Asymptotic theory and approximate normality	256
	12.1 Introduction	256
	12.2 Consistency and asymptotic normality	257
	12.2.1 Asymptotic normality, $\theta \in \mathbb{R}$	261
	12.2.2 Asymptotic normality: $\theta \in \mathbb{R}^3$	264
	12.2.3 Asymptotic normality of $g(\theta) \in \mathbb{R}^{p}$	266
	12.2.4 Asymptotic normality under model misspecification	267

		12.2.5 Asymptotic normality of M-estimators	268
		12.2.6 The non-iid case	271
	12.3	Approximate normality	271
		12.3.1 Estimation of the approximate variance	273
		12.3.2 Approximate normality of M-estimators	274
	12.4	Wald tests and confidence regions	276
		12.4.1 Wald test statistics	276
		12.4.2 Wald confidence intervals and regions	278
	12.5	Likelihood ratio test statistic	280
		12.5.1 Likelihood ratio test: $\theta \in \mathbb{R}$	280
		12.5.2 Likelihood ratio test for $\theta \in \mathbb{R}^s$, and $g(\theta) \in \mathbb{R}^p$	281
	12.6	Rao-score test statistic	281
	12.7	Exercises	283
13	Tool	s of the trade	286
	13.1	Introduction	286
	13.2	Equivalence of tests and confidence intervals	286
	13.3	Transformation of variables	287
	13.4	Mean and variance conditional identities	288
	13.5	Relevant inequalities	289
		13.5.1 Jensen's inequality for convex functions	289
		13.5.2 Cauchy-Schwarz inequality	291
	13.6	Asymptotic probability theory	291
		13.6.1 Convergence in distribution and probability	292
		13.6.2 Properties	294
		13.6.3 Slutsky's theorem	295
		13.6.4 Delta theorem	295
	13.7	Exercises	297
14	Func	lamental paradigms and principles of inference	299
	14.1	Introduction	299
	14.2	Sufficiency principle	300
		14.2.1 Finding sufficient statistics	301
		14.2.2 Examples of the sufficiency principle	303
	14.3	Conditionality principle	304
	14.4	The likelihood principle	306
		14.4.1 Relationship with sufficiency and conditionality	308
	14.5	Statistical significance versus statistical evidence	309
	14.6	Exercises	311
15	Misc	rellanea	313
-	15.1	Notation	313
	15.2	Acronyms	315
	15.3	Do you think like a frequentist or a Bayesian?	315
	15.4	Some useful distributions	316

CONTENTS xi

15.4.1 Discrete distributions	316
15.4.2 Continuous distributions	318
15.5 Software extras	321
15.5.1 R function Plkhci for likelihood ratio confidence intervals	321
15.5.2 R function Profile for calculation of profile	
likelihoods	322
15.5.3 SAS macro Plkhci for likelihood ratio confidence	
intervals	322
15.5.4 SAS macro Profile for calculation of profile	
likelihoods	323
15.5.5 SAS macro DeltaMethod for application of	
the delta method	323
15.6 Automatic differentiation	323
Appendix: Partial solutions to selected exercises	325
Bibliography	337
Index	345

Preface

Likelihood has a fundamental role in the field of statistical inference, and this text presents a fresh look at the pragmatic concepts, properties, and implementation of statistical estimation and inference based on maximization of the likelihood. The supporting theory is also provided, but for readability is kept separate from the pragmatic content.

The properties of maximum likelihood inference that are presented herein are from the point of view of the classical frequentist approach to statistical inference. The Bayesian approach provides another paradigm of likelihood-based inference, but is not covered here, though connections to Bayesian methodology are made where relevant. Leaving philosophical arguments aside (but see Chapter 14), one of the basic choices to be made before any analysis is to determine the most appropriate paradigm to use in order to best answer the research question and to meet the needs of scientific colleagues or clients. This text will aid this choice, by showing the best of what can be done using maximum likelihood under the frequentist paradigm.

The level of presentation is aimed at the reader who has already been exposed to an undergraduate course on the standard tools of statistical inference such as linear regression, ANOVA and contingency table analysis, but who has discovered, through curiosity or necessity, that the world of real data is far more diverse than that assumed by these models. For this reason, these standard techniques are not given any special attention, and appear only as examples of maximum likelihood inference where applicable. It will be assumed that the reader is familiar with basic concepts of statistical inference, such as hypothesis tests and confidence intervals.

Much of this text is focused on the presentation of tools, tricks, and bits of R, SAS and ADMB code that will be useful in analyzing real data, and these are demonstrated through numerous examples. Pragmatism is the key motivator throughout. So, for example, software utilities have been provided to ease the computational burden of the calculation of likelihood ratio confidence intervals.

Explanation of SAS and R code is made at a level that assumes the reader is already familiar with basic programming in these languages, and hence is comfortable with their general syntax, and with tasks such as data manipulation. ADMB is a somewhat different beast, and (at the present time) will be totally unfamiliar to the majority of readers. It is used sparingly. However, when the desired model is sufficiently complex or non-standard, ADMB provides a powerful choice for its implementation.

This text is divided into three parts:

Part I: Preliminaries: Chapters 1-2

The preliminaries in this part can be skimmed by the reader who is already familiar with the basic notions and properties of maximum likelihood. However, it should be noted that the simple binomial example in Chapter 1 is used to introduce several key tools, including the Wald and likelihood ratio methods for tests and confidence intervals. Their implementation in R, SAS and ADMB is via general purpose code that is easily extended to more challenging models in later chapters. Chapter 2 looks at examples of maximum likelihood modelling of independent and identically distributed data. Despite being iid data, some of these examples are nonstandard and demonstrate curious phenomena, including likelihoods that have no maximum or have multiple maxima. This chapter also sets up the basic notation employed throughout subsequent chapters.

Part II: Pragmatics: Chapters 3-10

This part covers the relevant practical application of maximum likelihood, including cutting-edge developments in methodology for coping with nuisance parameters (e.g., GREML - generalized restricted maximum likelihood) and latent variable models. The well-established methodology for construction of hypothesis tests and confidence intervals is presented in Chapter 3. But, knowing how to do the calculations isn't the same as actually working with real data, and it is Chapter 4 that really explains how it should be done. This chapter includes model selection, bootstrapping, prediction, and coverage of techniques to handle nonstandard situations. Chapter 5 looks at methods for maximizing the likelihood (especially stubborn ones), and Chapter 6 gives a flavour of some common applications, including survival analysis, and mark-recapture models. Generalized linear models are covered in Chapter 7, with some attention to variants such as the simple over-dispersion form of quasi-likelihood, and the use of nonstandard link functions. Chapter 8 covers some of the general variants of likelihood that are in common use, including quasi-likelihood and generalized estimating equations. Chapter 9 looks at modified forms of likelihood in the presence of nuisance parameters, including conditional, restricted and integrated likelihood. Chapter 10 looks at the use of latent-variable models (e.g., mixed-effects and state-space models). For arbitrary forms of such models, this is one place where ADMB comes to the fore.

Part III: Theoretical foundations: Chapters 11-14

The theory and associated tools that are required to formally establish the properties of maximum likelihood methodology are provided here. This part

provides completeness for those readers who wish to understand the true meaning of statistical concepts such as efficiency and large-sample asymptotics. In addition, Chapter 14 looks at some of the fundamental issues underlying a statistical paradigm based on likelihood.

Chapter 15 contains a collection of notation, descriptions of common statistical distributions, and details of software utilities. This text concludes with partial solutions to a selection of the exercises from the end of each chapter.

This book includes an accompanying website. Please visit www.wiley.com/go/Maximum_likelihood

Acknowledgements

I am extremely thankful to the many cohorts of statistics students at the University of Auckland who have perused and critiqued the parts of this text that have been used in my statistical inference course. This work was greatly assisted by a University of Auckland Research Fellowship. My greatest thanks are for the unwavering support of Professor Marti Anderson at Massey University, Auckland, and for her dedication at reading through the entire first draft.

Russell B. Millar Auckland, March 2011

Part I PRELIMINARIES

1

A taste of likelihood

When it is not in our power to follow what is true, we ought to follow what is most probable. – René Descartes

1.1 Introduction

The word *likelihood* has its origins in the late fourteenth century (Simpson and Weiner 1989), and examples of its usage include as an indication of probability or promise, or grounds for probable inference. In the early twentieth century, Sir Ronald Fisher (1890–1962) presented the 'absolute criterion' for parameter estimation (Fisher 1912), and some nine years later he gave this criterion the name *likelihood* (Fisher 1921, Aldrich 1997). Fisher's choice of terminology was ideal, because the centuries-old interpretation of the word *likelihood* is also applicable to the formal statistical definition of likelihood that is used throughout this book.

Here, likelihood is used within the traditional framework of frequentist statistics, and maximum likelihood (ML) is presented as a general-purpose tool for inference, including the evaluation of statistical significance, calculation of confidence intervals (CIs), model assessment, and prediction. The frequentist theory underpinning the use of maximum likelihood is covered in Part III, where it is seen that maximum likelihood estimators (MLEs) have optimal properties for sufficiently large sample sizes. It is for this reason that maximum likelihood is the most widely used form of traditional parametric inference. The pragmatic use of ML inference is the primary focus of this book and is covered in Part II. The reader who is already comfortable with the concept of likelihood and its basic properties can proceed to Part II directly.

Maximum Likelihood Estimation and Inference: With Examples in R, SAS and ADMB, First Edition. Russell B. Millar. © 2011 John Wiley & Sons, Ltd. Published 2011 by John Wiley & Sons, Ltd.

4 PRELIMINARIES

Likelihood is also a fundamental concept underlying other statistical paradigms, especially the Bayesian approach. Bayesian inference is not considered here, but consideration of the philosophical distinctions between frequentist and Bayesian statistics is examined in Chapter 14. In addition, it is seen that some maximum likelihood methodology can be motivated using Bayesian considerations. This includes techniques for prediction (Section 4.6), and the use of integrated likelihood (Section 9.3).

A simple binomial example (Example 1.1) is used in Section 1.2 to motivate and demonstrate many of the essential properties of likelihood that are developed in later chapters. In this example, the likelihood is simply the probability of observing y = 10 successes from 100 trials. The fundamental conceptual point is that likelihood expresses the probability of observing 10 successes as a function of the unknown success probability p. That is, the likelihood function does not consider other values of y. It takes the knowledge that y = 10 was the observed number of successes and it uses the binomial probability of the outcome y = 10, evaluated at different possible values of p, to judge the relative likelihood of those different values of p.

1.2 Motivating example

Throughout this book, adding a zero subscript to a parameter (e.g. p_0) is used generically to denote a specified value of the parameter. This is typically either its true unknown value, or a hypothesized value.

1.2.1 ML estimation and inference for the binomial

Example 1.1 applies ML methodology to the binomial model in order to obtain the MLE of the binomial probability, the standard error of the MLE, and confidence intervals. This example is revisited and extended in subsequent chapters. For example, Sections 4.2.2 and 4.3.1 look at issues concerning approximate normality of the MLE, and Example 4.10 considers prediction of a new observation from the binomial distribution.

Example 1.1. Binomial. A random sample of one hundred trials was performed and ten resulted in success. What can be inferred about the unknown probability of success, p_0 ?

For any potential value of p ($0 \le p \le 1$) for the probability of success, the probability of *y* successes from *n* trials is given by the binomial probability formula (Section 15.4.1). With y = 10 successes from n = 100 trials, this is

$$L(p) = \operatorname{Prob}(10 \text{ successes})$$

= $\frac{100!}{90! \ 10!} \ p^{10}(1-p)^{90}$
= $1.731 \times 10^{13} \times \ p^{10}(1-p)^{90}, \quad 0 \le p \le 1.$ (1.1)



Figure 1.1 Binomial likelihood for 10 successes from 100 trials.

The above probability is the likelihood, and has been denoted L(p) to make its dependence on p explicit.

A plot of L(p) (Figure 1.1) shows it to be unimodal with a peak at p = 0.1. This is the MLE and will be denoted \hat{p} . For the binomial model, the MLE of the probability of success is always the observed proportion of successes $\hat{p} = y/n$ (Example 2.5).

The curve in Figure 1.1 looks somewhat like the bell-shaped curve of the normal density function. However, it is not a density (it is a likelihood function) and nor is it bell-shaped. On close inspection it can be seen that the curve is slightly right-skewed.

Box 1.1

In the above example, the MLE \hat{p} is simply a point-estimate of p_0 , and is of limited use without any sense of how reliable it is. For example, it would be more meaningful to have a range of plausible values of the unknown p_0 , or to know if some pre-specified value, e.g. $p_0 = 0.5$, was reasonable. Such questions can be addressed by examining the shape of the likelihood function, or more usually, the shape of the log-likelihood function.

The (natural) log of the likelihood function is used far more predominantly in likelihood inference than the likelihood function itself, for several good reasons:

- 1. The likelihood and log-likelihood are both maximized by the MLE.
- Likelihood values are often extremely small (but can also be extremely large) depending on the model and amount of data. This can make numerical optimization of the likelihood highly problematic, compared to optimization of the log-likelihood.

6 PRELIMINARIES

- 3. The plausibility of parameter values is quantified by ratios of likelihood (Section 2.3), corresponding to a difference on the log scale.
- 4. It is from the log-likelihood (and its derivatives) that most of the theoretical properties of MLEs are obtained (see Part III).

The theoretical properties alluded to in Point 4 are the basis for the two most commonly used forms of likelihood inference – inference based on the likelihood ratio (LR) and inference based on asymptotic normality of the MLE. These two forms of likelihood-based inference are asymptotically equivalent (Section 12.5) in the sense that they lead to the same conclusions for sufficiently large sample sizes. However, in real situations there can be a non-negligible difference between these two approaches (Section 4.3).

Using the likelihood ratio approach in the context of Example 1.1, an interval of plausible values of the unknown parameter p_0 is obtained as all values p for which the log-likelihood is above a certain threshold. In Section 3.4 it is shown that the threshold can be chosen so that the resulting interval has desirable frequentist properties. In the continuation of Example 1.1 below, the threshold is chosen so that the resulting interval for parameter p.

The curvature of the log-likelihood is of fundamental importance in both the theory and practice of likelihood inference. The curvature is quantified by the second derivative, that is, the change in slope. When evaluated at the MLE, the second derivative is negative (because the slope changes from being positive for $p < \hat{p}$ to negative for $p > \hat{p}$) and the larger its absolute value the more sharply curved the loglikelihood is at its maximum. Intuitively, a sharply curved log-likelihood is desirable because this narrows the range over which the log-likelihood is close to its maximum value, that is, it narrows the range of plausible parameter values. In Section 3.2 it is seen that the variance of the MLE can be estimated by the inverse of the negative of the second derivative of the log-likelihood. This is particularly convenient in practice because some optimization algorithms evaluate the second derivative of the objective function as part of the algorithmic calculations (see Section 5.2). In the maximum likelihood context, the objective function is the log-likelihood, and the estimated variance of the MLE is an easily-calculated byproduct from such optimizers. The approximate normality of MLEs enables confidence intervals and hypothesis tests to be performed using well-established techniques.

The likelihood ratio and curvature-based methods of likelihood inference are demonstrated in the following continuation of Example 1.1.

Example 1.1 continued. The log-likelihood function for p, 0 , is

$$l(p) = \log L(p)$$

= $\log \left(\frac{100!}{90! \ 10!}\right) + 10 \log p + 90 \log(1-p)$
= $30.48232 + 10 \log p + 90 \log(1-p)$, (1.2)



Figure 1.2 Binomial log-likelihood for 10 successes from 100 trials, and 95 % likelihood ratio confidence interval.

and the maximized value of this log-likelihood is $l(\hat{p}) = l(0.1) \approx -2.03$.

In Section 3.4 it is seen that an approximate 95 % likelihood ratio confidence interval for parameter p is given by all values p_0 for which $l(p_0)$ is within about 1.92 of the maximized value of the log-likelihood. (The value 1.92 arises as one half of the 0.95 quantile of a chi-square distribution with one degree of freedom.) So, in this case, the interval is given by all values of p_0 for which $l(p_0)$ is -3.95 or higher. The confidence interval can be read from Figure 1.2, or obtained numerically for greater accuracy. This interval is (0.051, 0.169) to the accuracy of three decimal places. From the equivalence between confidence intervals and hypothesis tests (Section 13.2) it can be concluded that the null hypothesis $H_0: p = p_0$ will be rejected at the 5 % level for any value of p_0 outside of the interval (0.051, 0.169).

To perform inference based on the curvature of the log-likelihood, the second derivative of the log-likelihood is required. This second derivative is given in Equation (11.15), and for n = 100 trials and y = 10 successes it is

$$l''(p) = \frac{\partial^2 l(p)}{\partial p^2} = -\frac{10}{p^2} - \frac{90}{(1-p)^2} \,. \tag{1.3}$$

Evaluating this second derivative at the MLE $\hat{p} = 0.1$ gives

$$l''(\hat{p}) = -\frac{10}{0.01} - \frac{90}{0.81} \approx -1111.111$$
.

The inverse of the negative of l''(0.1) is exactly 0.0009, and according to likelihood theory (Sections 3.2 and 12.2), this is the approximate variance of \hat{p} . The approximate standard error is therefore $\sqrt{0.0009} = 0.03$.

Recall that for a binomial experiment, the true variance of \hat{p} is $p_0(1 - p_0)/n$, which is estimated by $\hat{p}(1 - \hat{p})/n$. This estimate of variance is also 0.0009, the same as that obtained from using $-1/l^{"}(0.1)$. (In fact, for the binomial the two variance estimates are always the same, for any values of *n* and *y*.)

8 PRELIMINARIES

For sufficiently large *n*, the distribution of \hat{p} can be approximated by a normal distribution, thereby permitting approximate tests and confidence intervals for *p* to be performed using familiar techniques. These are often called Wald tests or intervals, due to the influential work of Abraham Wald in establishing the large-sample approximate normality of MLEs (e.g. Wald 1943). The $(1 - \alpha)100$ % Wald confidence interval for *p* can be obtained using the familiar formula that calculates the upper (or lower) bounds as the point estimate plus (or minus) $z_{1-\alpha/2}$ times the estimated standard error (\hat{se}), where $z_{1-\alpha/2}$ is the $(1 - \alpha/2)$ quantile of the standard normal distribution. Thus, the approximate 95 % Wald confidence interval is

$$\widehat{p} \pm z_{0.975} \,\widehat{\mathrm{se}}(\widehat{p}),\tag{1.4}$$

where $z_{0.975} \approx 1.96$ and $\hat{se}(\hat{p}) = 0.03$. This interval is (0.041, 0.159). Equivalently, this interval is the collection of the values of p_0 such the null hypothesis H_0 : $p = p_0$ is not rejected at the 5 % level by the Z-statistic. This is the values of p_0 that satisfy the inequality

$$|Z| = \left| \frac{\widehat{p} - p_0}{\widehat{\operatorname{se}}(\widehat{p})} \right| < z_{0.975} .$$
(1.5)

Although the Wald CI and test statistic in (1.4) and (1.5) may be the most commonly taught and used methods of such inference for the binomial model, it is hoped that this text will convince the reader to avoid Wald (i.e. approximate normality) methodology whenever it is practicably feasible. See the next section for more on this.

Box 1.2

1.2.2 Approximate normality versus likelihood ratio

The Wald form of confidence interval used in (1.4) is based on the approximate normal distribution of \hat{p} . This is the most commonly used method for constructing approximate confidence intervals because of its intuitive appeal and computational ease. It was shown earlier that the likelihood ratio can be used as an alternative method for constructing confidence intervals – which should be used?

From a pragmatic point of view, there is considerable intuitive appeal in the Wald construction of a 95% (say) confidence interval, with bounds given by 1.96 standard errors each side of the point estimate. This form of CI will be the most familiar to anyone with a basic grounding in frequentist statistics. However, when the LR and Wald intervals differ substantially, it is generally the case that the LR approach is superior, in the sense that the CIs obtained using likelihood ratio will have actual coverage probability closer to the a priori chosen value of $(1-\alpha)$ (see

Section 4.3.1). In fact, the results of Brown *et al.* (2001) question the popular usage of the Wald CI for binomial inference because of its woeful performance, even for some values of *n* and *p* for which the normal approximation to the binomial distribution is generally considered reasonable (typically, $\min(n\hat{p}, n(1 - \hat{p})) \ge 5$). Unfortunately, the LR confidence interval is not as widely used because it requires (a little) knowledge of likelihood theory, but more importantly because it can not generally be calculated explicitly.

Application of Wald tests and construction of CIs extends to multi-parameter inference, but becomes more cumbersome and unfamiliar when simultaneous inference about two or more parameters is required. It is then that LR-based inference tends to be more commonly used. In particular, multi-parameter inference is typical of model selection problems, and in this area LR-based inference dominates. Also, it should be noted that model selection criterion such as Akaike's Information Criterion (AIC) (Section 4.4.1) make direct use of the likelihood.

In addition to the Wald and LR intervals, there are several other competing methods for constructing approximate confidence intervals for the probability parameter p in a binomial experiment. These include the Wilson score (see Box 3.1, Example 12.10, and Exercise 12.7), Agresti-Coull, and the misnamed 'exact' CIs. The comparisons performed by Agresti and Coull (1998) and Brown *et al.* (2002) suggest that the LR and Wilson score CIs are to be preferred.

Box 1.3

Summary

To conclude, Example 1.1 demonstrates likelihood inference in a nutshell. Much of the rest of this book is devoted to providing pragmatic guidance on the use (and potential abuse) of inferential methods based on likelihood ratios and approximate normality of MLEs, and their application to more complex and realistic models. These concepts extend naturally to models with two or more parameters, although the implementation can become challenging. For example, in a model where the number of parameters is s > 2, the second derivative of the log-likelihood is an *s*-dimensional square matrix (the Hessian) and the negative of its inverse provides an approximate variance matrix for the MLEs.

1.3 Using SAS, R and ADMB

This book is not just about understanding maximum likelihood inference, it is also very much about *doing* it with real data. Examples in SAS and R (Ihaka and Gentleman 1996, R Development Core Team 2010) are provided throughout Part II, along with a smattering of examples demonstrating Automatic Differentiation Model Builder (ADMB, ADMB-project (2008a, or any later version)).

10 PRELIMINARIES

Unlike the SAS and R environments, ADMB is a tool specifically designed for complex optimization problems. Due to the learning curve required to use ADMB, its use is difficult to justify if existing functionality within SAS or R can be used instead. Other than the quick demonstration of ADMB later in this chapter, it is used sparingly until Chapter 10 where it becomes the best choice for the general-purpose fitting of latent variable models. Some of its additional capabilities are noted in Sections 4.2.3 and 5.4.2.

The SAS examples presented in this text were implemented using SAS for Windows version 9.2. The SAS procedures used throughout are found in the statistics module SAS/STAT (SAS Institute 2008), with the exception that occasional use was made of the nonlinear optimizer PROC NLP which is in the operations research module SAS/OR. Some users of SAS/STAT may find that their licence does not extend to SAS/OR and hence will not be able to use PROC NLP. For this reason, PROC NLP is used sparingly and alternative SAS code is given where possible.

SAS procedures typically produce a lot of output by default. The output often includes a lot of superfluous information such as details about the contents of the data-set being used, computational information, and unwanted summary statistics. Throughout, the Output Delivery System (ODS) in the SAS software has been used to select only the required parts of the output produced by the SAS procedure.

Delwiche and Slaughter (2003, or any later edition) provides an excellent introduction to SAS. For ease of readability, the SAS code presented herein follows their typographical convention. This convention is to write SAS keywords in uppercase, and to use lowercase for variable names, data-set names, comments, etc. Note that SAS code is not case sensitive.

The R examples were run using R for Windows version 2.12.0. R is freely available under the terms of the Free Software Foundation's GNU General Public License (see http://www.R-project.org. Most of the R functions used herein are incorporated in the default installation of R. Others are available within specified R library packages, and can be easily loaded from within the R session.

ADMB is freely available via the ADMB project (http://www.admbproject.org), where full instructions for using ADMB can also be found. A short description of automatic differentiation is given in Section 15.6. In brief, ADMB is implemented by programming the objective function within an ADMB template file. The objective function is just the (negative) log-likelihood (and in latent variable models the density function of the latent variables also needs to be specified). An executable file is then created from the template file. Fortunately, much of the detail in creating the executable can be hidden behind convenient user interfaces. The ADMB examples in this book were run from within R using the interface provided by the PBSadmb package.

In many applications of ML inference it will be possible to make use of existing SAS procedures and R functions that are appropriate to the type of data being modelled, notwithstanding that this convenience often comes at the loss of flexibility. Rather than using existing functionality that is specific to the binomial model, the implementations of Example 1.1 presented below demonstrate a selection of the

general-purpose tools available in SAS and R, and the use of ADMB. In particular, calculation of likelihood ratio confidence intervals is an application of profile likelihood (Section 3.6), and the examples below make use of general-purpose code for this purpose.

1.3.1 Software resources

Several small pieces of code have been written to facilitate techniques described in this text. These are listed in Section 15.5, along with a brief description of their functionality. These software resources are freely available for download from http://www.stat.auckland.ac.nz/~millar. This web resource also contains the complete code, and the data, for all examples used in this text.

1.4 Implementation of the motivating example

The code used below demonstrates how an explicit log-likelihood function is maximized within each of SAS, R and ADMB, and the calculation of the Wald and likelihood-ratio confidence intervals. Some efficiencies could have been gained by taking advantage of built-in functionality within the software. For example, in the SAS example, the binomial model could have been expressed using the statement MODEL $\gamma \sim \text{BINOMIAL}(n,p)$, but the general-purpose likelihood specification has been used here for illustration. In R, various functionality (e.g. the mle function in package stat4, or maxLik function in the package of the same name) could have been used to shortcut some of the required code. However, the savings are minimal, and it is instructive to see the individual programming steps.

The first term of the binomial log-likelihood given in Equation (1.2) is a constant, and hence is irrelevant to maximization of the log-likelihood. However, it is good practice always to include the constant terms because it removes possible sources of confusion when fits of different model types are being compared (e.g. using Akaike's information criterion), or when verifying the fit of a model by using an alternative choice of software. Inclusion of the constant terms in the log-likelihood is becoming standard in most software applications of ML, but do not take this for granted.

The description of the code that is presented below is relatively complete, but this level of explanation is too unwieldy to be used throughout the remainder of this text. For more explanation on programming details and syntax, the reader should refer to the abundant online resources and documentation for each of these software.

1.4.1 Binomial example in SAS

The SAS code below uses PROC NLMIXED to implement Example 1.1, and produces the output shown in Figure 1.3.

Parameter estimates									
		Standard							
Parameter	Estimate	Error	DF	tValue	Pr> t	Alpha	Lower	Upper	Gradient
р	0.1000	0.03000	1E6	3.33	0.0009	0.05	0.04120	0.1588	8.566E-7

Figure 1.3 The parameter estimates table from PROC NLMIXED, including the 95 % Wald confidence interval (0.0412,0.1588).

```
DATA binomial;
y=10; n=100;
RUN;
*Select only the parameter estimates table;
ODS SELECT ParameterEstimates;
PROC NLMIXED DF=1E6 DATA=binomial;
PARMS p=0.5;
BOUNDS 0<p<1;
loglhood=LOG(COMB(n,y))+y*log(p)+(n-y)*log(1-p);
MODEL y~GENERAL(loglhood);
RUN;
```

Some features of the above code are:

- The default output includes several tables, including tables of log-likelihood values and fit statistics. The Output Delivery System statement ODS SELECT ParameterEstimates; is used to select only the required table.
- By default, NLMIXED calculates Wald intervals using a *t*-distribution with degrees of freedom equal to the number of observations (rows in the dataset). To get the normal-based Wald interval in (1.4), the value for the degrees of freedom needs to be set to a large number. In this case, it was set to one million using the procedure option DF=1E6.
- The PARMS statement is an optional statement used to explicitly list the parameters and their initial values.
- The BOUNDS statement is an optional statement used to specify the range of the parameter values (i.e. the parameter space).
- The model is specified using the MODEL statement. Here, the model is given as GENERAL (loglhood) to specify that PROC NLMIXED should maximize the value of the log-likelihood, loglhood, as specified by the preceding programming statement.
- In the SAS output in Figure 1.3, Gradient gives the slope of the loglikelihood upon termination of the optimization. It should be near zero. If not, then convergence of the optimizer to a maximum of the log-likelihood may not have been achieved.