

JAMES R. CARPENTER
MICHAEL G. KENWARD

Multiple Imputation and its Application



STATISTICS IN PRACTICE

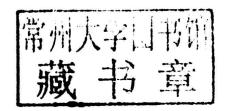
Multiple Imputation and its Application

James R. Carpenter

and

Michael G. Kenward

Department of Medical Statistics London School of Hygiene and Tropical Medicine, UK





This edition first published 2013 © 2013 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

Carpenter, James R.

Multiple imputation and its application / James R. Carpenter, Michael G. Kenward. - 1st ed. p. : cm.

Includes bibliographical references and index.

ISBN 978-0-470-74052-1 (hardback)

I. Kenward, Michael G., 1956- II. Title.

[DNLM: 1. Data Interpretation, Statistical. 2. Biomedical Research-methods. WA 950] 610.72'4-dc23

2012028821

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

ISBN: 978-0-470-74052-1

Cover photograph courtesy of Harvey Goldstein

Set in 10/12pt Times by Laserwords Private Limited, Chennai, India Printed and bound in Singapore by Markono Print Media Pte Ltd

Multiple Imputation and its Application

Statistics in Practice

Series Advisory Editors

Marian Scott University of Glasgow, UK

Stephen Senn *CRP-Santé, Luxembourg*

Wolfgang Jank University of Maryland, USA

Founding Editor

Vic Barnett
Nottingham Trent University, UK

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 clear explanations 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 fields of employment and research environments. Subject areas covered include medicine and pharmaceuticals; industry, finance and commerce; public services, and the earth and environmental sciences.

The books also provide support to students studying applied statistics courses in these areas. The demand for applied statistics graduates in these areas has led to such courses becoming increasingly prevalent at universities and colleges.

It is our aim to present judiciously chosen and well-written textbooks to meet everyday practical needs. Feedback from readers will be valuable in monitoring our success.

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

Preface

No study of any complexity manages to collect all the intended data. Analysis of the resulting partially collected data must therefore address the issues raised by the missing data. Unfortunately, the inferential consequences of missing data are not simply restricted to the proportion of missing observations. Instead, the interplay between the substantive questions and the reasons for the missing data is crucial. Thus, there is no simple, universal, solution.

Suppose, for the substantive question at hand, the inferential consequences of missing data are nontrivial. Then the analyst must make a set of assumptions about the reasons, or mechanisms, causing data to be missing, and perform an inferentially valid analysis under these assumptions. In this regard, analysis of a partially observed dataset is the same as any statistical analysis; the difference is that when data are missing we cannot assess the validity of these assumptions in the way we might do in a regression analysis, for example. Hence, sensitivity analysis, where we explore the robustness of inference to different assumptions about the reasons for missing data, is important.

Given a set of assumptions about the reasons data are missing, there are a number of statistical methods for carrying out the analysis. These include the EM algorithm, inverse probability weighting, a full Bayesian analysis and, depending on the setting, a direct application of maximum likelihood. These methods, and those derived from them, each have their own advantages in particular settings. Nevertheless, we argue that none shares the practical utility, broad applicability and relative simplicity of Rubin's Multiple Imputation (MI).

Following an introductory chapter outlining the issues raised by missing data, the focus of this book is therefore MI. We outline its theoretical basis, and then describe its application to a range of common analysis in the medical and social sciences, reflecting the wide application that MI has seen in recent years. In particular, we describe its application with nonlinear relationships and interactions, with survival data and with multilevel data. The last three chapters consider practical sensitivity analyses, combining MI with inverse probability weighting, and doubly robust MI.

Self-evidently, a key component of an MI analysis, is the construction of an appropriate method of imputation. There is no unique, ideal, way in which this should be done. In particular, there there has been some discussion in the literature about the relative merits of the joint modelling and full conditional

xii PREFACE

specification approaches. We have found that thinking in terms of joint models is both natural and convenient for formulating imputation models, a range of which can then be (approximately) implemented using a full conditional specification approach. Differences in computational speed between joint modelling and full conditional specification are generally due to coding efficiency, rather than intrinsic superiority of one method over the other.

Throughout the book we illustrate the ideas with several examples. The code used for these examples, in various software packages, is available from the book's home page, which is at http://www.wiley.com/go/multiple_imputation, together with exercises to go with each chapter.

We welcome feedback from readers; any comments and corrections should be e-mailed to mi@lshtm.ac.uk. Unfortunately, we cannot promise to respond individually to each message.

Data acknowledgements

We are grateful to the following:

AstraZeneca for permission to use data from the 5-arm asthma study in examples in Chapters 1, 3, 7 and 10;

GlaxoSmithKline for permission to use data from the dental pain study in Chapter 4, and the RECORD study in Chapter 12;

Mike English (Director, Child and Newborn Health Group, Kemri-Wellcome Trust Research Programme, Nairobi, Kenya) for permission to use data from a multifaceted intervention to implement guidelines and improve admission paediatric care in Kenyan district hospitals, in Chapter 9;

Peter Blatchford for permission to use data from the Class Size Study (Blatchford *et al*, 2002) in Chapter 9, and

Sarah Schroter for permission to use data from the study to improve the quality of peer review in Chapter 10.

In Chapters 1, 5, 8, 10 and 11 we have analysed data from the Youth Cohort Time Series for England, Wales and Scotland, 1984-2002 First Edition, Colchester, Essex, published by and freely available from the UK Data Archive, Study Number SN 5765. We thank Vernon Gayle for introducing us to these data.

In Chapter 6 we have analysed data from the Alzheimer's Disease Neuro-imaging Initiative (ADNI) database (adni.loni.ucla.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in the analysis or in the writing of this book. A complete listing of ADNI investigators can be found at: http://adni.loni.ucla.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf.

The ADNI was launched in 2003 by the National Institute on Aging (NIA), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), the Food and Drug Administration (FDA), private pharmaceutical companies and nonprofit organisations, as a \$60 million, five-year public-private partnership. The primary goal of ADNI has been to test whether serial magnetic resonance imaging (MRI), positron emission tomography (PET), other biological markers, and clinical and neuropsychological assessment can be combined to measure the

progression of mild cognitive impairment (MCI) and early Alzheimer's disease (AD). Determination of sensitive and specific markers of very early AD progression is intended to aid researchers and clinicians to develop new treatments and monitor their effectiveness, as well as lessen the time and cost of clinical trials.

The Principal Investigator of this initiative is Michael W. Weiner, MD, VA Medical Center and University of California San Francisco. ADNI is the result of the efforts of many co-investigators from a broad range of academic institutions and private corporations, and subjects have been recruited from over 50 sites across the US and Canada. The initial goal of ADNI was to recruit 800 adults, ages 55 to 90, to participate in the research, with approximately 200 cognitively normal older individuals to be followed for 3 years, 400 people with MCI to be followed for 3 years and 200 people with early AD to be followed for 2 years. For up-to-date information, see www.adni-info.org.

Data collection and sharing for this project was funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: Abbott; Alzheimer's Association; Alzheimer's Drug Discovery Foundation; Amorfix Life Sciences Ltd.; AstraZeneca; Bayer HealthCare; BioClinica, Inc.; Biogen Idec, Inc.; Bristol-Myers Squibb Company; Eisai, Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; GE Healthcare; Innogenetics, N.V.; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC; Johnson & Johnson Pharmaceutical Research & Development, LLC; Medpace, Inc.; Merck & Co., Inc.; Meso Scale Diagnostics, LLC; Novartis Pharmaceuticals Corporation; Pfizer, Inc.; Servier; Synarc, Inc.; and Takeda Pharmaceutical Company. The Canadian Institutes of Health Research is providing funds to support ADNI clinical sites in Canada. Private sector contributions are facilitated by the Foundation for the National Institutes of Health (www.fnih.org). The grantee organisation is the Northern California Institute for Research and Education, and the study is coordinated by the Alzheimer's Disease Cooperative Study at the University of California, San Diego. ADNI data are disseminated by the Laboratory for Neuro-imaging at the University of California, Los Angeles. This research was also supported by NIH grants P30 AG010129 and K01 AG030514.

In Chapter 7 we have analysed data from the 1958 National Childhood Development Study. This is published, and freely available from the UK Data Archive, Study Number SN 5565 (waves 0-3) and SN 5566 (wave 4). We thank Ian Plewis for introducing us to these data.

Acknowledgements

No book of this kind is written in a vacuum, and we are grateful to many friends and colleagues for research collaborations, stimulating discussions and comments on draft chapters.

In particular we would like to thank members of the Missing Data Imputation and Analysis (MiDIA) group, including (in alphabetical order) Jonathan Bartlett, John Carlin, Rhian Daniel, Dan Jackson, Shaun Seaman, Jonathan Sterne, Kate Tilling and Ian White.

We would also like to acknowledge many years of collaboration with Geert Molenberghs, James Roger and Harvey Goldstein.

James would like to thank Mike Elliott, Rod Little, Trivellore Raghunathan and Jeremy Taylor for facilitating a visit to the Institute for Social Research and Department of Biostatistics at the University of Michigan, Ann Arbor, in Summer 2011, when the majority of the first draft was written.

Thanks to Tim Collier for the anecdote in §1.3.

We also gratefully acknowledge funding support from the ESRC (3-year fellowship for James Carpenter, RES-063-27-0257, and follow-on funding RES-189-25-0103) and MRC (grants G0900724, G0900701 and G0600599).

We would also like to thank Richard Davies and Kathryn Sharples at Wiley for their encouragement and support.

Lastly, thanks to our families for their forbearance and understanding over the course of this project.

Despite the encouragement and support of those listed above, the text inevitably contains errors and shortcomings, for which we take full responsibility.

James Carpenter and Mike Kenward London School of Hygiene & Tropical Medicine

Glossary

Indices and symbols

i	indexes units, often individuals, unless defined otherwise
j	indexes variables in the data set, unless defined otherwise
n	total number of units in the data set, unless defined otherwise
p	depending on context, number of variables in a
	data set or number of parameters in a statistical model
X, Y, Z	random variables
$Y_{i,j}$ θ	i^{th} observation on j^{th} variable, $i = 1,, n, j = 1,, p$.
θ	generic parameter
θ	generic parameter column vector, typically $p \times 1$
β, γ, δ	regression coefficients
β	column vector of regression coefficients, typically $p \times 1$.

Matrices

Ω	matrix, typically of dimension $p \times p$.
$\mathbf{\Omega}_{i,j}$ $\mathbf{\Omega}^{T}$	i, j^{th} element of Ω
$\mathbf{\Omega}^T$	transpose of Ω , so that $\Omega_{i,j}^T = \Omega_{i,j}$.
$\mathbf{Y}_j = (Y_{1,j}, \dots, Y_{n,j})^T$	$n \times 1$ column vector of observations on variable j .
$\operatorname{tr}(\mathbf{\Omega})$	sum of diagonal elements of Ω , ie $\sum \Omega_{i,i}$
	known as the trace of the matrix.

Abbreviations

AIPW	Augmented Inverse Probability Weighting
CAR	Censoring At Random
CNAR	Censoring Not At Random
EM	Expectation Maximisation
FCS	Full Conditional Specification
FEV_1	Forced Expiratory Volume in 1 second (measured in litres)
FMI	Fraction of Missing Information
IPW	Inverse Probability Weighting
FCS FEV ₁ FMI	Full Conditional Specification Forced Expiratory Volume in 1 second (measured in litres) Fraction of Missing Information

此为试读,需要完整PDF请访问: www.ertongbook.com

xviii GLOSSARY

MAR Missing At Random

MCAR Missing Completely At Random

MI Multiple Imputation
MNAR Missing Not At Random
POD Partially Observed Data
POM Probability Of Missingness

S.E. Standard error

Probability distributions

f(.) probability distribution function F(.) cumulative distribution function

'|' to be verbalised 'given', as in f(Y|X)

'the probability distribution function of Y given X'

Contents

	Pref	face	X
	Data	a acknowledgements	xiii
	Acknowledgements		
	Glos	ssary	xvi
PA	RT	I FOUNDATIONS	1
1	Intr	roduction	3
	1.1	Reasons for missing data	4
	1.2	Examples	7
	1.3	Patterns of missing data	7
		1.3.1 Consequences of missing data	9
	1.4	Inferential framework and notation	10
		1.4.1 Missing Completely At Random (MCAR)	11
		1.4.2 Missing At Random (MAR)	12 17
		1.4.3 Missing Not At Random (MNAR)	21
	1.5	1.4.4 Ignorability Using observed data to inform assumptions about the	2
	1.5	missingness mechanism	21
	1.6	Implications of missing data mechanisms for regression analyses	24
	1.0	1.6.1 Partially observed response	24
		1.6.2 Missing covariates	28
		1.6.3 Missing covariates and response	30
		1.6.4 Subtle issues I: The odds ratio	30
		1.6.5 Implication for linear regression	32
		1.6.6 Subtle issues II: Subsample ignorability	33
		1.6.7 Summary: When restricting to complete records is valid	34
	1.7	Summary	35
2	The	multiple imputation procedure and its justification	37
	2.1		3
	2.2	Intuitive outline of the MI procedure	38
		=	

vi	C	ONTENTS	
	2.3	The generic MI procedure	44
	2.4	Bayesian justification of MI	46
	2.5	Frequentist inference	48
		2.5.1 Large number of imputations	49
		2.5.2 Small number of imputations	49
	2.6	Choosing the number of imputations	54
	2.7	Some simple examples	55
	2.8	MI in more general settings	62
		2.8.1 Survey sample settings	70
		Constructing congenial imputation models	70
		Practical considerations for choosing imputation models	71
	2.11	Discussion	73
P/	ART	II MULTIPLE IMPUTATION FOR CROSS	
SI	ECTI	ONAL DATA	75
3	Mult	tiple imputation of quantitative data	77
	3.1	Regression imputation with a monotone missingness pattern	77
		3.1.1 MAR mechanisms consistent with a monotone pattern	79
		3.1.2 Justification	81
	3.2	Joint modelling	81
		3.2.1 Fitting the imputation model	82
	3.3	Full conditional specification	85
	2.4	3.3.1 Justification	86
	3.4	Full conditional specification versus joint modelling	87
	3.5	Software for multivariate normal imputation	88
	3.6	Discussion	88
4		tiple imputation of binary and ordinal data	90
	4.1	Sequential imputation with monotone missingness pattern	90
	4.2	Joint modelling with the multivariate normal distribution	92
	4.3	Modelling binary data using latent normal variables	94
		4.3.1 Latent normal model for ordinal data	98
	4.4	General location model	103
	4.5	Full conditional specification	103
	4 :21	4.5.1 Justification	103
	4.6	Issues with over-fitting	104
	4.7	Pros and cons of the various approaches	109
	4.8	Software	110
	4.9	Discussion	111

5 Multiple imputation of unordered categorical data

Multivariate normal imputation for categorical data

5.1 Monotone missing data

5.2

112

112

114

			CONTENTS	vi
	5.3	Maximum indicant model		114
		5.3.1 Continuous and categorical variable		117
		5.3.2 Imputing missing data		119
		5.3.3 More than one categorical variable		120
	5.4	General location model		121
	5.5	FCS with categorical data		122
	5.6	Perfect prediction issues with categorical data		124
	5.7	Software		126
	5.8	Discussion		126
6	Non	llinear relationships		127
	6.1	Passive imputation		128
	6.2	No missing data in nonlinear relationships		130
	6.3	Missing data in nonlinear relationships		133
		6.3.1 Predictive Mean Matching (PMM)		133
		6.3.2 Just Another Variable (JAV)		134
		6.3.3 Joint modelling approach		135
		6.3.4 Extension to more general models and i	missing data	
		patterns		138
		6.3.5 Metropolis-Hastings sampling		140
		6.3.6 Rejection sampling		141
		6.3.7 FCS approach		143
	6.4	Discussion		145
7		eractions		147
	7.1	Interaction variables fully observed		147
	7.2	8		151
	7.3	General nonlinear relationships		155
		Software		163
	7.5	Discussion		164
D/	рт	III ADVANCED TOPICS		1/5
F	INI	III ADVANCED TOPICS		165
8		vival data, skips and large datasets		167
	8.1	Time-to-event data		167
		8.1.1 Imputing missing covariate values		169
		8.1.2 Survival data as categorical		173
		8.1.3 Imputing censored survival times		177
	8.2	Nonparametric, or 'hot deck' imputation		180
		8.2.1 Nonparametric imputation for survival of	lata	182
	8.3	Multiple imputation for skips		184
	8.4	Two-stage MI		188
	8.5	Large datasets		190
		8.5.1 Large datasets and joint modelling		190

viii	C	CONTENTS	
		8.5.2 Shrinkage by constraining parameters	192
		8.5.3 Comparison of the two approaches	195
	8.6	Multiple imputation and record linkage	195
	8.7	Measurement error	197
	8.8	Multiple imputation for aggregated scores	200
		Discussion	202
9	Mult	ilevel multiple imputation	203
	9.1	Multilevel imputation model	203
	9.2	MCMC algorithm for imputation model ,	214
	9.3	Imputing level-2 covariates using FCS	220
	9.4	Individual patient meta-analysis	222
		9.4.1 When to apply Rubin's rules	224
	9.5	Extensions	225
		9.5.1 Random level-1 covariance matrices	226
	0.7	9.5.2 Model fit	228
	9.6	Discussion	228
10	Sens	itivity analysis: MI unleashed	229
	10.1	Review of MNAR modelling	230
	10.2	Framing sensitivity analysis	233
	10.3	Pattern mixture modelling with MI	235
		10.3.1 Missing covariates	240
		10.3.2 Application to survival analysis	241
	10.4	Pattern mixture approach with longitudinal data via MI	246
		10.4.1 Change in slope post-deviation	247
	10.5	Piecing together post-deviation distributions from	240
	10.7	other trial arms	249
	10.6	Approximating a selection model by importance weighting 10.6.1 Algorithm for approximate sensitivity analysis by	257
		re-weighting	259
	10.7	Discussion	268
11		ading survey weights	269
		Using model based predictions	269
	11.2	Bias in the MI variance estimator	271
		11.2.1 MI with weights	274
	110	11.2.2 Estimation in domains	276
		A multilevel approach	277
		Further developments	280
	11.5	Discussion	281
12	Rob	ust multiple imputation	282
		Introduction	282
	12.2	Theoretical background	284

		CONTENTS	1X
12.2.1	Simple estimating equations		284
12.2.2	The Probability Of Missingness (POM) mod	lel	285
12.2.3	Augmented inverse probability weighted esti	imating	
	equation		286
	multiple imputation		287
12.3.1	Univariate MAR missing data		287
12.3.2	Longitudinal MAR missing data		289
12.4 Simulat			292
	Univariate MAR missing data		292
	Longitudinal monotone MAR missing data		293
	Longitudinal nonmonotone MAR missing da		293
	Nonlongitudinal nonmonotone MAR missing	g data	297
	Results and discussion		297
	ECORD study		302
12.6 Discuss	sion		304
Appendix A Ma	arkov Chain Monte Carlo		306
Appendix B Pro	bability distributions		310
B.1 Posterio	or for the multivariate normal distribution		313
Bibliography			316
Index of Author	rs		327
Index of Examp	bles		332
Index			334