

# Synthetic Datasets for Statistical Disclosure Control

Theory and Implementation

# **Lecture Notes in Statistics**

**201**

Edited by P. Bickel, P.J. Diggle, S. Fienberg, U. Gather,  
I. Olkin, S. Zeger



Jörg Drechsler

# Synthetic Datasets for Statistical Disclosure Control

Theory and Implementation

 Springer

Jörg Drechsler  
Department for Statistical Methods  
Institute for Employment Research  
Regensburger Straße 104  
90478 Nürnberg  
Germany  
[Joerg.Drechsler@iab.de](mailto:Joerg.Drechsler@iab.de)

ISSN 0930-0325  
ISBN 978-1-4614-0325-8                      e-ISBN 978-1-4614-0326-5  
DOI 10.1007/978-1-4614-0326-5  
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2011931290

© Springer Science+Business Media, LLC 2011

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are 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 on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

*To my mother and my father (in loving  
memory) for their love and support*



# Foreword

The topic of Jörg Drechsler's work is, in my view, extremely important because it addresses two conflicting demands that are becoming ever more important and complex with the increasing sophistication of our society. First, there is the demand to have access to the vast amounts of publicly supported data collected on all of us. Second, there is the demand to preserve the confidentiality of critical information about individuals in the data being released.

For a specific example of the first demand, in the United States there is the recent call to use the vast collection of medical data, routinely collected on patients from hospitals, pharmacies, etc., to conduct "comparative effectiveness research" in order to find the best combination of medical treatments for individuals. The search for answers to such questions, and therefore the request for publicly available microdata, i.e., data on individuals, is legitimate. Nevertheless, the release of such data threatens the privacy of patients.

The second demand, therefore, is for any released data to preserve confidential information from the individuals whose data are being released, whether because of explicit or implicit guarantees made to them. Even the release of one piece of confidential information can have relatively dire consequences when combined with publicly available information. For another U.S. example, with a person's name and birth date, both of which are available essentially to anyone, all an "intruder" needs is a social security number (taxpayer number) to open credit card accounts, obtain loans, charge hospital bills, open Internet and cell phone accounts, etc. – with all records and debts attached to that social security number. The result is that the holder of that social security number can have a disastrous credit rating that is essentially impossible to correct, even after thousands of dollars in expenses and many years of trying. This "stolen identity" problem is just one example of the untoward effects of the release of confidential information, which may include life-altering consequences, such as being denied mortgages on home purchases.

The work that Jörg Drechsler is pursuing in this book addresses both demands by trying to find ways to benefit society by releasing microdata, here multiply imputed synthetic microdata, that simultaneously preserve individuals' confidential information and yet allow valid inferences at some level of detail through the use of



specialized methods for combining the analyses of the resulting multiply imputed datasets. The topic is a statistically challenging one that needs much development, and I'm sure that this book will be a critical stimulus to this development. Jörg is to be congratulated for this great contribution.

Cambridge, Massachusetts, March 2011

*D. B. Rubin*

# Acknowledgements

This book would never have been possible without the help of many colleagues and friends, and I am very grateful for their wonderful support. First, I want to thank my Ph.D. advisor, Susanne Rässler, for introducing me to the world of multiple imputation and suggesting I join a research project on synthetic data at the Institute for Employment Research (IAB) that finally became the cornerstone of my thesis and eventually of this book. Her remarkable enthusiasm helped me pass some of the local minima of my dissertation function, and without her I would never have met and eventually worked with some of the greatest researchers in the field.

I am very grateful to Trivellore Raghunathan for joining my dissertation committee and providing helpful suggestions for the revision of my thesis for this book. Although I only had two weeks during a visit at the University of Michigan to benefit from his expertise, I learned a lot in that short period of time and I am still deeply impressed by his ability to grasp complex research problems within seconds but even more importantly by his capacity to instantly come up with often simple and straightforward solutions for seemingly (at least for me) unsolvable problems.

I also want to thank John Abowd for inviting me to participate in weekly videoconferences with the leading experts on synthetic data in the United States. When I started my research, I was the only one involved in that topic in Europe and following the discussions and learning from the experience of these experts during these weekly meetings was extremely helpful for my endeavor. To Don Rubin, one of the founding fathers of synthetic data for data confidentiality, I am thankful for his invitation to present my work at Harvard and for fruitful discussions on some of my papers on the topic that later found their way into my thesis. I feel especially honored that he accepted writing the foreword for this book. Bill Winkler deserves my gratitude for providing the extensive list of references on microdata confidentiality included in the appendix of this book. John Kimmel and Marc Strauss at Springer provided great support while I worked on turning my thesis into an acceptable contribution for the Springer Lecture Notes in Statistics Series. Anne-Sophie Charest contributed very helpful comments on an earlier version of this book.

At the IAB I am especially thankful to Hans Kiesl, Thomas Büttner, and Stefan Bender. Hans always helped me out when my lack of background in survey statis-

tics once again became too obvious. Thomas joined me in the dissertation journey. It was a great relief to have a fellow sufferer. And both of them provided helpful discussions on the details of multiple imputation and unforgettable road trips framing JSMS and other conferences around the world. Stefan was very supportive of my research from the very beginning. He stood up for my work when others were still merely laughing at the idea of generating synthetic datasets, even though he was and probably still is skeptical about the idea himself. He helped me find my way in the jungle of official statistics and assisted me in any way he could.

My deepest gratitude is to Jerry Reiter, with whom I had the pleasure to work on several projects that later became part of my thesis. Chapters 6, 7, and especially Chapter 9 in this book borrow heavily from joint papers that were a direct result of these projects. Almost everything I know on the theoretical concepts behind synthetic datasets I owe to Jerry. He has been and continues to be a great mentor and friend.

Finally, I want to thank my mother, Ursula Drechsler, her partner Jochen Paschedag, and the rest of my family for their wonderful support and care. Even though spending three years developing fake data must have seemed bizarre to them, they were always interested in the progress of my work and helped me whenever they could. Most importantly, I would never have survived this trip without the constant love of my fiancée, Veronika. There is no way I can thank her enough for all her patience and understanding for numerous weekends and evenings I spent in front of the computer. She always cheered me up when deadlines were approaching surprisingly fast and the simulations still didn't provide the results they were supposed to show. I thank her for bringing more colors to my life.

Nürnberg, April 2011

*Jörg Drechsler*

# Contents

<b>Foreword</b> .....	vii
<b>Acknowledgements</b> .....	ix
<b>Acronyms</b> .....	xv
<b>List of Figures</b> .....	xvii
<b>List of Tables</b> .....	xix
<b>1 Introduction</b> .....	1
<b>2 Background on Multiply Imputed Synthetic Datasets</b> .....	7
2.1 The history of multiply imputed synthetic datasets .....	7
2.2 Advantages of multiply imputed synthetic datasets compared with other SDC methods .....	10
<b>3 Background on Multiple Imputation</b> .....	13
3.1 Two general approaches to generate multiple imputations .....	14
3.1.1 Joint modeling .....	14
3.1.2 Fully conditional specification (FCS) .....	15
3.1.3 Pros and cons of joint modeling and FCS .....	18
3.2 Real data problems and possible ways to handle them .....	18
3.2.1 Imputation of semi-continuous variables .....	19
3.2.2 Bracketed imputation .....	19
3.2.3 Imputation under linear constraints .....	20
3.2.4 Skip patterns .....	20
<b>4 The IAB Establishment Panel</b> .....	23

<b>5</b>	<b>Multiple Imputation for Nonresponse</b> .....	27
5.1	Inference for datasets multiply imputed to address nonresponse ....	27
5.1.1	Univariate estimands .....	27
5.1.2	Multivariate estimands .....	28
5.2	Analytical validity for datasets multiply imputed to address nonresponse .....	30
5.3	Multiple imputation of the missing values in the IAB Establishment Panel .....	31
5.3.1	The imputation task .....	31
5.3.2	Imputation models .....	32
5.3.3	Evaluating the quality of the imputations .....	33
<b>6</b>	<b>Fully Synthetic Datasets</b> .....	39
6.1	Inference for fully synthetic datasets .....	40
6.1.1	Univariate estimands .....	40
6.1.2	Multivariate estimands .....	40
6.2	Analytical validity for fully synthetic datasets .....	41
6.3	Disclosure risk for fully synthetic datasets .....	42
6.4	Application of the fully synthetic approach to the IAB Establishment Panel .....	44
6.4.1	The imputation procedure .....	46
6.4.2	Measuring the analytical validity .....	47
6.4.3	Assessing the disclosure risk .....	48
<b>7</b>	<b>Partially Synthetic Datasets</b> .....	53
7.1	Inference for partially synthetic datasets .....	53
7.1.1	Univariate estimands .....	54
7.1.2	Multivariate estimands .....	55
7.2	Analytical validity for partially synthetic datasets .....	56
7.3	Disclosure risk for partially synthetic datasets .....	56
7.3.1	Ignoring the uncertainty from sampling .....	57
7.3.2	Accounting for the uncertainty from sampling .....	58
7.4	Application of the partially synthetic approach to the IAB Establishment Panel .....	59
7.4.1	Measuring the analytical validity .....	60
7.4.2	Assessing the disclosure risk .....	61
7.5	Pros and cons of fully and partially synthetic datasets .....	62
<b>8</b>	<b>Multiple Imputation for Nonresponse and Statistical Disclosure Control</b> .....	65
8.1	Inference for partially synthetic datasets when the original data are subject to nonresponse .....	65
8.1.1	Univariate estimands .....	66
8.1.2	Multivariate estimands .....	67
8.2	Analytical validity and disclosure risk .....	68

- 8.3 Generating synthetic datasets from the multiply imputed IAB Establishment Panel ..... 68
  - 8.3.1 Selecting the variables to be synthesized ..... 68
  - 8.3.2 The synthesis task ..... 70
  - 8.3.3 Measuring the analytical validity ..... 71
  - 8.3.4 Caveats in the use of synthetic datasets ..... 76
  - 8.3.5 Assessing the disclosure risk ..... 78
- 9 A Two-Stage Imputation Procedure to Balance the Risk–Utility Trade-Off ..... 87**
  - 9.1 Inference for synthetic datasets generated in two stages ..... 88
    - 9.1.1 Fully synthetic data ..... 88
    - 9.1.2 Partially synthetic data ..... 90
  - 9.2 Analytical validity and disclosure risk ..... 90
  - 9.3 Application of the two-stage approach to the IAB Establishment Panel ..... 90
    - 9.3.1 Analytical validity for the panel from one-stage synthesis .. 91
    - 9.3.2 Disclosure risk for the panel from one-stage synthesis ..... 93
    - 9.3.3 Results for the two-stage imputation approach ..... 96
- 10 Chances and Obstacles for Multiply Imputed Synthetic Datasets .... 99**
- A Bill Winkler’s Microdata Confidentiality References ..... 103**
- B Binned Residual Plots to Evaluate the Imputations for the Categorical Variables ..... 119**
- C Simulation Study for the Variance-inflated Imputation Model ..... 127**
- Bibliography ..... 131**
- References ..... 131**
- Index ..... 137**

