

Generating Multiply Imputed Synthetic Datasets: Theory and Implementation

Dissertation

zur Erlangung des akademischen Grades

eines Doktors der Sozial- und Wirtschaftswissenschaften

(Dr. rer. pol.)

an der Fakultät Sozial- und Wirtschaftswissenschaften
der Otto-Friedrich-Universität Bamberg

vorgelegt von
Jörg Drechsler

Bamberg, im Februar 2010

Datum der Disputation: 10. Dezember 2009

Prüfungskommission:

1. Gutachter: Professor Dr. Susanne Rässler
2. Gutachter: Professor Trivellore Raghunathan, Ph.D.
3. Gutachter (Disputation): Professor Dr. Henriette Engelhardt-Wölfler

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

Acknowledgements

This work would never have been possible without the help of many colleagues and friends and I am very grateful for their wonderful support. Firstly, I want to thank my advisor Susanne Rässler for introducing me to the world of multiple imputation and suggesting to join a research project on synthetic data at the Institute for Employment Research (IAB) that finally became the corner stone of this thesis. Her remarkable enthusiasm helped me to 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. 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 understand complex research problems within seconds even when explained badly by some PhD student, but even more importantly 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 the weekly video conferences with the leading experts on synthetic data in the U.S. When I started my research, I was the only one involved in that topic in Europe and following the discussions and learning from their experience during these weekly meetings was extremely helpful for my endeavor. To Don Rubin, one of the founding fathers of synthetic data, I am thankful for inviting me to present my work at Harvard and for fruitful discussions on some of my papers on the topic. Bill Winkler deserves my gratitude for providing the extensive list of references on microdata confidentiality included in the Appendix of my thesis. At the IAB I am especially thankful to Hans Kiesl, Thomas Büttner, and Stefan Bender. Hans for always helping me out when my lack of background in survey statistics once again became too obvious. Thomas for joining me in the dissertation journey. It was a great relief to have a fellow sufferer. And both of them for helpful discussions on the details of multiple imputation and for unforgettable road trips framing JSMs and other conferences around the world. Stefan was very supportive for 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 sceptical

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 which later became part of my thesis. Almost everything I know on the theoretical concepts behind synthetic datasets I owe to him. He has been and continues to be a great mentor and friend.

Most importantly 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. Finally, I would never have survived this trip without the constant love of my girlfriend 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. Thanks for bringing more colors to my live.

Contents

1	Introduction	1
2	Background on Multiply Imputed Synthetic Datasets	7
2.1	The history of multiply imputed synthetic datasets	7
2.2	Advantages of multiply imputed synthetic datasets compared to other SDC methods	11
3	Multiple Imputation for Nonresponse	15
3.1	The concept of multiple imputation	15
3.2	Two general approaches to generate imputations for missing values	17
3.2.1	Joint modeling	17
3.2.2	Fully conditional specification (FCS)	18
3.2.3	Pros and cons of joint modeling and FCS	20
3.3	Real data problems and possible ways to handle them	21
3.3.1	Imputation of skewed continuous variables	21
3.3.2	Imputation of semi-continuous variables	22
3.3.3	Imputation under non-negativity constraints	22
3.3.4	Imputation under linear constraints	23
3.3.5	Skip patterns	23
4	The IAB Establishment Panel	25
5	Fully Synthetic Datasets	29
5.1	Inference for fully synthetic datasets	30
5.2	Data utility for fully synthetic datasets	31
5.3	Disclosure risk for fully synthetic datasets	32
5.4	Application of the fully synthetic approach to the IAB Estab- lishment Panel	34

5.4.1	The imputation procedure	36
5.4.2	Measuring the data utility	37
5.4.3	Assessing the disclosure risk	40
6	Partially Synthetic Datasets	47
6.1	Inference for partially synthetic datasets	48
6.2	Data utility for partially synthetic datasets	49
6.3	Disclosure risk for partially synthetic datasets	49
6.3.1	Ignoring the uncertainty from sampling	50
6.3.2	Accounting for the uncertainty from sampling	52
6.4	Application of the partially synthetic approach to the IAB Es- tablishment Panel	53
6.4.1	Measuring the data utility	54
6.4.2	Assessing the disclosure risk	55
6.5	Pros and cons of fully and partially synthetic datasets	57
7	Multiple Imputation for Nonresponse and Statistical Disclo- sure Control	59
7.1	Inference for partially synthetic datasets when the original data is subject to nonresponse	60
7.2	Data utility and disclosure risk	61
7.3	Multiple imputation of the missing values in the IAB Establish- ment Panel	61
7.3.1	The imputation task	62
7.3.2	Imputation models	62
7.3.3	Evaluating the quality of the imputations	63
7.4	Generating synthetic datasets from the multiply imputed IAB Establishment Panel	71
7.4.1	The synthesis task	72
7.4.2	Measuring the data utility	75
7.4.3	Caveats in the use of synthetic datasets	81
7.4.4	Assessing the disclosure risk	84
7.4.4.1	Log-linear modeling to estimate the number of matches in the population	86
7.4.4.2	Results from the disclosure risk evaluations	87
7.4.4.3	Disclosure risk for large establishments	88

7.4.4.4	Additional protection for the largest establishments in the survey	91
8	A Two Stage Imputation Procedure to Balance the Risk-Utility-Trade-Off	93
8.1	Inference for synthetic datasets generated in two stages	94
8.1.1	Fully synthetic data	94
8.1.2	Partially synthetic data	96
8.2	Data utility and disclosure risk	97
8.3	Application of the two stage approach to the IAB Establishment Panel	98
8.3.1	Data utility for the panel from one stage synthesis	98
8.3.2	Disclosure risk for the panel from one stage synthesis	99
8.3.3	Results for the two stage imputation approach	104
9	Chances and Obstacles for Multiply Imputed Synthetic Datasets	107
	Appendix	113
A.1	Bill Winkler's Microdata Confidentiality References	113
A.2	Binned residual plots to evaluate the imputations for the categorical variables	132
A.3	Simulation study for the variance inflated imputation model	139

List of Figures

5.1	The fully synthetic approach for the IAB Establishment Panel.	36
5.2	Included variables from the IAB Establishment Panel and the German Social Security Data.	38
5.3	Occurrence of establishments already included in the original survey by establishment size.	43
5.4	Distribution of the matching rates for different multiple response questions.	44
5.5	Histogram of the relative difference between original and imputed values for the variable <i>establishment size</i>	44
7.1	Observed and imputed data for <i>payroll</i> and <i>number of participants in further education</i>	66
7.2	Model checks for <i>turnover</i> and <i>number of participants in further education with college degree</i>	68
7.3	Ordered probit regression of <i>expected employment trend</i> on 39 explanatory variables and industry dummies.	78
7.4	Original point estimates against synthetic point estimates for the overall mean and the means in subgroups defined by establishment size class, industry code and region.	79
7.5	Box plots of CI overlaps for all continuous variables for the overall mean and the means in all subgroups defined by different stratifying variables.	80
7.6	QQ-plots for the <i>number of employees covered by social security</i> 2006 and 2007 and the employment trend between the two years.	83
7.7	Plots of F_t against \hat{F}_t for all establishments and for establishments with more than 100 employees.	88

A.1	Binned residual plots for the categorical variables with missing rates above 1%.	132
A.2	Binned residual plots for the categorical variables with missing rates above 1%.	133
A.3	Binned residual plots for the categorical variables with missing rates above 1%.	134
A.4	Binned residual plots for the categorical variables with missing rates above 1%.	135
A.5	Binned residual plots for the categorical variables with missing rates above 1%.	136
A.6	Binned residual plots for the categorical variables with missing rates above 1%.	137
A.7	Binned residual plots for the categorical variables with missing rates above 1%.	138

List of Tables

5.1	Results from the vocational training regression for one stage full synthesis.	39
5.2	How many records are sampled how often in the new samples? .	41
5.3	Establishments from the IAB Establishment Panel that also occur in at least one of the new samples.	42
6.1	Results from the vocational training regression for one stage partial synthesis.	55
7.1	Missing rates and means per quantile for <i>NB.PRE</i>	67
7.2	Expectations for the investments in 2007.	70
7.3	Regression results from a probit regression of <i>part time-employees (yes/no)</i> on 19 explanatory variables in West Germany.	76
7.4	Regression results from a probit regression of <i>part time-employees (yes/no)</i> on 19 explanatory variables in East Germany.	77
7.5	Regression results from a probit regression of <i>employment trend (increase/no increase)</i> on 19 explanatory variables in West Germany.	82
7.6	Probabilities to be included in the target sample and in the original sample depending on establishment size.	85
7.7	Average F_t and \hat{F}_t for different establishment size classes.	87
7.8	Disclosure risk summaries for the synthetic establishment panel wave 2007.	88
7.9	False match rate and true match risk for different levels of γ	89
7.10	Mode of the establishment size rank and average match rate for large establishments.	90
8.1	Average number of employees by industry for one stage synthesis.	99

8.2	Results from the vocational training regression for one stage partial synthesis revisited.	100
8.3	Confidence interval overlap for the average number of employees for one stage synthesis.	101
8.4	Confidence interval overlap for the vocational training regression for one stage synthesis.	102
8.5	Average confidence interval overlap for all 31 estimands for ten independent simulations of one stage synthesis.	103
8.6	Averages of the disclosure risk measures over ten simulations of one stage synthesis.	104
8.7	Average CI overlap and match risk for two stage synthesis based on ten simulations.	105
A.1	Simulation results for the variance inflated imputation model.	140
A.2	Simulation results if Y_1 is excluded from the imputation model.	142

Chapter 1

Introduction

National Statistical Institutes (NSIs) like the U.S. Census Bureau or the German Federal Statistical Office gather valuable information on many different aspects of the society. Broad access to this information is desirable to stimulate research in official statistics. However, most data obtained by the institutes are collected under the pledge of privacy and thus the natural interest of enabling as much research as possible with the collected data has to stand back behind the confidentiality guaranteed to the survey respondent. But not only legal aspects are relevant when considering disseminating data to the public. Respondents that feel their privacy is at risk might be less willing to provide sensitive information, might give incorrect answers or might even refuse to participate completely – with devastating consequences for the quality of the data collected (Lane, 2007). Traditionally, this meant that access to the data was strictly limited to researchers working for the NSI. With the increasing demand for access to the data on the micro-level from external researchers, accelerated by the improvements in computer technology, agencies started looking for possibilities to disseminate data that provide a high level of data quality while still guaranteeing confidentiality for the participating units.

Over the years a broad literature on statistical disclosure limitation (SDL) techniques for microdata evolved (see Bill Winkler’s famous list of microdata confidentiality references in the Appendix A.1). These techniques can be divided into two main categories: Approaches that protect the data by reducing the amount of information contained in the released file through coarsening of the data and approaches classified as data perturbation methods that try to maintain most of the originally collected information but protect the data

by changing some of the values on the micro level. Information reducing approaches protect the data by

- *categorizing continuous variables*: Building categories from the underlying continuous variables and reporting only in which category the unit falls, for example building age groups in five year intervals.
- *top coding*: Setting values above a certain threshold equal to the threshold, for example reporting the income for all individuals with income above 100,000 as "100,000+"
- *coarsening categorical variables*: Coarsening to a reduced number of categories, for example instead of providing information on the state level, only reporting whether a respondent lives in West or East Germany.
- *dropping variables*: Dropping some variables that are considered too sensitive (e.g. HIV-status) or are not enough protected by any of the above methods.

There is a vast literature on data perturbation methods and discussing all approaches including possible modifications is beyond the scope of this introduction. A detailed overview is given in the handbook on statistical disclosure control (Center of Excellence for Statistical Disclosure Control, 2009) issued by members of the CENEX-SDC project funded by Eurostat. A good reference for recent developments are the proceedings from the biannual conference *Privacy in Statistical Databases* (Springer LNCS 3050, 4302, 5262).

While the first methods developed in the eighties like swapping and adding noise mainly focused on disclosure protection and preserved only some univariate statistics like the population mean and the variance of a single variable, more sophisticated methods emerged in recent years. But these sophisticated methods often require different complicated adjustments for each estimate to get unbiased results, preserve only certain statistics like the vector of the means or the variance-covariance matrix, or are valid only under specific distributional assumptions like multivariate normality that are unrealistic for real datasets. Besides, most statistical agencies still only apply standard methods mainly because of their easy of implementation. Winkler (2007b) shows the devastating consequences on data quality for many of these easy to implement procedures while others fail to achieve their primary goal: protecting the data adequately.

Since many of the proposed data perturbation methods significantly reduce data quality and it is often impossible for the researcher using the perturbed data to judge, if the results are still at least approximately valid, there is a common mistrust among researchers against these methods. Still, strict legal requirements in many countries often force agencies to perturb their data before release, even though they know that data quality can be heavily affected. The situation is a little different in Germany where the required disclosure protection for datasets only used for scientific purposes, so called *scientific use files* is lower than for datasets that are available to anybody (*public use files*). For scientific use files, the German Federal Law on Statistics enables the release of de facto anonymous microdata. "Factual anonymity means that the data can be allocated to the respondent or party concerned only by employing an excessive amount of time, expenses and manpower" (Knoche, 1993). The concept of factual anonymity takes into account a rational thinking intruder, who calculates the costs and benefits of the re-identification of the data. Because factual anonymity depends on several conditions and is not further defined by law, it is necessary to estimate the costs and benefits of a re-identification for every dataset with a realistic scenario. Disseminating scientific use files under this law is much easier than under the usual requirement that a re-identification of a single unit should be impossible under any circumstance. For this reason the scientific use files available in Germany traditionally are protected using only a mixture of the non perturbative methods described above. Nevertheless, there is a common agreement that the dissemination of microdata on businesses is not possible using only non perturbative methods, since the risk of disclosure is much higher for these data than it is for microdata on individuals for several reasons:

- The underlying population is much smaller for businesses than it is for individuals.
- Variables like turnover or establishment size have very skewed distributions that make the identification of single units in the dataset very easy.
- There is a lot of information about businesses in the public domain already. This information can be used to identify records in the released dataset.
- The benefit from identifying a unit in an establishment survey might be

higher for a potential attacker than the benefit of identifying a unit in a household survey.

- In most business surveys the probability of inclusion is very high for large businesses (often close to 1) so there is no additional privacy protection from sampling for these units.

Since only few variables like turnover, region, and industry code are necessary to identify many businesses, no data on enterprises were disseminated for many years. In 2002 a joint project of the German Federal Statistical Office, several Statistical Offices of the Länder and the Institute for Applied Economic Research started investigating the possibilities of generating scientific use files for these data applying data perturbative methods for the first time in Germany. They came to the result that using these methods a release is possible and disseminated several survey datasets protected by either adding multiplicative noise or microaggregation (Statistisches Bundesamt, 2005). With the long history of releasing only unperturbed data, it is not surprising that acceptance of these datasets was rather limited in the following years. Many users of these data tend to believe the collected data is the direct truth and ignore all the additional uncertainty and possible bias introduced on the collection stage by measurement errors, coding mistakes, bad sampling design and especially steadily increasing nonresponse rates that make the implicit assumption of a missingness pattern that is missing completely at random (Rubin, 1987) of complete case analysis more and more questionable. The additional bias introduced by the perturbation method might be dwarfed by the bias already inherent in the data due to these facts. But also the selected perturbation methods might be a reason for the limited acceptance. Winkler (2007b) illustrates the negative consequences of univariate microaggregation, namely on correlations and although correction factors for estimations based on data perturbed by multiplicative noise are illustrated in the German *Handbuch zur Anonymisierung wirtschaftsstatistischer Mikrodaten* (Statistisches Bundesamt, 2005) for the linear model and the SIMEX Method (Lechner and Pohlmeier, 2005) can be used for nonlinear models, both are difficult to compute and are applicable only under some additional assumptions. The *Handbuch* shows that the SIMEX method produces biased results for a probit regression using simulated data. A further disadvantage the two methods share with most data perturbative methods is that logical constraints between variables are not

preserved.

This illustrates the common dilemma for data disseminating agencies: Fulfilling only one goal – no risk of disclosure or high data quality – is straightforward; release data generated completely at random or release the original unchanged data. In both cases at least one party will be unhappy about the results, but balancing the two goals is extremely difficult. A dataset that guarantees the confidentiality of the respondent but is not accepted by the research community due to data quality concerns is of little value and the question arises, if the high costs in time and money to produce these datasets are justified.

A new approach to address the trade-off between data utility and disclosure risk overcoming the problems discussed above was proposed by Rubin (1993): The release of multiply imputed synthetic datasets. Specifically, he proposed that agencies (i) randomly and independently sample units from the sampling frame to comprise each synthetic dataset, (ii) impute unknown data values for units in the synthetic samples using models fit with the original survey data, and (iii) release multiple versions of these datasets to the public. These are called fully synthetic datasets.

However, the quality of this method strongly depends on the accuracy of the model used to impute the "missing" values. If the model doesn't include all the relationships between the variables that are of interest to the analyst or if the joint distribution of the variables is misspecified, results from the synthetic datasets can be biased. Furthermore, specifying a model that considers all the skip patterns and constraints between the variables in a large dataset can be cumbersome if not impossible. To overcome these problems, a related approach suggested by Little (1993) replaces observed values with imputed values only for variables that bear a high risk of disclosure or for variables that contain especially sensitive information, leaving the rest of the data unchanged. This approach, discussed as generating partially synthetic datasets in the literature, has been adopted for some datasets in the US (Abowd and Woodcock, 2001, 2004; Kennickell, 1997; Abowd *et al.*, 2006).

The aim of this book is to give the reader a detailed introduction to the different approaches to generating multiply imputed synthetic datasets (MISD) by combining the theory with illustrative examples using a real dataset, the German IAB Establishment Panel. We start by giving an overview of the history on synthetic datasets and discussing the major advantages of this approach compared to other perturbation methods. Since the method is based on

the ideas of multiple imputation (Rubin, 1978), the next chapter recapitulates its basic concepts originally proposed to impute values missing due to nonresponse. Advantages and disadvantages of the two major imputation strategies (joint modeling and fully conditional specification (FCS)) are also addressed. The Chapters 5-8 on different synthetic data generation approaches are all organized in the same manner. First, the general ideas of the specific approach are discussed, then the point and variance estimates that provide valid inferences in this context are presented. Each section concludes with an extensive application to a real dataset. Since all applications are based on the German IAB Establishment Panel, this dataset is introduced in a separate chapter at the beginning of the main part of the book (Chapter 4). The discussed data generation approaches include generating fully synthetic datasets (Chapter 5), generating partially synthetic datasets (Chapter 6), and generating synthetic datasets when the original data is subject to nonresponse (Chapter 7).

Chapter 8 contains an extension to the standard synthetic data generation to better address the trade-off between data utility and disclosure risk: Imputation in two stages, where variables that drive the disclosure risk are imputed less often than others. Since in general data quality and disclosure risk both increase with the number of imputations, defining a different number of imputations for different variables can lead to datasets that maintain the desired data quality with reduced risk of disclosure. In this chapter, the new combining procedures that are necessary for the point and variance estimate are presented for fully and partially synthetic datasets and the IAB Establishment Panel is used to illustrate the impact of the number of imputations on the data quality and the disclosure risk and to show the possible advantage of using a two stage imputation approach. The book concludes with a glimpse into the future of synthetic datasets, discussing the potentials and possible obstacles of the approach and ways to address the concerns of data users and their understandable discomfort with using data that doesn't consist only of the originally collected values.

Chapter 2

Background on Multiply Imputed Synthetic Datasets

2.1 The history of multiply imputed synthetic datasets

In 1993 the Journal of Official Statistics published a special issue on data confidentiality. Two articles in this volume lay the fundament for the development of multiply imputed synthetic datasets (MISD). In his discussion *Statistical Disclosure Limitation* Rubin suggested for the first time to generate synthetic datasets based on his ideas of multiple imputation for missing values (Rubin, 1987). He proposed to treat all the observations from the sampling frame that are not part of the sample as missing data and to impute them according to the multiple imputation framework. Afterwards, several simple random samples from these fully imputed datasets should be released to the public. Because the released dataset does not contain any real data, disclosure of sensitive information is very difficult. On the other hand, if the imputation models are selected carefully and the predictive power of the models is high, most of the information contained in the original data will be preserved. This approach is now called generating fully synthetic datasets in the literature.

In the same issue Little suggested a closely related approach that is also based on the idea of replacing sensitive information by multiple imputation. The major difference is that only part of the data is replaced. These could be either some sensitive variables like income or turnover or key variables like age, place of birth, and sex that could be jointly used to identify a single unit

in the dataset. With this approach, now called generating partially synthetic datasets, it is not mandatory to replace all units for one variable. The replacement can be tailored only to the records at risk. It might be sufficient for example to replace the income only for units with a yearly income above 100,000 EUR to protect the data. This method guarantees that only those records that need to be protected are altered. Leaving unchanged values in the dataset will generally lead to higher data quality, but releasing unchanged values obviously poses a higher risk of disclosure.

In 1994 Fienberg suggested generating synthetic datasets by bootstrapping from a "smoothed" estimate of the empirical cumulative density function of the survey data. This approach was further developed for categorical data in Fienberg *et al.* (1998). 10 years after the initial proposal the complete theory for deriving valid inferences from multiply imputed synthetic datasets was presented for the first time. Raghunathan *et al.* (2003) illustrated, why the standard combining procedures for multiple imputation (Rubin, 1987) are not valid in this context and developed the correct procedures for fully synthetic datasets. The procedures for partially synthetic datasets were presented by Reiter (2003). One year earlier Liu and Little suggested the selective multiple imputation of key variables (SMIKe), replacing a set of sensitive and nonsensitive cases by multiple draws from their posterior predictive distribution under a general location model.

Reiter also demonstrated the validity of the fully synthetic combining procedures under different sampling scenarios (Reiter, 2002), derived the combining procedures when using multiple imputation for missing data and for disclosure avoidance simultaneously (Reiter, 2004), developed significance tests for multi-component estimands in the synthetic data context (Reiter, 2005c), provided an empirical example for fully synthetic datasets (Reiter, 2005b) and presented a non parametric imputation method based on CART models to generate synthetic data (Reiter, 2005d). Recent work includes suggestions for the adjustment of survey weights (Mitra and Reiter, 2006), selecting the number of imputations when using multiple imputation for missing data and disclosure control (Reiter, 2008b), measuring the risk of identity disclosure for partially synthetic datasets (Reiter and Mitra, 2009; Drechsler and Reiter, 2008), and a two stage imputation strategy to better address the trade off between data utility and disclosure risk (Reiter and Drechsler, 2010). A new imputation strategy based on kernel density estimation for variables with very

skewed or even multi-modal distributions has been suggested by Woodcock and Benedetto (2009), while Winkler (2007a) proposed the use of different EM-Algorithms to generate synthetic data subject to convex constraints. The attractive features of synthetic datasets are further discussed by Fienberg and Makov (1998); Abowd and Lane (2004); Little *et al.* (2004); An and Little (2007) and Domingo-Ferrer *et al.* (2009).

It took several years before the ground braking ideas proposed in 1993 were ever applied to any real dataset. The U.S. Federal Reserve Board was the first agency to protect data in its Survey of Consumer Finances by replacing monetary values at high risk of disclosure with multiple imputations, releasing a mixture of these imputed values and the unreplaced, collected values (Kennickell, 1997). Abowd and Woodcock (2001) illustrated the possibilities of protecting longitudinal, linked datasets with data from the French National Institute of Statistics and Economic Studies (INSEE). A very successful implementation of a partially synthetic dataset is the data used behind *On the Map*, illustrating commuting patterns, i.e. where people live and work, for the entire U.S. via maps available to the public on the web (<http://lehdmap.did.census.gov/>). Since the point of origin (where people live) is already in the public domain, only the destination points are synthesized. Machanavajjhala *et al.* (2008) developed a sophisticated synthesizer that maximizes the level of data protection based on the ideas of differential privacy (Dwork, 2006) while still guaranteeing a very high level of data utility. The most ambitious synthetic data project up to date is the generation of a public use file for the Survey of Income and Programm Participation (SIPP) funded by the U.S. Census Bureau and the Social Security Administration (SSA). The variables from the SIPP are combined with selected variables from the International Revenue Service's (IRS) lifetime earnings data, and the SSA's individual benefit data. Almost all of the approximately 625 variables contained in this longitudinal, linked dataset were synthesized. In 2007, four years after the start of the project a beta version of the file was released to the public (www.sipp.census.gov/sipp/synthdata.html). Abowd *et al.* (2006) summarize the steps involved in creating this public use file and provide a detailed disclosure risk and data utility evaluation that indicates that confidentiality is guaranteed while data utility is high for many estimates of interest.

The Census Bureau also protects the identities of people in group quarters (e.g., prisons, shelters) in the public use files of the American Communities Survey by

replacing demographic data for people at high disclosure risk with imputations. Partially synthetic, public use datasets are in the development stage in the U.S. for the Longitudinal Business Database, the Longitudinal Employer-Household Dynamics survey, and the American Communities Survey veterans and full sample data. Recently a statement by the American Statistical Association on data access and personal privacy explicitly mentioned distributing synthetic datasets as an appropriate method of disclosure control (<http://www.amstat.org/news/statementondataaccess.cfm>).

Outside the U.S. the ideas for generating multiply imputed synthetic dataset have been ignored for many years except for some small simulation studies at ISTAT in Italy (Poletti, 2003; Franconi and Stander, 2002, 2003; Poletti *et al.*, 2002). They suggest generating model based synthetic datasets. The main difference to the methods described in this book is that they do not propose multiple imputation and therefore do not correct for the additional variance from imputation. In 2006 the German Institute for Employment Research launched a research project to generate synthetic datasets of its longitudinal establishment survey for release as a scientific use file. In the first phase of the project the fully and partially synthetic approach were tested on a subset of the data (Drechsler *et al.*, 2008b,a). Drechsler *et al.* (2008a) also discuss the advantages and disadvantages of the two approaches in terms of data utility and disclosure risk. Since the evaluations during the first stage of the project indicated that the dataset could be sufficiently protected by the partial synthetic approach, the second stage of the project focused on the generation of a partially synthetic dataset for the complete last wave of the survey. The release of this dataset, the first outside the U.S., is planned for spring 2010. The growing interest in synthetic datasets in Europe is also documented by the report on synthetic data files requested by Eurostat 2008 and published by Domingo-Ferrer *et al.* (2009). Outside Europe statistical agencies in Australia, Canada, and New Zealand (Graham and Penny, 2005; Graham *et al.*, 2009) also are investigating the approach.

2.2 Advantages of multiply imputed synthetic datasets compared to other SDC methods

Generally the aim of this approach is to preserve the joint distribution of the data. Most data perturbation methods either preserve only univariate statistics or only some predefined multivariate statistics like the mean and the variance-covariance matrix in previously defined subgroups. However, most of these methods for statistical disclosure control (SDC) are used to generate datasets for public release on the microdata level and it is impossible to anticipate all analyses potential users will perform with the data. For example one analyst might remove some outliers before running her regressions and it is completely unclear what the effects of SDC methods that only preserve statistics in predefined subsets of the data will be for this reduced dataset. Besides, for some analyses it might be desirable to preserve more than just the first two moments of the distribution, e.g., maintain interaction and nonlinear effects.

Furthermore, many SDC methods are only applicable either to categorical variables or to continuous variables. This means that often a combination of different techniques is required to fully protect a dataset before release. Methods based on multiple imputation on the other hand can be applied to categorical and continuous variables likewise rendering the use of different methods that might require different adjustments by the data analyst unnecessary.

For fully synthetic datasets the actual disclosure risk is further reduced, since the synthetic data is generated for new samples from the population and the intruder never knows, if a unit in the released data was actually included in the original data. Partially synthetic datasets on the other hand have the advantage that the synthesis can be tailored specifically to the records at risk. For some datasets it might only be necessary to synthesize certain subsets of the dataset. Obviously, the decision which records will remain unchanged is a delicate task and a careful disclosure risk evaluation is necessary in this context.

On the other hand, as with any perturbation method, limited data utility is a problem of synthetic data. Only the statistical properties explicitly captured by the model used by the data protector are preserved. A logical question at this point is why not directly publish the statistics one wants to preserve rather than release a synthetic micro dataset. Possible defenses against this

argument are:

- Synthetic data are normally generated by using more information on the original data than is specified in the model whose preservation is guaranteed by the data protector releasing the synthetic data.
- As a consequence of the above, synthetic data may offer utility beyond the models they explicitly preserve.
- It is impossible to anticipate all possible statistics an analyst might be interested in. So access to the micro dataset should be granted.
- Not all users of a public use file will have a sound background in statistics. Some of the users might only be interested in some descriptive statistics and won't be able to generate the results if only the parameters are provided.
- The imputation models in most applications can be very complex, because different models are fitted for every variable and often for different subsets of the dataset. This might lead to hundreds of parameters just for one variable. Thus, it is much more convenient even for the skilled user of the data to have the synthesized dataset available.
- The most important reason for not releasing the parameters is that the parameters themselves could be disclosive in some occasions. For that reason, only some general statements about the generation of the public use file should be released. For example, these general statements could provide information, which variables were included in the imputation model, but not the exact parameters. So the user can judge if her analysis would be covered by the imputation model, but she will not be able to use the parameters to disclose any confidential information.

But the most important advantage is that imputation based synthetic data can tackle many real data problems, other SDC methods cannot handle:

First, most of the data collected by agencies are subject to nonresponse and besides the fact that missing data can lead to biased estimates if not treated correctly by the analyst, many SDC methods can not be applied to SDC methods containing missing values. Since generating multiply imputed synthetic datasets is based on the ideas of multiple imputation for handling item

nonresponse in surveys, it is straight forward to impute missing values before generating synthetic datasets. Reiter (2004) developed methods for simultaneous use of multiple imputation for missing data and disclosure limitation.

Second, model based imputation procedures offer more flexibility if certain constraints need to be preserved in the data. For example non-negativity constraints and linear constraints like *total number of employees* \geq *number of part time employees* can be directly incorporated on the model building stage. Almost all SDC methods fail to preserve linear constraints unless the exact same perturbation is applied to all variables for one unit, which in turn significantly increases the risk of disclosure.

Third, skip patterns, e.g. a battery of questions are only asked if they are applicable, are very common in surveys. Especially, if the skip patterns are hierarchical, it is very difficult to guarantee that perturbed values are consistent with these patterns. With the fully conditional specification approach (see also Section 3.2.2) that sequentially imputes one variable at a time by defining conditional distributions to draw from, it is possible to generate synthetic datasets that are consistent with all these rules.

Lastly, as Reiter (2008a) points out, the MI approach can be relatively transparent to the public analyst. Meta-data about the imputation models can be released and the analyst can judge based on this information if the analysis he or she seeks to perform will give valid results with the synthetic data. For other SDC approaches it is very difficult to decide, how much a particular analysis has been distorted.

Chapter 3

Multiple Imputation for Nonresponse¹

For many datasets, especially for non mandatory surveys, missing data are a common problem. Deleting units that are not fully observed, using only the remaining units is a popular, easy to implement approach in this case. This can possibly lead to severe bias if the strong assumption of a missing pattern that is completely at random (MCAR) is not fulfilled (see for example Rubin (1987)). Imputing the missing values can overcome this problem. However, ad hoc methods like, e.g., mean imputation can destroy the correlation between the variables. Furthermore, imputing missing values only once (single imputation) generally doesn't account for the fact that the imputed values are only estimates for the true values. After the imputation process, they are often treated like truly observed values leading to an underestimation of the variance in the data and by this to p -values that are too significant. Multiple imputation was suggested by Rubin (1978) to overcome these problems.

3.1 The concept of multiple imputation

Multiple imputation, introduced by Rubin (1978) and discussed in detail in Rubin (1987; 2004), is an approach that retains the advantages of imputation while allowing the uncertainty due to imputation to be directly assessed. With multiple imputation, the missing values in a dataset are replaced by $m > 1$ simulated versions, generated according to a probability distribution for the

¹Most of this chapter is taken from Drechsler and Rässler (2008) and Drechsler (2009).

true values given the observed data. More precisely, let Y_{obs} be the observed and Y_{mis} the missing part of a dataset Y , with $Y = (Y_{mis}, Y_{obs})$, then missing values are drawn from the Bayesian posterior predictive distribution of $(Y_{mis}|Y_{obs})$, or an approximation thereof. Typically, m is small, such as $m = 5$. Each of the imputed (and thus completed) datasets is first analyzed by standard methods designed for complete data; the results of the m analyses are then combined to produce estimates, confidence intervals, and test statistics that reflect the missing-data uncertainty properly. In this chapter, we discuss analysis with scalar parameters only, for multidimensional quantities see Little and Rubin (2002), Section 10.2.

To understand the procedure of analyzing multiply imputed datasets, think of an analyst interested in an unknown scalar parameter Q , where Q could be, e.g. the population mean or a regression coefficient in a linear regression.

Inferences for this parameter for datasets with no missing values usually are based on a point estimate q , a variance estimate u , and a normal or Student's t reference distribution. For analysis of the imputed datasets, let q_i and u_i for $i = 1, 2, \dots, m$ be the point and variance estimates achieved from each of the m completed datasets. To get a final estimate over all imputations, these estimates have to be combined using the combining rules first described by Rubin (1978).

For the point estimate, the final estimate simply is the average of the m point estimates $\bar{q}_m = \frac{1}{m} \sum_{i=1}^m q_i$. Its variance is estimated by $T = \bar{u}_m + (1 + m^{-1})b_m$, where $\bar{u}_m = \frac{1}{m} \sum_{i=1}^m u_i$ is the "within-imputation" variance, $b_m = \frac{1}{m-1} \sum_{i=1}^m (q_i - \bar{q}_m)^2$ is the "between-imputation" variance, and the factor $(1 + m^{-1})$ reflects the fact that only a finite number of completed-data estimates q_i are averaged together to obtain the final point estimate. The quantity $\hat{\gamma} = (1 + m^{-1})b_m/T$ estimates the fraction of information about Q that is missing due to nonresponse.

Inferences from multiply imputed data are based on \bar{q}_m , T , and a Student's t reference distribution. Thus, for example, interval estimates for Q have the form $\bar{q}_m \pm t(1 - \alpha/2)\sqrt{T}$, where $t(1 - \alpha/2)$ is the $(1 - \alpha/2)$ quantile of the t distribution. Rubin and Schenker (1986) provide the approximate value $\nu_{RS} = (m - 1)\hat{\gamma}^{-2}$ for the degrees of freedom of the t distribution, under the assumption that with complete data, a normal reference distribution would have been appropriate. Barnard and Rubin (1999) relax the assumption of Rubin and Schenker (1986) to allow for a t reference distribution with complete

data, and suggest the value $\nu_{BR} = (\nu_{RS}^{-1} + \hat{\nu}_{obs}^{-1})^{-1}$ for the degrees of freedom in the multiple-imputation analysis, where $\hat{\nu}_{obs} = (1 - \hat{\gamma})(\nu_{com})(\nu_{com} + 1)/(\nu_{com} + 3)$ and ν_{com} denotes the complete-data degrees of freedom.

3.2 Two general approaches to generate imputations for missing values

Over the years, two different methods emerged to generate draws from $P(Y_{mis}|Y_{obs})$: joint modeling and fully conditional specification (FCS), often also referred to as sequential regression multivariate imputation (SRMI) or chained equations. The first assumes that the data follow a specific distribution, e.g. a multivariate normal distribution. Under this assumption a parametric multivariate density $P(Y|\theta)$ can be specified with θ representing parameters from the assumed underlying distribution. Within the Bayesian framework, this distribution can be used to generate draws from $(Y_{mis}|Y_{obs})$. Methods to create multivariate imputations using this approach have been described in detail by Schafer (1997a), e.g., for the multivariate normal, the log-linear, and the general location model.

FCS on the other hand does not require an explicit assumption for the joint distribution of the dataset. Instead, conditional distributions $P(Y_j|Y_{-j}, \theta_j)$ are specified for each variable separately. Thus imputations are based on univariate distributions allowing for different models for each variable. Missing values in Y_j can be imputed for example by a linear or a logistic regression of Y_j on Y_{-j} , depending on the scales of measurement of Y_j , where Y_{-j} denotes all columns of Y excluding Y_j . The process of iteratively drawing from the conditional distributions can be viewed as a Gibbs sampler that will converge to draws from the theoretical joint distribution of the data if this joint distribution exists.

3.2.1 Joint modeling

In general, it will not be possible to specify $P(Y_{mis}|Y_{obs})$ directly. Note however, that we can write

$$P(Y_{mis}|Y_{obs}) = \int P(Y_{mis}, \theta|Y_{obs})d\theta = \int P(Y_{mis}|Y_{obs}, \theta)P(\theta|Y_{obs})d\theta \quad (3.1)$$

Given this equation, imputations can be generated in two steps:

1. Generate random draws for the parameter θ from its observed-data posterior distribution $P(\theta|Y_{obs})$ given the observed values.
2. Generate random draws for Y_{mis} from its conditional predictive distribution $P(Y_{mis}|Y_{obs}, \theta)$ given the actual parameter θ from step 1.

With joint modeling the second step usually is straight forward. The distribution of $(Y_{mis}|Y_{obs}, \theta)$ can be obtained from the underlying model. For example a multivariate normal density can be assumed for the complete data. But the first step usually requires Markov Chain Monte Carlo techniques, since the observed-data posterior distribution for $(\theta|Y_{obs})$ seldom follows standard distributions, especially if the missing pattern is not monotone. Therefore, often simple random draws from the complete-data posterior $f(\theta|Y_{obs}, Y_{mis})$ are performed. This means that even for joint modeling convergence of the Markov Chain has to be monitored and it is not guaranteed that it will ever converge. Though the probability of non-convergence might be much lower in this context than with FCS, it is still possible and Schafer (1997a) provides examples where the necessary stationary distribution can never be obtained.

3.2.2 Fully conditional specification (FCS)

With FCS the problem of drawing from a k -variate distribution is replaced by drawing k times from much easier to derive univariate distributions. Every variable in the dataset is treated separately using a regression model suitable for that specific variable. Thus, continuous variables can be imputed using a normal model, binary variables can be imputed with a logit model and so on. Here, we can specify $P(\theta|Y_{obs})$ directly and no iterations are necessary, because we don't have to draw from possibly awkward multivariate distributions. For example, if we want to impute a continuous variable Y , we can assume $Y|X \sim N(\mu, \sigma^2)$, where X denotes all variables that are used as explanatory variables for the imputation. The two step imputation approach described above can now be applied as follows:

Let n be the number of observations in the observed part of Y . Let k be the number of regressors to be included in the regression. Let $\hat{\sigma}^2$ and $\hat{\beta}$ be the variance and the beta-coefficient estimates obtained from ordinary least square regressions using only the observed data. Finally, let X_{obs} be the

matrix of regressors for the observed part of Y and X_{mis} be the matrix of regressors for the fraction of the data where Y is missing. Imputed values for Y_{mis} can now be generated using the following algorithm:

Step 1: Draw new values for $\theta = (\sigma^2, \beta)$ from $P(\theta|Y_{obs})$, i.e.,

- draw $\sigma^2|X \sim (Y_{obs} - X_{obs}\hat{\beta})'(Y_{obs} - X_{obs}\hat{\beta})\chi_{n-k}^{-2}$,
- draw $\beta|\sigma^2, X \sim N(\hat{\beta}, (X'_{obs}X_{obs})^{-1}\sigma^2)$.

Step 2: Draw new values for Y_{mis} from $P(Y_{mis}|Y_{obs}, \theta)$, i.e.,

- draw $Y_{mis}|\beta, \sigma^2, X \sim N(X_{mis}\beta, \sigma^2)$.

Note that we are drawing new values for the parameters directly from the observed-data posterior distributions. This means, we don't need Markov Chain Monte Carlo techniques to obtain new values from the complete-data posterior distribution of the parameters. However, there are more variables with missing data. Thus, we generate new values for Y_{mis} by drawing from $P(Y_{mis}|\beta, \sigma^2, X)$ and the matrix of regressors X might contain imputed values from an earlier imputation step. These values have to be updated now, based on the new information in our recently imputed variable Y . Hence, we have to sample iteratively from the fully conditional distribution for every variable in the dataset. This iterative procedure essentially can be seen as a Gibbs sampler for which the iterative draws will converge to draws from the joint distribution, if the joint distribution exists.

In a more detailed notation, for multivariate Y , let $Y_j|Y_{-j}$ be the distribution of Y_j conditioned on all rows of Y except Y_j and θ_j be the parameter specifying the distribution of $Y_j|Y_{-j}$. If Y consists of k rows, and each Y_j is univariate, then the t th iteration of the method consists of the following successive draws:

$$\begin{aligned} \theta_1^{(t)} &\sim P(\theta_1|Y_1^{obs}, Y_2^{(t-1)}, \dots, Y_k^{(t-1)}) \\ Y_1^{(t)} &\sim P(Y_1^{mis}|Y_2^{(t-1)}, \dots, Y_k^{(t-1)}, \theta_1^{(t)}) \\ &\vdots \\ \theta_k^{(t)} &\sim P(\theta_k|Y_k^{obs}, Y_1^{(t)}, Y_2^{(t)}, \dots, Y_{k-1}^{(t)}) \\ Y_k^{(t)} &\sim P(Y_k^{mis}|Y_1^{(t)}, \dots, Y_{k-1}^{(t)}, \theta_k^{(t)}) \end{aligned}$$

Since imputations are generated sequentially variable by variable, this approach is also called sequential regression multivariate imputation (SRMI,

Raghunathan *et al.* (2001)). The sampler will converge to the desired joint distribution of $(Y_{mis}|Y_{obs})$, but only if this joint distribution really exists. In practice it is often impossible to verify this, thus its existence is implicitly assumed. This is problematic, since it will always be possible to draw from the conditional distributions and we will not get any hint that the Gibbs sampler actually never converges.

3.2.3 Pros and cons of joint modeling and FCS

In general, imputing missing values by joint modeling is faster and the imputation algorithms are simpler to implement. Furthermore, if the underlying joint distribution can be specified correctly, joint modeling will guarantee valid results with the imputed dataset. However, empirical data will seldom follow a standard multivariate distribution, especially if they consist of a mix of numerical and categorical variables. Besides, FCS provides a flexible tool to account for bounds, interactions, skip patterns or constraints between different variables (see Section 3.3). It will be very difficult to handle these restrictions that are very common in survey data by joint modeling. In practice the imputation task is often centralized at the methodological department of the statistical agency and imputation experts will fill in missing values for all the surveys conducted by the agency. Imputed datasets that don't fulfill simple restrictions like non-negativity or other logical constraints will never be accepted by subject matter analysts from other departments. Thus, preserving these constraints is a central element of the imputation task.

Overall, joint modeling will be preferable, if only a limited number of variables need to be imputed, no restrictions have to be maintained and the joint distribution can be approximated reasonably well with a standard multivariate distribution. For more complex imputation tasks only fully conditional specification will enable the imputer to preserve constraints inherent in the data. In this case, convergence of the Gibbs sampler should be carefully monitored. A simple way to detect problems with the iterative imputation procedure, is to store the mean of every imputed variable for every iteration of the Gibbs sampler. A plot of the imputed means over the iterations can indicate if there is only the expected random variation between the iterations or if there is a trend between the iterations indicating problems with the model. Of course no observable trend over the iterations is only a necessary and not a sufficient

condition for convergence, since the monitored estimates can stay stable for hundreds of iterations before drifting off to infinity. Nevertheless, this is a straightforward method to identify flawed imputation models. More complex methods to monitor convergence are discussed in Arnold *et al.* (1999).

3.3 Real data problems and possible ways to handle them

The basic concept of multiple imputation is straightforward to apply and multiple imputation software like IVEware in SAS (Raghunathan *et al.*, 2002), `mice` (Van Buuren and Oudshoorn, 2000) and `mi` (Su *et al.*, 2009) in *R*, `ice` in Stata (Royston, 2005) (for FCS), and the stand alone packages NORM, CAT, MIX, and PAN (Schafer, 1997b)(for joint modeling) further reduce the modeling burden for the imputer. However, simply applying standard imputation procedures to real data can lead to biased or inconsistent imputations. Several additional aspects have to be considered in practice, when imputing real data. Unfortunately most of the standard software with the positive exceptions of IVEware and the new `mi` package in *R* can only handle some of these aspects:

3.3.1 Imputation of skewed continuous variables

One problem that especially arises when modeling business data is that most of the continuous variables like *turnover* or *number of employees* are heavily skewed. To control for this skewness, we suggest to transform each continuous variable by taking the cubic root before the imputation. We prefer the cubic root transformation over the log transformation that is often used in the economic literature to model skewed variables like turnover, because the cubic root transformation is less sensitive to deviations between the imputed and the original values in the right tail of the distribution. Since the slope of the exponential function increases exponentially whereas the slope of $f(x) = x^3$ increases only quadratically, a small deviation in the right tail of the imputed transformed variable has more severe consequences after backtransformation for the log transformed variable than for the variable transformed by taking the cubic root.

3.3.2 Imputation of semi-continuous variables

Another problem with modeling continuous variables that often arises in surveys, is the fact that many of these variables in fact are semi-continuous, i.e. they have a spike at one point of the distribution, but the remaining distribution can be seen as a continuous variable. For most variables, this spike will occur at zero. To give an example, in our dataset the establishments are asked how many of their employees obtained a college degree. Most of the small establishments do not require such high skilled workers. In this case, we suggest to adopt the two step imputation approach proposed by Raghunathan *et al.* (2001): In the first step we impute whether the missing value is zero or not. For that, missing values are imputed using a logit model with outcome 1 for all units with a positive value for that variable. In the second step a standard linear model is applied only to the units with observed positive values to predict the actual value for the units with a predicted positive outcome in step one. All values for units with outcome zero in step one are set to zero.

3.3.3 Imputation under non-negativity constraints

Many survey variables can never be negative in reality. This has to be considered during the imputation process. A simple way to achieve this goal is to redraw from the imputation model for those units with imputed values that are negative until all values fulfill the non-negativity constraint. In practice, usually an upper bound z has to be defined for the number of redraws for one unit, since it is possible that the probability to draw a positive value for this unit from the defined model is very low. The value for this unit is set to zero, if z draws from the model never produced a positive value. However, there is a caveat with this approach. Redrawing from the model for negative values is equivalent to drawing from a truncated distribution. If the truncation point is not at the very far end of the distribution, i.e. the model is misspecified, even simple descriptive analyses like the mean of the imputed variable will significantly differ from the true value of the complete data. For this reason, this approach can only be applied, if the probability to draw negative values from the specified model is very low and we only want to prevent that some very unlikely unrealistic values are imputed. If the fraction of units that would have to be corrected with this approach is too high, the model needs to be revised. Usually it is helpful to define different models for different subgroups of

the data. To overcome the problem of generating too many negative values, a separate model for the units with small values should be defined.

3.3.4 Imputation under linear constraints

In many surveys the outcome of one variable by definition has to be equal to or above the outcome of another variable. For example, the total number of employees always has to be at least as high as the number of part-time employees. When imputing missing values in this situation, Schenker *et al.* (2006) suggest the following approach: Variables that define a subgroup of another variable are always expressed as a proportion, i.e. all values for the subgroup variable are divided by the total before the imputation and thus are bounded between zero and one. A logit transformation of the variables guarantees that the variables will have values in the full range $] -\infty, \infty[$ again. Missing values for these transformed variables can be imputed with a standard imputation approach based on linear regressions. After the imputation all values are transformed back to get proportions again and finally all values are multiplied with the totals to get back the absolute values. To avoid problems on the bounds of the proportions, we suggest setting proportions greater than 0.999999 to 0.999999 before the logit transformation and to use the two step imputation approach described in Section 3.3.2 to determine zero values.

3.3.5 Skip patterns

Skip patterns, e.g. a battery of questions are only asked if they are applicable, are very common in surveys. Although it is obvious that they are necessary and can significantly reduce the response burden for the survey participant, they are a nightmare for anybody involved in data editing and imputation or statistical disclosure control. Especially, if the skip patterns are hierarchical, it is very difficult to guarantee that imputed values are consistent with these patterns. With fully conditional specification, it is straightforward to generate imputed datasets that are consistent with all these rules. The two step approach described in Section 3.3.2 can be applied to decide if the questions under consideration are applicable. Values are imputed only for the units selected in step one. Nevertheless, correctly implementing all filtering rules is a labor intensive task that can be more cumbersome than defining good imputation models. Furthermore, the filtering can lead to variables that are answered

by only a small fraction of the respondents and it can be difficult to develop good models based on a small number of observations.

Chapter 4

The IAB Establishment Panel

Since the establishment survey of the German Institute for Employment Research (IAB) is used throughout this book to illustrate the different aspects of multiply imputed synthetic datasets, a short introduction to this dataset should prelude the body of this book.

The IAB Establishment Panel¹ is based on the German employment register aggregated via the establishment number as of 30 June of each year. The basis of the register, the German Social Security Data (GSSD) is the integrated notification procedure for the health, pension and unemployment insurances, which was introduced in January 1973. This procedure requires employers to notify the social security agencies about all employees covered by social security. As by definition the German Social Security Data only include employees covered by social security - civil servants and unpaid family workers for example are not included - approx. 80% of the German workforce are represented. However, the degree of coverage varies considerably across the occupations and the industries.

Since the register only contains information on employees covered by social security, the panel includes establishments with at least one employee covered by social security. The sample is drawn using a stratified sampling design. The stratification cells are defined by ten classes for the size of the establishment, 16 classes for the region², and 17 classes for the industry.³ These cells are also

¹The approach and structure of the establishment panel are described for example by Fischer *et al.* (2008) and Kölling (2000).

²Before 2006 the stratification by region contained 17 classes since two separate classes were used for East and West Germany.

³Between 2000 and 2004 20 industry classes were used, before 2000 the sample was

used for weighting and extrapolation of the sample. The survey is conducted by interviewers from TNS Infratest Sozialforschung. For the first wave, 4,265 establishments were interviewed in West Germany in the third quarter of 1993. Since then the Establishment Panel has been conducted annually - since 1996 with over 4,700 establishments in East Germany in addition. In the wave 2008 more than 16,000 establishments participated in the survey. The response rate of units that have been interviewed repeatedly is over 80%. Each year, the panel is accompanied by supplementary samples and follow-up samples to include new or reviving establishments and to compensate for panel mortality. The questionnaire contains a set of core questions that are asked annually with detailed information about employment development, business policy, vocational training, personnel structure and personnel movements, investments, wages & salaries and adherence to collective agreements. Information on further training, working time, public funding and innovations is asked every other year. Changing additional questions relevant for the current political debate complete the survey.

Considered one of most important business surveys in Germany, there is high demand for access to these data from external researchers. Because of the sensitive nature of the data, researchers desiring direct access to the data have to work on site at the IAB. Alternatively, researchers can submit code for statistical analyses to the IAB research data center, whose staff run the code on the data and send the results to the researchers. To help researchers develop code, the IAB provides access to a publicly available "dummy dataset" with the same structure as the Establishment Panel. The dummy dataset comprises random numbers that only mirror the variable type and the range of the variable without attempts to preserve the joint distributional properties of the variables in the original data. The consequence is that analysis code developed using the dummy dataset often will not run on the original data and it can happen that the code has to be send back to the researcher for revisions several times. For all analyses done with the genuine data, researchers can publicize their analyses only after IAB staff check for potential violations of confidentiality.

Releasing scientific use files of the Establishment Panel would allow more researchers to access the data with fewer burdens, stimulating research on German business data. It also would free up staff time from running code and

stratified by 16 industry classes.

conducting confidentiality checks. Because there are so many sensitive variables in the dataset, standard disclosure limitation methods like swapping or microaggregation would have to be applied with high intensity, which would severely compromise the utility of the released data. Therefore, the IAB decided to develop synthetic data. The first release of a synthetic dataset from the wave 2007 of the panel is planned for spring 2010.

Chapter 5

Fully synthetic datasets¹

In 1993, Rubin suggested to create fully synthetic datasets based on the multiple imputation framework. His idea was to treat all units in the population that have not been selected in the sample as missing data, impute them according to the multiple imputation approach and draw simple random samples from these imputed populations for release to the public. Most surveys are conducted using complex sampling designs. Releasing simple random samples simplifies research for the potential user of the data, since the design doesn't have to be incorporated in the model. It is not necessary however to release simple random samples. If a complex design is used, the analyst accounts for the design in the within variance u_i , $i = 1, \dots, m$.

For illustration, think of a dataset of size n , sampled from a population of size N . Suppose further, the imputer has information about some variables X for the whole population, for example from census records, and only the information from the survey respondents for the remaining variables Y . Let Y_{inc} be the observed part of the population and Y_{exc} the nonsampled units of Y . For simplicity, assume that there are no item-missing data in the observed dataset. The approach also applies if there are missing data. The synthetic datasets can be generated in two steps: First, construct m imputed synthetic populations by drawing Y_{exc} m times independently from the posterior predictive distribution $f(Y_{exc}|X, Y_{inc})$ for the $N - n$ unobserved values of Y . If the released data should contain no real data for Y , all N values can be drawn from this distribution. Second, take simple random samples from these populations

¹Most of this chapter is taken from Drechsler *et al.* (2008b) and Drechsler and Reiter (2009).

and release them to the public. The second step is necessary as it might not be feasible to release m whole populations for the simple matter of data-size. In practice, it is not mandatory to generate complete populations. The imputer can make random draws from X in a first step and only impute values of Y for the drawn X . The analysis of the m simulated datasets follows the same lines as the analysis after multiple imputation (MI) for missing values in regular datasets as described in Section 3.1.

5.1 Inference for fully synthetic datasets

To understand the procedure of analyzing fully synthetic datasets, think of an analyst interested in an unknown scalar parameter Q , where Q could be, e.g., the mean of a variable, the correlation coefficient between two variables, or a regression coefficient in a linear regression. Inferences for this parameter for datasets with no missing values usually are based on a point estimate q , an estimate for the variance of q , u and a normal or Student's t reference distribution. For analysis of the imputed datasets, let q_i and u_i for $i = 1, \dots, m$ be the point and variance estimates for each of the m completed datasets. The following quantities are needed for inferences for scalar Q :

$$\bar{q}_m = \sum_{i=1}^m q_i/m \quad (5.1)$$

$$b_m = \sum_{i=1}^m (q_i - \bar{q}_m)^2 / (m - 1) \quad (5.2)$$

$$\bar{u}_m = \sum_{i=1}^m u_i/m . \quad (5.3)$$

The analyst then can use \bar{q}_m to estimate Q and

$$T_f = (1 + m^{-1})b_m - \bar{u}_m \quad (5.4)$$

to estimate the variance of \bar{q}_m . The difference in the variance estimate compared to the variance estimate for standard multiple imputation (see Section 3.1) is due to the additional sampling from the synthetic units for fully synthetic datasets. Hence, the variance b_m between the datasets already reflects the variance within each imputation. When n is large, inferences for scalar Q can be based on t-distributions with degrees of freedom $\nu_f = (m - 1)(1 - r_m^{-1})^2$,

where $r_m = ((1 + m^{-1})b_m/\bar{u}_m)$. Derivations of these methods are presented in Raghunathan *et al.* (2003). Extensions for multivariate Q are presented in Reiter (2005c).

A disadvantage of this variance estimate is that it can become negative. For that reason, Reiter (2002) suggests a slightly modified variance estimator that is always positive: $T_f^* = \max(0, T_f) + \delta(\frac{n_{syn}}{n}\bar{u}_m)$, where $\delta = 1$ if $T_f < 0$, and $\delta = 0$ otherwise. Here, n_{syn} is the number of observations in the released datasets sampled from the synthetic population.

5.2 Data utility for fully synthetic datasets

It is important to quantify the analytic usefulness of the synthetic datasets. Existing utility measures are of two types: (i) comparisons of broad differences between the original and released data, and (ii) comparisons of differences in specific models between the original and released data. Broad difference measures essentially quantify some statistical distance between the distributions of the original and released data, for example a Kullback-Leibler or Hellinger distance. As the distance between the distributions grows, the overall quality of the released data generally drops.

A very useful measure for specific estimands is the interval overlap measure of Karr *et al.* (2006). For any estimand, we first compute the 95% confidence intervals for the estimand from the synthetic data, (L_s, U_s) , and from the collected data, (L_o, U_o) . Then, we compute the intersection of these two intervals, (L_i, U_i) . The utility measure is

$$I = \frac{U_i - L_i}{2(U_o - L_o)} + \frac{U_i - L_i}{2(U_s - L_s)}. \quad (5.5)$$

When the intervals are nearly identical, corresponding to high utility, $I \approx 1$. When the intervals do not overlap, corresponding to low utility, $I = 0$. The second term in (5.5) is included to differentiate between intervals with $(U_i - L_i)/(U_o - L_o) = 1$ but different lengths. For example, for two synthetic data intervals that fully contain the collected data interval, the measure I favors the shorter interval. The synthesis is successful if we obtain large values of I for many estimands. To compute one number summaries of utility, we can average the values of I over all estimands. This utility measure provides more information than a simple comparison of the two point estimates from the

different datasets because it also considers the standard error of the estimate. Estimates with large standard errors might still have a high confidence interval overlap and by this a high data utility even if their point estimates differ considerably from each other, because the confidence intervals will increase with the standard error of the estimate. For more details on this method see Karr *et al.* (2006).

There do not exist published broad utility measures that account for all m synthetic datasets. The U.S. Census Bureau has adapted an approach described by Woo *et al.* (2009) which is based on how well one can discriminate between the original and disclosure protected data. In this approach, the agency stacks the original and synthetic datasets in one file and estimates probabilities of being “assigned” to the original data conditional on all variables in the dataset. When the probabilities are close to 0.5 for all records in the original and synthetic data, the distributions of the variables are similar—this fact comes from the literature on propensity scores (Rosenbaum and Rubin, 1983)—so that the synthetic data have high utility. This approach is especially useful as a diagnostic for deficiencies in the synthesis methods (variables with significant coefficients in the logistic regression have different distributions in the original and synthetic data).

5.3 Disclosure risk for fully synthetic datasets

In general, the disclosure risk for fully synthetic datasets is very low, since all values are synthetic values. Still, it is not necessarily zero: For example in most establishment surveys the probability of inclusion depends on the size of the establishment and sometimes can be close to 1 for the largest establishments. Since the released synthetic samples will have to be stratified, too to take advantage of the efficiency gained by stratification, the additional protection offered in the fully synthetic approach by drawing new samples from the sampling frame can be very modest for larger establishments. A possible intruder can be confident that large establishments in the released synthetic data represent establishments that were also included in the original survey. The same argument holds for the release of synthetic census data.

Besides this actual risk of disclosure the perceived risk of disclosure also needs to be considered. The released data might look like the data from a potential survey respondent an intruder was looking for. And once the intruder thinks,

he identified a single respondent and the estimates are reasonable close to the true values for that unit, it is no longer important that the data are all made up. The potential respondent will feel that his privacy is at risk. Nevertheless the disclosure risk in general will be very low since the imputation models would have to be almost perfect and the intruder faces the problem that he never knows (i) if the imputed values are anywhere near the true values and (ii) if the target record is included in one of the different synthetic samples.

For this reason the theory on disclosure risk for fully synthetic datasets is far less developed than the theory for partially synthetic datasets (see Section 6.3). Only recently Abowd and Vilhuber (2008) proposed some measures based on the ideas of differential privacy from the computer science literature. To understand the concept of differential privacy, we need some further definitions. Let D_{rel} be the released dataset. Let N be the hypothetical population –unknown to the intruder– from which D_{rel} was supposedly generated. According to Dwork (2006) ϵ -differential privacy is fulfilled if

$$\max \left| \ln \left(\frac{Pr(D_{rel}|N^1)}{Pr(D_{rel}|N^2)} \right) \right| \leq \epsilon \quad (5.6)$$

where ϵ is a predefined threshold and the maximum is taken over all N^1, N^2 that only differ in a single row. The basic idea is that if the ratio is too large, the intruder gains too much information from the released data, since it is far more likely that D_{rel} was generated from N^1 and not from N^2 . The data releasing agency can decide which level of ϵ it is willing to accept. Abowd and Vilhuber (2008) show that this definition of disclosure risk is closely related to the risk of inferential disclosure from the SDC literature that measures the risk by the information gain about a single respondent from the released data compared to the a priori information before the release. The paper also illustrates that synthesizing categorical variables under a Multinomial-Dirichlet model can fulfill the requirements of ϵ -differential privacy.

The definition of ϵ -differential privacy is very appealing since it can be defined ex ante – the agency only needs to select an SDC method that can guarantee ϵ -differential privacy – and the agency can also select the level of privacy guaranteed by defining ϵ . Still, the measure is based on the very strong assumption that the intruder knows all records in the dataset except one and measures how much information the intruder can reveal about this one record. To keep this information low, strong requirements for the SDC method are necessary, namely that the transition matrix between the observed and the

released data doesn't contain any zeros, i.e. any point in the outcome space of a variable must be reachable with positive probability from any given observed value through the transition function between the original and the disclosure protected data implicitly specified by the SDC method. For many datasets this would mean that some very unlikely or even unrealistic events must be reachable with positive probability. Thus, the gain in data protection can come at a very high price in terms of data quality. For this reason Machanavajjhala *et al.* (2008) defined (ϵ, δ) -probabilistic differential privacy, where $1 - \delta$ is the probability that (5.6) holds. This measure has been developed for the Multinomial-Dirichlet model. Further research is necessary to investigate the applicability of this approach to other synthesis models.

5.4 Application of the fully synthetic approach to the IAB Establishment Panel

To generate fully synthetic datasets for the IAB Establishment Panel, we need information from the sampling frame of the Establishment Panel. We obtain this information by aggregating the German Social Security Data (GSSD) to the establishment level. From this aggregated dataset, we can sample new records that provide the basis for the generation of the synthetic datasets. As noted earlier, the German Social Security Data contains information on all employees covered by social security. The notifications of the GSSD include for every employee, among other things, the workplace and the establishment identification number. By aggregating records with the same establishment identification number it is possible to generate establishment information from the GSSD. As we use the 1997 wave of the IAB Establishment Panel for our analysis, data are taken and aggregated from the GSSD for June, 30th 1997 (see Figure 5.2 for all characteristics used). We use the establishment identification number again to match the aggregated establishment characteristics from the GSSD with the IAB Establishment Panel.

In this simulation, we only impute values for a set of variables from the 1997 wave of the IAB Establishment Panel. As it is not feasible to impute values for the millions of establishments contained in the German Social Security Data for 1997, we sample from this frame, using the same sampling design as for the IAB Establishment Panel: Stratification by establishment size, region and

industry. Every stratum contains the same number of units as the observed data from the 1997 wave of the Establishment Panel.

Due to panel mortality a supplementary sample has to be drawn for the IAB Establishment Panel every year. In the 1997 wave, this supplementary sample primarily consisted of newly founded establishments because in that year the questionnaire had a focus on establishment births. Therefore, start-ups are overrepresented in the sample. Arguably, answers from these establishments differ systematically from the answers provided by establishments existing for several years. Drawing a new sample without taking this oversampling into account could lead to a sample after imputation that differs substantially from that in the Establishment Panel.

For simplicity reasons, we define establishments not included in the German Social Security Data before July 1995 as new establishments and delete them from the sampling frame and the Establishment Panel. For the 1997 wave of the Establishment Panel, this means a reduction from 8,850 to 7,610 observations.

Merging the GSSD and the IAB Establishment Panel using the establishment identification number reveals that 278 units from the panel are not included in the GSSD.² These units are also omitted leading to a final sample of 7,332 observations. Furthermore, we have to verify that the stratum parameters size, industry and region match in both datasets. Merging indicates that there are some differences between the two records. If the datasets differ, values from the GSSD are adopted.

Cross tabulation of the stratum parameters for the 7,332 observations in our sample provides a matrix containing the number of observations for each stratum. Now, a new dataset can be generated easily by drawing establishments from the German Social Security Data according to this matrix.

After matching, every dataset is structured as follows: Let N be the total number of units in the newly generated dataset, that is the number of units in the new sample n_s plus the number of units in the panel n_p , $N = n_s + n_p$. Let X be the matrix of variables with information for all observations in N . Then X consists of the variables *establishment size* (from the GSSD), *region* and *industry* and the other variables added from the German Social Security

²There are several possible reasons for this, e.g. re-organization of the firm leading to new establishment identification numbers, coding errors, or delays in the notifications for an establishment in the GSSD.

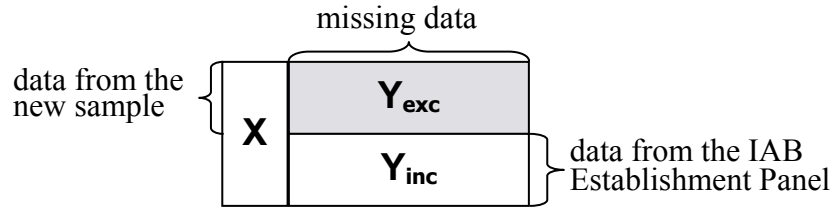


Figure 5.1: The fully synthetic approach for the IAB Establishment Panel.

Data. Note that the variable *establishment size* is included in both, the GSSD and the establishment panel. These two variables need not necessarily be identical, since they are reported at different points in time. However, we use the establishment size from the GSSD as a very strong predictor when synthesizing the establishment size in the establishment panel. Let Y be the selected variables from the Establishment Panel, with $Y = (Y_{inc}, Y_{exc})$, where Y_{inc} are the observed values from the Establishment Panel and Y_{exc} are the hypothetic missing data for the newly drawn values in X (see Figure 5.1).

Now, values for the missing data can be imputed as outlined in Chapter 3 by drawing Y_{exc} from the posterior predictive distribution $f(Y_{exc}|X, Y_{inc})$ for the $N - n_p$ unobserved values of Y . After the imputation procedure, all observations from the GSSD and all originally observed values from the establishment panel are omitted and only the imputed values for the panel are released. Results from an analysis on these released data can be compared with the results achieved with the real data.

5.4.1 The imputation procedure

For this simulation, we only generate 10 synthetic datasets. Previous research has shown that releasing large numbers of fully synthetic datasets improves synthetic data inferences (Reiter, 2005b). The usual advice from multiple imputation for missing data - release five multiply-imputed datasets - tends not to work well for fully synthetic data because the fractions of "missing" information are large. Drechsler *et al.* (2008b) obtain higher analytic validity by generating 100 fully synthetic datasets using the two stage imputation approach described in Chapter 8.

To generate the synthetic datasets we use the SRMI approach (see Section 3.2.2) as implemented in the software IVEware (Raghunathan *et al.*, 2002).

Since most of the continuous variables like *establishment size* are heavily skewed, these variables are transformed by taking the cubic root before imputation to get rid of the skewness. In general, all variables are used as predictors in the imputation models in hopes of reducing problems from uncongeniality (Meng, 1994). Uncongeniality refers to the situation when the model used by the analyst of the data differs from the model used for the imputation. This can lead to biased results, if the analyst's model is more complex than the imputation model and the imputation model omitted important relationships present in the original data. Since the true data generating model is usually unknown and an imputation model that is more complex than the true model only causes some loss in efficiency, the standard imputation strategy should be to include as many variables as possible in the imputation model (Little and Raghunathan (1997)). In the multinomial logit model for the categorical variables some explanatory variables are dropped for multicollinearity reasons. For the imputation procedure we use 26 variables from the GSSD and reduce the number of panel variables to be imputed to 48 to avoid multicollinearity problems (Figure 5.2 provides a broad description of the information contained in these variables).

5.4.2 Measuring the data utility

To evaluate the quality of the synthetic data, we compare analytic results achieved with the original data with results from the synthetic data. Basis is an analysis by Thomas Zwick: "Continuing Vocational Training Forms and Establishment Productivity in Germany" published in the German Economic Review, Vol. 6(2), pp. 155-184 in 2005. Since this analysis is used for validity evaluations in several chapters of the book, we provide a detailed description here.

Zwick analyses the productivity effects of different continuing vocational training forms in Germany. He argues that vocational training is one of the most important measures to gain and keep productivity in a firm. For his analysis he uses the waves 1997 to 2001 from the IAB Establishment Panel.

In 1997 and 1999 the Establishment Panel included the following additional question that was asked if the establishment did support continuous vocational training in the first part of 1997 or 1999 respectively: "For which of the following internal or external measures were employees exempted from

Information contained in the IAB Establishment Panel (wave 1997)	Information contained in the German Social Security Data (from 1997)
<p>Available for establishments in the survey</p> <ul style="list-style-type: none"> - number of employees in June 1996 - qualification of the employees - number of temporary employees - number of agency workers - working week (full-time and overtime) - the firm's commitment to collective agreements - existence of a works council - turnover, advance performance and export share - investment total - overall wage bill in June 1997 - technological status - age of the establishment - legal form and corporate position - overall company-economic situation - reorganisation measures - company further training activities - additional information on new foundations 	<p>Available for all German establishments with at least one employee covered by social security</p> <ul style="list-style-type: none"> - number of full-time and part-time employees - short-time employment - mean of the employees age - mean of wages from full-time employees - mean of wages from all employees - occupation - schooling and training - number of employees by gender - number of German employees
<p>Covered in both datasets</p> <ul style="list-style-type: none"> ➤ establishment number, branch and size ➤ location of the establishment ➤ number of employees in June 1997 	

Figure 5.2: Included variables from the IAB Establishment Panel and the German Social Security Data.

work or were costs completely or partly taken over by the establishment?" Possible answers were: formal internal training, formal external training, seminars and talks, training on the job, participation at seminars and talks, job rotation, self-induced learning, quality circles, and additional continuous vocational training. Zwick examines the productivity effects of these training forms and demonstrates that formal external training, formal internal training and quality circles do have a positive impact on productivity. Especially for formal external courses the productivity effect can be measured even two years after the training.

To detect why some firms offer vocational training and others not, Zwick runs a probit regression using the 1997 wave of the establishment panel. In the regression, Zwick uses two variables (*investment in IT* and the *codetermination of the employees*) that are only included in the 1998 wave of the establishment panel. Moreover, he excludes some observations based on information from other years. As we impute only the 1997 wave eliminating newly founded establishments, we have to rerun the regression, using all observations except

Table 5.1: Results from the vocational training regression for one stage full synthesis.

	original data	synthetic data	CI overlap
Redundancies expected	0.253***	0.293***	0.848
Many emp. exp. on maternity leave	0.262**	0.240	0.770
High qualification need exp.	0.646***	0.601***	0.227
Appr. tr. react. on skill shortages	0.113*	0.149*	0.930
Training react. on skill shortages	0.540***	0.532***	0.620
Establishment size 20-199	0.684***	0.649***	0.857
Establishment size 200-499	1.352***	1.215***	0.457
Establishment size 500-999	1.346***	1.404***	0.382
Establishment size 1000 +	1.955***	1.753***	0.932
Share of qualified employees	0.787***	0.812***	0.437
State-of-the-art tech. equipment	0.171***	0.186***	0.712
Collective wage agreement	0.255***	0.293***	0.901
Apprenticeship training	0.490***	0.423***	0.534
industry, East Germany dummies	Yes		

Notes: *** Significant at the 0.1% level, ** Significant at the 1% level,

* Significant at the 5% level

Source: IAB Establishment Panel 1997 without newly founded establishments and establishments not represented in the GSSD; regression according to Zwick (2005)

for newly founded establishments and deleting the two variables which are not part of the 1997 wave. We find that the results from the adjusted regression differ only slightly from the original regression. All the variables significant in Zwick's analysis are still significant. Only for the variable *high number of maternity leaves expected*, the significance level decreases from 1% to 5%.

For his analysis, Zwick runs the regression only on units with no missing values for the regression variables, losing all the information on establishments that did not respond to all variables used. This might lead to biased estimates if the assumption of a missing pattern that is completely at random (see for example Rubin (1987)) does not hold. For that reason, we compare the regression results from the synthetic datasets that by definition have no missing values, with the results, Zwick would have achieved if he would have run his regression on a dataset with all the missing values multiply imputed.

Comparing results from Zwick's regression run on the original data and on

the synthetic data are presented in Table 5.1. The last column of the table measures data utility by looking at the overlap between the confidence intervals for the estimates from the original data and the confidence intervals for the estimates from the synthetic data as described in Section 5.2. All variables in the regression except for the industry dummies that are part of the sampling design are synthesized. Since all imputation models (except for some categorical variables) are based on all variables in the dataset, the imputation model for the vocational training variable contains all the variables that are used in the regression. All estimates are close to the estimates from the real data and except for the variable *high number of maternity leaves expected*, that is not significant on any given significance level in the synthetic data, remain significant on the same level when using the synthetic data. The confidence interval overlap is high for most estimates, but it drops below 50% for four of the thirteen variables. Only for the dummy variable that indicates establishments with 200 to 499 employees and the dummy variable for establishments with more than 1,000 employees the absolute deviation between the estimates from the two datasets is higher than 0.1 (0.138 and 0.202 respectively). Obviously Zwick would have come nearly to the same conclusions in his analysis, if he would have used the synthetic data instead of the real data. See Drechsler *et al.* (2008b) for a two stage imputation approach that could further improve the quality of the synthetic data.

These results indicate that valid statistical inferences can be achieved using the synthetic datasets, but is the confidentiality of the survey respondents guaranteed? In our case disclosure of potentially sensitive information can be possible, when the following two conditions are fulfilled:

1. An establishment is included in the original dataset and in at least one of the newly drawn samples.
2. The original values and the imputed values for this establishment are nearly the same.

5.4.3 Assessing the disclosure risk

To determine the disclosure risk in our setting, we assume that the intruder would search for records that appear in more than one of the 10 new samples. Since the intruder doesn't know, if any establishment in the synthetic datasets

Table 5.2: How many records are sampled how often in the new samples?

Occurrence in sample(s)	number of records	percentage
1	45,553	82.75%
2	5,600	10.17%
3	1,805	3.28%
4	873	1.59%
5	507	0.92%
6	320	0.58%
7	164	0.30%
8	99	0.18%
9	45	0.08%
10	86	0.16%
Total	55,052	100%

is also included in the original dataset, he may use the probability of inclusion in the synthetic datasets as an estimator for the probability that this record is also included in the original survey. For example, if an establishment is included in all 10 new samples, the probability that this establishment is also included in the original sample will be very high, since we use the original sampling design for the 10 new samples.

Table 5.2 displays how often different records occur in the synthetic samples. Overall 55,052 establishments are sampled in the synthetic datasets. The vast majority are sampled only once or twice. Only roughly 7% of the establishments are sampled at least three times and less than 1% are sampled more than six times. But even if the intruder is able to identify records that are sampled more than once, which in itself is a difficult task, since almost all values are imputed and thus differ from sample to sample, he or she can not be sure whether this record really is included in the original survey. Table 5.3 displays how often the records from the original survey actually occur in the synthetic samples. 61.0 percent of the establishments included in the original survey do not occur in any of the 10 new drawn samples. 14.9 percent are contained in one of the 10 samples while only 5.5 percent can be found more than five times. Larger establishments have a higher probability of inclusion in the original survey (for some of the cells of the stratification matrix this probability is close to one). Since we use the same sampling design for drawing new establishments for our synthetic datasets, this means that larger estab-

Table 5.3: Establishments from the IAB Establishment Panel that also occur in at least one of the new samples.

Occurrence in sample(s)	number of records	percentage
None	4,469	61.0%
1	1,091	14.9%
2	535	7.3%
3	362	4.9%
4	275	3.8%
5	199	2.7%
6	144	2.0%
7	89	1.2%
8	53	0.7%
9	32	0.4%
10	83	1.1%
Total	7,332	100%

lishments also have a higher probability to be included in the original survey and in at least one of the new samples. Keeping that in mind, having only 25% of establishments between 200-999 employees and 49% of establishments with 1000+ employees in at least one of the new samples is a very good result in terms of data confidentiality (see Figure 5.3).

Comparing Table 5.2 and 5.3 we can see that only for the records that occur in all 10 datasets the probability that these records are also included in the original survey is very high. 96.5% (83 of the 86 records) of the establishments are contained in the original survey. But this probability decreases quickly. It is 71.1%, 53.5% and 54.3% for establishments that occur in 9, 8 and 7 samples respectively. For establishments that occur less than 7 times, the probability is always lower than 50%.

But even if a record is correctly identified, the intruder will only benefit from the identification, if the imputed values of these establishments are close to the original ones. The second step of our evaluation therefore takes a closer look at the establishments from the survey that appear at least once in the newly drawn samples. Using only these establishments the differences between original and imputed values can be detected. For each synthetic record that is also included in the original survey, we compare the imputed value to the true value. Binary variables tend to have a matching rate between 60 per-

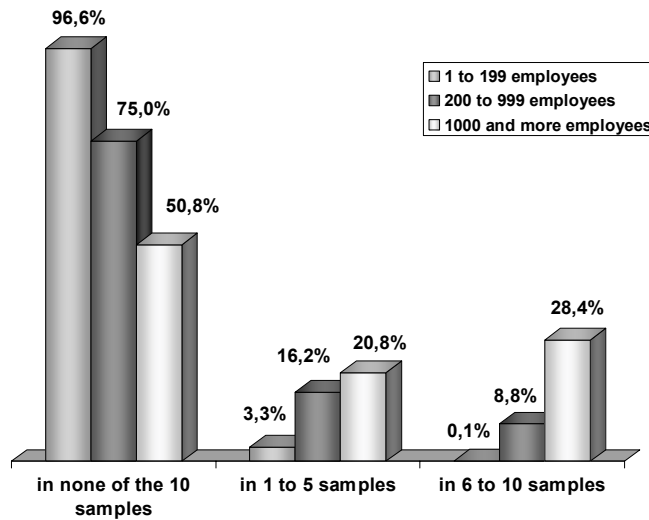


Figure 5.3: Occurrence of establishments already included in the original survey by establishment size.

cent and 90 percent, i.e. for 60 to 90 percent of these synthetic records the imputed binary value is the same as the true value from the survey. Multiple response questions with few categories show a high rate of identical answers in the total item block, too. But with an increase in the number of categories this rate decreases rapidly. For example, for an imputed multiple response variable consisting of 4 categories, the probability of having the same values for all 4 categories is about 57 percent. This probability decreases to about 6 percent if the number of categories increases to 13 (see Figure 5.4).

Imputed numeric variables always differ more or less from the original value. To evaluate the uncertainty for an intruder wanting to identify an establishment using the imputed data, we examine the variable *establishment size* for the 83 establishments that appear in all 10 datasets. The average relative difference between the imputed and the original values is 21%. A plot of the distribution of the relative difference for each record in each synthetic dataset shows that there are outliers for which the imputed values are two, three or even four times higher than the original ones (see Figure 5.5). Thus, for an intruder who wants to identify an establishment using his knowledge of the true size of the establishment, the imputed variable *establishment size* will hardly be of any use.

Summing up the second step, we find that for establishments, which are repre-

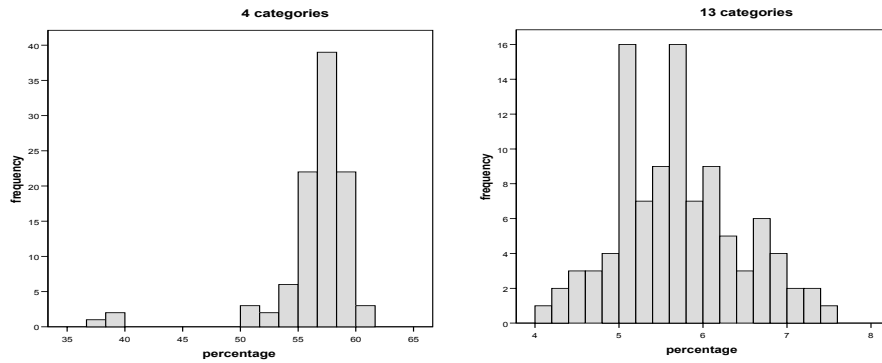


Figure 5.4: Distribution of the matching rates for different multiple response questions.

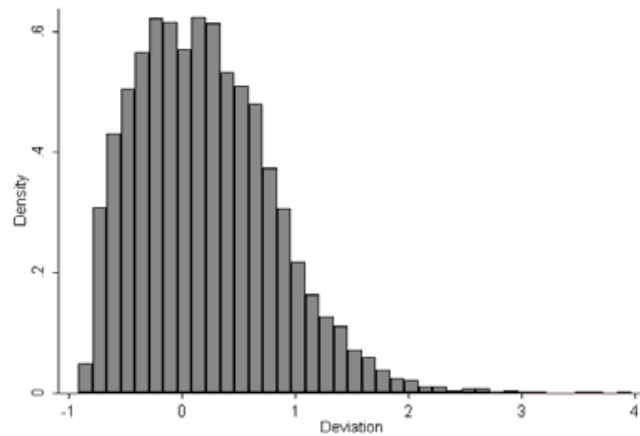


Figure 5.5: Histogram of the relative difference between original and imputed values for the variable *establishment size*.

sented in both datasets, up to 90 percent of some imputed binary variables are identical to the original values. But just one binary variable won't be sufficient to identify a single establishment. Using more binary variables, the risk of identical values will decrease quickly. If, for example, we assume the intruder needs five binary variables for identification and the variables are independently distributed, the risk will be $0.9^5 = 0.59$. Still, this only holds, if the establishment she or he is looking for is really included in the synthetic data which is very unlikely to begin with. Normally an intruder needs variables with more information than just two categories for a successful re-identification. But as shown for the variable *establishment size*, the chance of identifying an establishment by combining information from numeric and categorical variables is very low.

These results together with the results for the data utility in Section 5.4.2 indicate that a release of the described subset of the data would be possible. Of course the data utility for different estimates should be evaluated in detail for different kinds of estimates before an actual release.

Chapter 6

Partially Synthetic Datasets¹

As of this writing, no agency adopted the fully synthetic approach discussed in the previous chapter, but some agencies have adopted a variant of Rubin's original approach, suggested by Little (1993): release datasets comprising the units originally surveyed with some collected values, such as sensitive values or values of key identifiers, replaced with multiple imputations. These are called *partially synthetic datasets*. For example, the U.S. Federal Reserve Board protects data in the Survey of Consumer Finances by replacing large monetary values with multiple imputations (Kennickell (1997)). In 2007 the U.S. Census Bureau released a partially synthetic, public use file for the Survey of Income and Program Participation (SIPP) that includes imputed values of social security benefits information and dozens of other highly sensitive variables (www.sipp.census.gov/sipp/synth_data.html). The Census Bureau also protects the identities of people in group quarters (e.g., prisons, shelters) in the public use files of the American Communities Survey by replacing demographic data for people at high disclosure risk with imputations. Partially synthetic, public use datasets are in the development stage in the U.S. for the Longitudinal Business Database, the Longitudinal Employer-Household Dynamics survey, and the American Communities Survey veterans and full sample data.

¹Most of this chapter is taken from Drechsler *et al.* (2008a) and Drechsler and Reiter (2008).

6.1 Inference for partially synthetic datasets

Following Reiter (2003, 2004), let $Z_j = 1$ if unit j is selected to have any of its observed data replaced, and let $Z_j = 0$ otherwise. Let $Z = (Z_1, \dots, Z_s)$, where s is the number of records in the observed data. Let $Y = (Y_{rep}, Y_{nrep})$ be the data collected in the original survey, where Y_{rep} includes all values to be replaced with multiple imputations and Y_{nrep} includes all values not replaced with imputations. Let $Y_{rep}^{(i)}$ be the replacement values for Y_{rep} in synthetic dataset i . Each $Y_{rep}^{(i)}$ is generated by simulating values from the posterior predictive distribution $f(Y_{rep}^{(i)}|Y, Z)$, or some close approximation to the distribution such as those of Raghunathan *et al.* (2001). The agency repeats the process m times, creating $D^{(i)} = (Y_{nrep}, Y_{rep}^{(i)})$ for $i = 1, \dots, m$, and releases $\mathbf{D} = \{D^{(1)}, \dots, D^{(m)}\}$ to the public.

To get valid inferences, secondary data users can use the combining rules presented by Reiter (2003). Let Q be an estimand, such as a population mean or regression coefficient. Suppose that, given the original data, the analyst would estimate Q with some point estimator q and the variance of q with some estimator u . Let $q^{(i)}$ and $u^{(i)}$ be the values of q and u in synthetic dataset $D^{(i)}$, for $i = 1, \dots, m$. The analyst computes $q^{(i)}$ and $u^{(i)}$ by acting as if each $D^{(i)}$ is the genuine data. The following quantities are needed for inferences for scalar Q :

$$\bar{q}_m = \sum_{i=1}^m q_i/m \quad (6.1)$$

$$b_m = \sum_{i=1}^m (q_i - \bar{q}_m)^2/(m-1) \quad (6.2)$$

$$\bar{u}_m = \sum_{i=1}^m u_i/m . \quad (6.3)$$

The analyst then can use \bar{q}_m to estimate Q and

$$T_p = b_m/m + \bar{u}_m \quad (6.4)$$

to estimate the variance of \bar{q}_m .

Similar to the variance estimator for multiple imputation of missing data, b_m/m is the correction factor for the additional variance due to using a finite number of imputations. However, the additional b_m , necessary in the missing data context, is not necessary here, since \bar{u}_m already captures the variance of

Q given the observed data. This is different in the missing data case, where \bar{u}_m is the variance of Q given the completed data and $\bar{u} + b_m$ is the variance of Q given the observed data.

When n is large, inferences for scalar Q can be based on t -distributions with degrees of freedom $\nu_m = (m-1)(1+r_m^{-1})^2$, where $r_m = (m^{-1}b_m/\bar{u}_m)$. Methods for multivariate inferences are developed in Reiter (2005c). The variance estimate T_p can never be negative, so no adjustments are necessary for partially synthetic datasets.

6.2 Data utility for partially synthetic datasets

To evaluate the data utility of partially synthetic datasets, we can use the same methods as for fully synthetic datasets. Namely, measuring the confidence interval overlap between confidence intervals obtained from the synthetic data and confidence intervals obtained from the original data or measuring how well one can discriminate between the original and the synthetic data based on the ideas of propensity score matching. See Section 5.2 for details.

6.3 Disclosure risk for partially synthetic datasets

The disclosure risk is higher for partially synthetic datasets than it is for fully synthetic datasets, especially if the intruder knows that some unit participated in the survey, since true values remain in the dataset and imputed values are generated only for the survey participants and not for the whole population. Thus for partially synthetic datasets assessing the risk of disclosure is an equally important evaluation step as assessing the data utility. It is essential that the agency identifies and synthesizes all variables that bear a risk of disclosure. A conservative approach would be, to also impute all variables that contain the most sensitive information. Once the synthetic data is generated, careful checks are necessary to evaluate the disclosure risk for these datasets. Only if the datasets prove to be useful both in terms of data utility and in terms of disclosure risk, a release should be considered.

As noted above, the risk of disclosure significantly increases, if the intruder knows, who participated in a survey. Thus, it is important to distinguish between a scenario, in which the intruder knows that the target she is looking for is in the data and a scenario, in which the intruder has some external information, but does not know, whether any of the targets she is looking for, actually is included in the survey. For most surveys the latter case will be a more realistic assumption, but there might be situations in which it is publicly known who participated in a survey or the agency might want to release a complete synthetic population. We therefore start by presenting methods to evaluate the disclosure risk under the conservative assumption that the intruder has full information about survey participation and afterwards discuss necessary extensions to account for the additional sampling uncertainty, if the intruder does not have any response knowledge. Both methods only evaluate the risk of *identification disclosure*, i.e. the risk that a unit is correctly identified in the released data. Methods to evaluate the risk of *inferential disclosure*, i.e. the amount of additional information an intruder might obtain about a unit for which he or she already knows that it participated in the survey, still need to be developed for partially synthetic datasets. The concept of differential privacy described in Section 5.3 might be useful in this context. Future research is still needed on this topic.

6.3.1 Ignoring the uncertainty from sampling

To evaluate disclosure risks if the intruder knows which units are included in the released data, we can compute probabilities of identification by following the approach of Reiter and Mitra (2009). Related approaches are described by Duncan and Lambert (1989), Fienberg *et al.* (1997), and Reiter (2005a). Roughly, in this approach we mimic the behavior of an ill-intentioned user of the released data who possesses the true values of the quasi-identifiers for selected target records (or even the entire database). To illustrate, suppose the malicious user has a vector of information, \mathbf{t} , on a particular target unit in the population corresponding to a unit in the m released simulated datasets, $\mathbf{D} = \{D^{(1)}, \dots, D^{(m)}\}$. Let t_0 be the unique identifier (e.g., establishment name) of the target, and let d_{j_0} be the (not released) unique identifier for record j in \mathbf{D} , where $j = 1, \dots, s$. Let M be any information released about the simulation models.

The malicious user's goal is to match unit j in \mathbf{D} to the target when $d_{j0} = t_0$, and not to match when $d_{j0} \neq t_0$ for any $j \in \mathbf{D}$. Let J be a random variable that equals j when $d_{j0} = t_0$ for $j \in \mathbf{D}$ and equals $s + 1$ when $d_{j0} = t_0$ for some $j \notin \mathbf{D}$. The malicious user thus seeks to calculate the $Pr(J = j|\mathbf{t}, \mathbf{D}, M)$ for $j = 1, \dots, s + 1$. She then would decide whether or not any of the identification probabilities for $j = 1, \dots, s$ are large enough to declare an identification. Note that in this scenario $Pr(J = s + 1|\mathbf{t}, \mathbf{D}, M) = 0$ because the intruder knows that the target record she is looking for is included in the released data. Because the malicious user does not know the actual values in Y_{rep} , she should integrate over its possible values when computing the match probabilities. Hence, for each record in \mathbf{D} we compute

$$Pr(J = j|\mathbf{t}, \mathbf{D}, M) = \int Pr(J = j|\mathbf{t}, \mathbf{D}, Y_{rep}, M)Pr(Y_{rep}|\mathbf{t}, \mathbf{D}, M)dY_{rep}. \quad (6.5)$$

This construction suggests a Monte Carlo approach to estimating each $Pr(J = j|\mathbf{t}, \mathbf{D}, M)$. First, sample a value of Y_{rep} from $Pr(Y_{rep}|\mathbf{t}, \mathbf{D}, M)$. Let Y^{new} represent one set of simulated values. Second, compute $Pr(J = j|\mathbf{t}, \mathbf{D}, Y_{rep} = Y^{new}, M)$ using exact or, for continuous synthesized variables, distance-based matching assuming Y^{new} are collected values. This two-step process is iterated R times, where ideally R is large, and (1) is estimated as the average of the resultant R values of $Pr(J = j|\mathbf{t}, \mathbf{D}, Y_{rep} = Y^{new}, M)$. When M has no information, the malicious user can treat the simulated values as plausible draws of Y_{rep} .

To illustrate, suppose that region and employee size are the only quasi-identifiers in a survey of establishments. A malicious user seeks to identify an establishment in a particular region of the country with 125 employees. The malicious user knows that this establishment is in the sample. Suppose that the agency releases m datasets after simulating only employment size, without releasing information about the imputation model. In each $D^{(i)}$, the malicious user would search for all establishments matching the target on region and having synthetic employee size within some interval around 125, say 110 to 140. The agency selects the intervals for employment size based on its best guess of the amount of uncertainty that intruders would be willing to tolerate when estimating true employee sizes. Let $N^{(i)}$ be the number of records in $D^{(i)}$ that meet these criteria. When no establishments with all of those characteristics are in $D^{(i)}$, set $N^{(i)}$ equal to the number of establishments

in the region, i.e., match on all non-simulated quasi-identifiers. For any j ,

$$Pr(J = j|\mathbf{t}, \mathbf{D}, M) = (1/m) \sum_i (1/N^{(i)}) (Y_j^{new,i} = \mathbf{t}), \quad (6.6)$$

where $(Y_j^{new,i} = \mathbf{t}) = 1$ when record j is among the $N^{(i)}$ matches in $D^{(i)}$ and equals zero otherwise. Similar computations arise when simulating region and employee size: the malicious user exactly matches on the simulated values of region and distance-based matches on employee size to compute the probabilities.

Following Reiter (2005a), we quantify disclosure risk with summaries of these identification probabilities. It is reasonable to assume that the malicious user selects as a match for \mathbf{t} the record j with the highest value of $Pr(J = j|\mathbf{t}, \mathbf{D}, M)$, if a unique maximum exists. We consider three risk measures: the *expected match risk*, the *true match risk*, and the *false match rate*. To calculate these, we need some further definitions. Let c_j be the number of records in the dataset with the highest match probability for the target t_j for $j = 1, \dots, s$; let $I_j = 1$ if the true match is among the c_j units and $I_j = 0$ otherwise. Let $K_j = 1$ when $c_j I_j = 1$ and $K_j = 0$ otherwise. The *expected match risk* can now be defined as $\sum_j (1/c_j) I_j$. When $I_j = 1$ and $c_j > 1$, the contribution of unit j to the expected match risk reflects the intruder randomly guessing at the correct match from the c_j candidates. The *true match risk* equals $\sum_j K_j$. Finally, let $F_j = 1$ when $c_j(1 - I_j) = 1$ and $F_j = 0$ otherwise; and, let s equal the number of records with $c_j = 1$. The *false match rate* equals $\sum F_j/s$. It is important to note that these summary statistics are helpful to summarize the overall disclosure risk for the complete data, but the real advantage of the suggested measures is the fact that the identification probabilities are calculated on the record level. This enables disclosure risk evaluations for specified subgroups of the data. In some situations only a few records in the dataset might be correctly identified, but all identified records belong to the same subgroup. In this case, the overall measure that indicates a low disclosure risk might be misleading since the risk of disclosure e.g. for the largest establishments in the dataset might still be very high.

6.3.2 Accounting for the uncertainty from sampling

If the intruder does not know, if the target, he or she is looking for participated in the survey, the fact that the survey usually only comprises a sample of the

population adds an additional layer of protection to the released data. In this case we can use the extensions to the measures described above suggested by Drechsler and Reiter (2008). We simply have to replace $N_{\mathbf{t},i}$ in (6.6) with $F_{\mathbf{t}}$, the number of records in the population that match the target on region and establishment size in the above example. When the intruder and the agency do not know $F_{\mathbf{t}}$, it can be estimated using the approach in Elamir and Skinner (2006), which assumes that the population counts follow an all-two-way-interactions log-linear model. The agency can determine the estimated counts, $\hat{F}_{\mathbf{t}}$, by fitting this log-linear model with D_{obs} . Alternatively, since D_{obs} is in general not available to intruders, the agency can fit a log-linear model with each D_i , resulting in the estimates $\hat{F}_{\mathbf{t},i}$ for $i = 1, \dots, m$. We note that in this scenario $Pr(J = s + 1 | \mathbf{t}, \mathbf{D}, M) = 1 - \sum_{j=1}^s Pr(J = j | \mathbf{t}, \mathbf{D}, M)$.

For some target records, the value of $N_{\mathbf{t},i}$ might exceed $F_{\mathbf{t}}$ (or $\hat{F}_{\mathbf{t}}$ if it is used). It should not exceed $\hat{F}_{\mathbf{t},i}$, since $\hat{F}_{\mathbf{t},i}$ is required to be at least as large as $N_{\mathbf{t},i}$. For such cases, we presume that the intruder sets $Pr(J = s + 1 | \mathbf{t}, \mathbf{D}, M) = 0$ and picks one of the matching records at random. To account for this case, we can re-write (6.5) for $j = 1, \dots, s$ as

$$Pr(J = j | \mathbf{t}, \mathbf{D}, M) = (1/m) \sum_i \min(1/F_{\mathbf{t}}, 1/N_{\mathbf{t},i}) (Y_{ij}^{\text{new}} = \mathbf{t}) . \quad (6.7)$$

We can use the three summary statistics of the identification probabilities described in Section 6.3.1, with the important difference that we also have to consider $Pr(J = s + 1 | \mathbf{t}, \mathbf{D}, M)$, the probability for a match outside the sample. In many cases this will be the highest match rate. It is reasonable to assume that the intruder does not match whenever $Pr(J = s + 1 | \mathbf{t}, \mathbf{D}, M)$ is the maximum probability for the target. If this assumption is considered to strong, the data disseminating agency can define a threshold γ and assume that the intruder matches to the released data only when $Pr(J = s + 1 | \mathbf{t}, \mathbf{D}, M) \leq \gamma$, where $0 \leq \gamma \leq 1$.

6.4 Application of the partially synthetic approach to the IAB Establishment Panel

To achieve results that can be compared to the results in Section 5.4, we use the same subset of variables from the wave 1997 as in the fully synthetic application (see Section 5.4 for a description of the variables selected).

For the partially synthetic datasets, we replace only two variables (the *number of employees* and the *industry*, coded in 16 categories) with synthetic values, since these are the only two variables that might lead to disclosure in the analyses we use to evaluate the data utility of the synthetic datasets. If we intended to release the complete data to the public, some other variables would have to be synthesized, too. Identifying all the variables that provide a potential disclosure risk is an important and labour intensive task. Nevertheless, the two variables mentioned above definitely impose a high risk of disclosure, since they are easily available in public databases and especially large firms can be identified without difficulty using only these two variables. We define a multinomial logit model for the imputation of the industry code and a linear model stratified by four establishment size classes defined by quartiles for the number of employees. For the partially synthetic datasets, we use the same number of variables in the imputation model as in the fully synthetic data example (26 from the German Social Security Data (GSSD)), 48 from the establishment panel), but the original sample is used and no additional samples are drawn from the GSSD. We generate the same number of synthetic datasets, but the modeling is performed using own coding in *R*.

6.4.1 Measuring the data utility

For an evaluation of the utility of the partially synthetic data, we compare analytic results achieved with the original data with results from the synthetic data. The regression results in Table 6.1 are again based on the analysis by Zwick (2005) described in detail in Section 5.4.2.²

All estimates are very close to the estimates from the real data and except for the variables *many employees expected on maternity leave* and *apprenticeship training reaction on skill shortages* for which the significance level increases from 1% to 0.1% and from 5% to 1% respectively, remain significant on the same level when using the synthetic data. With an average of 0.925 over all 13 estimates, the confidence interval overlap is very high. Only the effect of the largest establishment size class is slightly underestimated leading to a reduced

²For simplicity, we impute all missing values first and treat one fully imputed dataset as the original data. Since missing rates are low for all variables used in the regression, results for the original data only change in the third digit compared to the results in Table 5.1. See Chapter 7 on how to correctly generate synthetic datasets from data that is subject to nonresponse.

Table 6.1: Results from the vocational training regression for one stage partial synthesis.

	original data	synthetic data	CI overlap
Redundancies expected	0.250***	0.259***	0.956
Many emp. exp. on maternity leave	0.267**	0.316***	0.869
High qualification need exp.	0.648***	0.653***	0.982
Appr. tr. react. on skill shortages	0.115*	0.121**	0.969
Training react. on skill shortages	0.539***	0.547***	0.962
Establishment size 20-199	0.682***	0.695***	0.920
Establishment size 200-499	1.350***	1.335***	0.936
Establishment size 500-999	1.344***	1.344***	0.994
Establishment size 1000 +	1.956***	1.754***	0.685
Share of qualified employees	0.789***	0.803***	0.948
State-of-the-art tech. equipment	0.170***	0.175***	0.962
Collective wage agreement	0.257***	0.275***	0.894
Apprenticeship training	0.488***	0.496***	0.953
industry, East Germany dummies	Yes		

Notes: *** Significant at the 0.1% level, ** Significant at the 1% level,
* Significant at the 5% level

Source: IAB Establishment Panel 1997 without newly founded establishments and establishments not represented in the GSSD; regression according to Zwick (2005)

overlap of 0.685. For all other estimates, the overlap is above 0.85, indicating a very high quality of the synthetic data. Obviously Zwick would have come to the same conclusions in his analysis, if he would have used the partially synthetic data instead of the real data.

6.4.2 Assessing the disclosure risk

To evaluate the risk of disclosure we apply the disclosure risk measures described in Section 6.3.1, i.e. we assume, the intruder knows, who participated in the survey. We further assume, the intruder knows the true values for the number of employees and industry. This is a conservative scenario but gives, in some sense, an upper bound on the risk for this level of intruder knowledge. For an application of the disclosure risk measures without response knowledge, see Section 7.4.4. Intruders might also know other variables on the file,

in which case the agency may need to synthesize them as well. The intruder computes probabilities using the approach outlined in Section 6.3.1. We assume that the agency does not reveal the synthesis model to the public, so that the only information in M is that establishment size and industry were synthesized. For a given target t , records from each $D^{(i)}$ must meet two criteria to be possible matches. First, the record's synthetic industry code exactly matches the target's true industry code. Second, the record's synthetic number of employees lies within an agency-defined interval around the target's true number of employees. Acting as the agency, we define the interval as follows. We divide the true number of employees (transformed by taking the cubic root) into twenty quantiles and calculate the standard deviation of the number of employees within each quantile. The interval is $t_e \pm sd_s$, where t_e is the target's true value and sd_s is the standard deviation of the quantile in which the true value falls. When there are no synthetic records that fulfill both matching criteria, the intruder matches only on the industry code. We use 20 quantiles because this is the largest number of groups that guarantees at least some variation within each group. Using a larger number of quantiles results in groups with only one value of employment, which forces exact matching for targets in those quantiles. On the other hand, using a small number of quantiles does not differentiate adequately between small and large establishments. For small establishments, we want the potential matches to deviate only slightly from the original values. For large establishments, we accept higher deviations.

Given this matching scenario the expected match risk and the true match risk both would be 139, i.e. the intruder would get 139 true correct single matches from the 7,332 records in her target file. The false match rate would be 98.1%. There is no obvious common pattern for the identified records. Neither for the region nor for the industry the distribution of the identified records differs significantly from the distribution in the underlying data. The identified records consist of very small and very large establishments. However, as one might expect, the actual risk of disclosure depends on establishment size. While only 1.38% of the establishments with less than 100 employees are identified, this rate increases to 1.87% for establishments with 100-1,000 employees and to 5.21% for establishments with more than 1,000 employees. Considering the fact that the intruder matches on 7,332 records and never knows which of the 7,330 single matches she obtains actually are correct matches the risk is very moderate. Especially since these measures are based on the very conservative

assumptions that (i) the intruder knows who participated in the survey and (ii) has exact information on the industry code and the establishment size for all the survey participants. If the agency deems the risk of disclosure still too high, it might broaden the industry codes or suppress this information completely in the released file. Another possibility would be to use less detailed models for the large establishments to ensure a higher level of perturbation for these records. As an alternative, the agency might consider releasing fully synthetic datasets instead.

6.5 Pros and cons of fully and partially synthetic datasets

Obviously there are advantages and disadvantages for both, the partially and the fully synthetic approach. The fully synthetic approach provides a very high level of disclosure protection rendering the identification of single units in the released data almost impossible. Partially synthetic datasets can not offer such a high level of protection per se, since true values remain in the data and synthetic values are only generated for units that participated in the survey. This means that evaluating the disclosure risk is an equally important step as evaluating the data quality for partially synthetic datasets.

Nevertheless, partially synthetic datasets have the important advantage that in general the data utility will be higher, since only for some variables the true values have to be replaced with imputed values, so by definition the joint distribution for all the unchanged variables will be exactly the same as in the original dataset. The quality of the synthetic datasets will highly depend on the quality of the underlying models and for some variables it will be very hard to define good models, especially if logical constraints and skip patterns should be preserved. But if these variables do not contain any sensitive information or information that might help identify single respondents, why bother to find these models? Why bother to perturb these variables first place? Furthermore, the risk of biased imputations will increase with the number of variables that are imputed, if the SRMI approach (see Section 3.2.2) is used for imputations. For, if one of the variables is imputed based on a *bad* model, the biased imputed values for that variable could be the basis for the imputation of another variable and this variable again could be used for the imputation of another

one and so on. So a small bias could increase to a really problematic bias over the imputation process. A comparison of the results in Sections 5.4.2 and 6.4.1 underline these thoughts. The partially synthetic datasets provide higher data quality in terms of lower deviation from the true estimates and higher confidence interval overlap between estimates from the original data and estimates from the synthetic data almost for all estimates. Still, this increase of data utility comes at the price of an increase in the risk of disclosure. Although the disclosure risk for fully synthetic datasets might not be zero, the disclosure risk will definitely be higher if true values remain in the dataset and the released data is based only on survey participants. Thus, it is important to make sure that all variables that might lead to disclosure are imputed in a way that confidentiality is guaranteed. This means that a variety of disclosure risk checks are necessary before the data can be released, but this is a problem common to all perturbation methods that are based only on the information from the survey respondents. Agencies willing to release synthetic public use files will have to consider carefully, which approach suites best for their datasets. If the data consists only of all small number of variables and imputation models are easy to set up, the agencies might consider releasing fully synthetic datasets, since these datasets will provide the highest confidentiality protection, but if there are many variables in the data considered for release and the data contain a lot of skip patterns, logical constraints and questions that are asked only to a small subgroup of survey respondents, the agencies might be better off to release partially synthetic datasets and include a detailed disclosure risk study in their evaluation of the quality of the datasets considered for release.

Chapter 7

Multiple Imputation for Nonresponse and Statistical Disclosure Control

Most if not all surveys are subject to item nonresponse and even registers can contain missing values, if implausible values are set to missing during the data editing process. Since the generation of partially synthetic datasets is based on the ideas of multiple imputation, it is reasonable to use the approach to impute missing values and generate synthetic values simultaneously. The imputation of missing values is not an issue for fully synthetic datasets, since the original data is only used for model building.

At a first glance, it seems logical, to impute missing values and generate synthetic values in one step, using the same model from the originally observed values. However, as Reiter (2004) points out, this can lead to biased imputations, if only a subset of the data, e.g. the income for units with income above \$100,000, should be replaced with synthetic values, but the imputation model for the missing values is based on the entire dataset. To allow for different models, Reiter (2004) suggests imputation in two stages. On the first stage, all missing values are imputed m times using the standard multiple imputation approach for nonresponse (see Chapter 3). On the second stage, all values that need to be replaced are synthesized r times in every first stage nest leading to a total of $M = m * r$ datasets that are released to the public. Each released dataset includes a label indicating from which first stage imputed dataset it was generated.

7.1 Inference for partially synthetic datasets when the original data is subject to non-response

The two stage imputation described above generates two sources of variability. The first, when missing values are imputed, the second, when sensitive or identifying variables are replaced with synthetic values. Neither the combining rules for the imputation of missing values described in Section 3.1 nor the combining rules for partially synthetic datasets described in Section 6.1 correctly reflect these two sources of variability. Reiter (2004) derived the combining rules necessary to obtain valid inferences in this two stage setting:

Again, let Q be an estimand, such as a population mean or regression coefficient. Suppose that, given the original data, the analyst would estimate Q with some point estimator q and the variance of q with some estimator u . Let $q_i^{(l)}$ and $u_i^{(l)}$ be the values of q and u in synthetic dataset $D_i^{(l)}$, for $l = 1, \dots, m$ and $i = 1, \dots, r$. The analyst computes $q_i^{(l)}$ and $u_i^{(l)}$ by acting as if each $D_i^{(l)}$ is the genuine data. The following quantities are needed for inferences for scalar Q :

$$\bar{q}_M = \sum_{l=1}^m \sum_{i=1}^r q_i^{(l)} / (mr) = \sum_{l=1}^m \bar{q}^{(l)} / m \quad (7.1)$$

$$\bar{b}_M = \sum_{l=1}^m \sum_{i=1}^r (q_i^{(l)} - \bar{q}^{(l)})^2 / m(r-1) = \sum_{l=1}^m b^{(l)} / m \quad (7.2)$$

$$B_M = \sum_{l=1}^m (\bar{q}^{(l)} - \bar{q}_M)^2 / (m-1) \quad (7.3)$$

$$\bar{u}_M = \sum_{i=1}^m \sum_{i=1}^r u_i^{(l)} / (mr) . \quad (7.4)$$

The analyst then can use \bar{q}_M to estimate Q and

$$T_M = (1 + 1/m)B_M - \bar{b}_M/r + \bar{u}_M \quad (7.5)$$

to estimate the variance of \bar{q}_M .

When n is large, inferences for scalar Q can be based on t -distributions with degrees of freedom

$$\nu_M = \left(\frac{((1 + 1/m)B_M)^2}{(m-1)T_M^2} + \frac{(\bar{b}_M/r)^2}{m(r-1)T_M^2} \right)^{-1} \quad (7.6)$$

Methods for multivariate inferences are developed in Kinney and Reiter (2010). Similar to the variance estimate for fully synthetic datasets, T_M can become negative, since \bar{b}_M/r is subtracted. In this case Reiter (2008b) suggests to use the conservative variance estimator $T_M^{adj} = (1 + 1/m)B_m + \bar{u}_M$. This estimator is equivalent to the variance estimator for multiple imputation for missing data. Consequently the degrees of freedom is given by:

$$\nu_M^{adj} = (m - 1)(1 + m\bar{u}_M/((m + 1)B_M))^2 \quad (7.7)$$

Generally negative variances can be avoided by increasing m and r .

7.2 Data utility and disclosure risk

To evaluate the data utility in this setting, we can use the same measures as for fully synthetic or partially synthetic datasets. Namely, measuring the confidence interval overlap between confidence intervals obtained from the synthetic data and confidence intervals obtained from the original data or measuring how well one can discriminate between the original and the synthetic data based on the ideas of propensity score matching (see Section 5.2). The difference to the standard one stage synthesis is that we compare the synthetic datasets with the datasets imputed on stage one.

For disclosure risk evaluations the disclosure risk measures described in Section 6.3 can be used. Depending on the scenario, measures that assume the intruder knows who participated in a survey (see Section 6.3.1) or measures that consider the additional uncertainty from sampling (see Section 6.3.2) can be applied.

7.3 Multiple imputation of the missing values in the IAB Establishment Panel¹

In the remainder of this chapter, we describe all the steps that were necessary to generate a scientific use file of the wave 2007 of the IAB Establishment Panel that will be released in fall 2009. We start by illustrating the extensive imputation task required to impute all missing values in the dataset. We briefly

¹Most of this section is taken from Drechsler (2009).

discuss, how we selected the variables to be synthesized. We also describe the synthesis process and the models we implemented for the synthesis. Finally, we present results from the data utility and disclosure risk evaluations that we preformed before the actual release.

7.3.1 The imputation task

Most of the 284 variables included in the wave 2007 of the panel are subject to nonresponse. Only 26 variables are fully observed. However, missing rates vary considerably between variables and are modest for most variables. 65.8% of the variables have missing rates below 1%, 20.4% of the variables have missing rates between 1% and 2%, 15.1% rates between 2% and 5% and only 12 variables have missing rates above 5%. The five variables with missing rates above 10% are *subsidies for investment and material expenses* (13.6%), *payroll* (14.4%), *intermediate inputs as proportion of turnover* (17.4%), *turnover in the last fiscal year* (18.6%), and *number of workers who left the establishment due to restructuring measures* (37.5%). Obviously, the variables with the highest missing rates contain information that is either difficult to provide like *number of workers who left the establishment due to restructuring measures* or considered sensitive like *turnover in the last fiscal year*. The variable *number of workers who left the establishment due to restructuring measures* is only applicable to 626 establishments in the dataset, who declared they had restructuring measures in the last year. Of these 626 only 391 establishments provided information on the number of workers that left the establishment due to these measures. Clearly, it is often difficult to tailor exactly which workers left as a result of the measures and which left for other reasons. This might be the reason for the high missing rates. The low number of observed values is also problematic for the modeling task, so this variable should be used with caution in the imputed dataset.

7.3.2 Imputation models

Since the dataset contains a mixture of categorical variables and continuous variables with skewed distributions and a variety of often hierarchical skip patterns and logical constraints, it is impossible to apply the joint modeling approach described in Section 3.2.1. We apply the fully conditional specification approach described in Section 3.2.2, iteratively imputing one variable at

a time, conditioning on the other variables available in the dataset. For the imputation we basically rely on three different imputation models. The linear model for the continuous variables, the logit model for binary variables and the multinomial logit for categorical variables with more than two categories. Multiple imputation procedures for these models are described in Raghunathan *et al.* (2001). In general, all variables that don't contain any structural missings are used as predictors in the imputation models in hopes of reducing problems from uncongeniality (Meng, 1994). In the multinomial logit model for the categorical variables the number of explanatory variables is limited to 30 variables found by stepwise regression to speed up the imputation process. To improve the quality of the imputation we define several separate models for the variables with high missing rates like *turnover* or *payroll*. Independent models are fit for East and West Germany and for different establishment size classes.

All continuous variables are subject to non-negativity constraints and the outcome of many variables is further restricted by linear constraints. To complicate the imputation process most variables have huge spikes at zero and as mentioned before the filtering rules are often hierarchical. Simply applying standard imputation procedures can lead to biased or inconsistent imputations in this context. We therefore have to rely on a mixture of the adjustments presented in Section 3.3. Since the package `mi` was not available at the beginning of this project and other standard packages could not deal with all these problems or did not allow detailed model specification, we use own coding in *R* for the imputation routines to generate $m = 5$ datasets.

7.3.3 Evaluating the quality of the imputations

It is more difficult to evaluate the quality of the imputations for missing values than evaluating the quality of the imputations for statistical disclosure control (SDC). With the latter, we can simply compare any statistic obtained from the protected data with the same statistic obtained from the original data, since we have the exact information what the correct outcome should be. Methods for evaluating the data quality of synthetic datasets are described in Section 5.2. When imputing missing values, this information by definition is not available and the assumption that the response mechanism is ignorable (Rubin, 1987), necessary for obtaining valid imputations if the response mechanism is

not modeled directly, can not be tested with the observed data. A response mechanism is considered ignorable, if, given that the sampling mechanism is ignorable, the response probability only depends on the observed information.² If these conditions are fulfilled, the missing data is said to be *missing at random (MAR)* and imputation models only need to be based on the observed information. As a special case, the missing data is said to be *missing completely at random (MCAR)*, if the response mechanism does not depend on the data (observed or unobserved), which implies that the distribution of the observed data and the distribution of the missing data are identical. If the above requirements are not fulfilled, the missing data is said to be *missing not at random (MNAR)* and the response mechanism needs to be modeled explicitly. Little and Rubin (2002) provide examples for non-ignorable missing-data models.

As noted before, it is not possible to check, if the missing data is *MAR* with the observed data. But even if the *MAR* assumption can not be tested, this does not mean, the imputer can not test the quality of his or her imputations at all. Abayomi *et al.* (2008) suggest several ways of evaluating model based imputation procedures. Basically their ideas can be divided in two categories: On the one hand, the imputed data can be checked for reasonability. Simple distributional and outlier checks can be evaluated by subject matter experts for each variable to avoid implausible imputed values like a turnover of \$ 10 million for a small establishment in the social sector. On the other hand, since imputations usually are model based, the fit of these models can and indeed should be evaluated. Abayomi *et al.* (2008) label the former as *external diagnostic techniques*, since the imputations are evaluated using outside knowledge and the latter *internal diagnostic techniques*, since they evaluate the modeling based on model fit without the need of external information.

To automate the external diagnostics to some extent, Abayomi *et al.* (2008) suggest to use the Kolmogorov Smirnov test to flag any imputations for which the distribution of the imputed values significantly differs from the distribution of the observed values. Of course a significant difference in the distributions does not necessarily indicate problems with the imputation. Indeed, if the

²The additional requirement that the sampling mechanism is also ignorable (Rubin, 1987), i.e. the sampling probability only depends on observed data, is usually fulfilled in scientific surveys. The stratified sampling design of the IAB Establishment Panel also satisfies this requirement.

missing data mechanism is *MAR*, but not *MCAR* we would expect the two distributions to differ. The test is only intended to decrease the number of variables that need to be checked manually implicitly assuming that no significant difference between the original and the imputed data indicates no problem with the imputation model.

However, we are skeptical about this automated selection method, since the test is sensitive to the sample size, so the chance of rejecting the null hypothesis will be lower for variables with lower missing rates and variables that are answered only by a subset of the respondents. Furthermore it is unclear what significance level to choose and as noted above, rejection of the null hypothesis does not necessarily indicate an imputation problem, but not rejecting the null hypothesis is not a guarantee that we found a good imputation model either. However, this is implicitly assumed by this procedure.

For the continuous variables, we searched for possible flaws in the imputations by plotting the distributions for the original and imputed values for every variable. We checked, if any notable differences between these distributions can be justified by differences in the distributions of the covariates. Figure 7.1 displays the distributions for two representative variables based on kernel density estimation. Original values are represented with a solid line, imputed values with a dashed line. Both variables are reported on the log-scale. The left variable (*payroll*) represents a candidate that we did not investigate further, since the distributions almost match exactly. The right variable (*number of participants in further education (NB.PFE)*) is an example for a variable for which we tried to understand the difference between the distribution of the observed values and the distribution of the imputed values before accepting the imputation model.

Obviously, most of the imputed values for the variable *NB.PFE* are larger than the observed values for this variable. To understand this difference, we examined the dependence between the missing rate and the establishment size. In Table 7.1 we present the percentage of missing units in 10 establishment size classes defined by quantiles and the mean of *NB.PFE* within these quantiles. The missing rates are low up to the sixth establishment size class. Beyond that point the missing rates increase significantly with every class. The average number of further education participants increases steadily with every establishment size class with largest increases in the second half of the table. With these results in mind, it is not surprising that the imputed values for

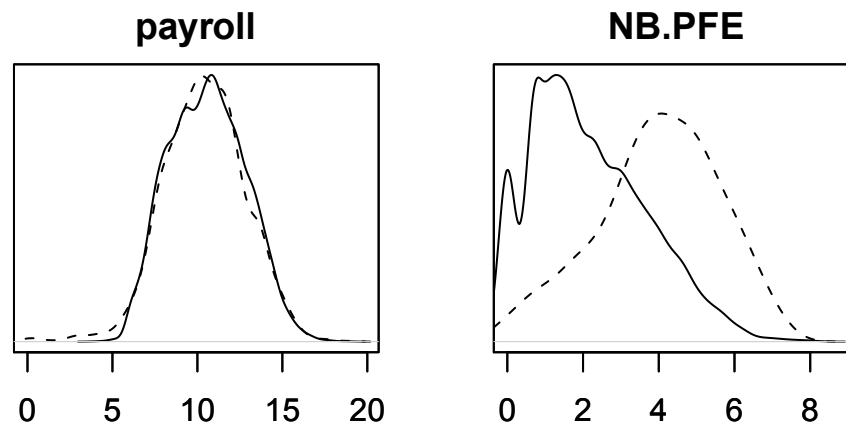


Figure 7.1: Observed (solid line) and imputed (dashed line) data for *payroll* and *number of participants in further education (NB.PFE)*. Both variables are reported on the log-scale.

that variable are often larger than the observed values.

We inspected several continuous variables by comparing the distributions of the observed and imputed values in our dataset and did not find any differences in the distributions that could not be explained by the missingness pattern. However, these comparisons are only meaningful, if enough observations are imputed. Otherwise the distributions between observed data and imputed data might look completely different, only because using kernel density estimation to produce a smooth distribution graph is not appropriate in this context. For this reason we restricted the density comparisons to variables with more than 200 imputed values above zero. For the remaining variables we plotted histograms to check for differences between the observed and imputed values and to detect univariate outliers in the imputed data.

We also investigated if any weighted imputed value for any variable lay above the maximum weighted observed value for that variable. Again, this would not necessarily be problematic, but we did not want to produce any unrealistic influential outliers. However, we did not find any weighted imputed value that was higher than the maximum of its weighted observed counterpart.

For the internal diagnostics, we used three graphics to evaluate the model fit: A Normal Q-Q plot, a plot of the residuals from the regression against the fitted values and a binned residual plot (Gelman and Hill, 2006). The Normal Q-Q plot indicates if the assumption of a normal distribution for the resid-

Table 7.1: Missing rates and means per quantile for *NB.PRE*.

est. size quantile	missing rate in %	$mean(NB.PFE)$ per quantile
1	0.09	1.61
2	0.00	2.49
3	0.57	3.02
4	0.36	4.48
5	0.44	6.09
6	0.37	9.53
7	0.85	15.48
8	1.16	26.44
9	3.18	56.39
10	6.66	194.09

uals is justified by plotting the theoretical quantiles of a normal distribution against the empirical quantiles of the residuals. The residual plot visualizes any unwanted dependencies between the fitted values and the residuals. The binned residual plot plots the average fitted value against the average residual within predefined bins. This is especially helpful for categorical variables since the output of a simple residual plot is difficult to interpret if the outcome is discrete.

Figure 7.2 again provides an example of one model (one of the models for the variable *turnover*) that we did not inspect any further and one model (for the variable *number of participants in further education with college degree (NB.PFE.COL)*), for which we checked the model for necessary adjustments.

For both variables the assumption that the residuals are more or less normally distributed seems to be justified. For the variable *turnover*, the two residual plots further confirm the quality of the model. Only a small amount of residuals fall outside of the grey dotted 95% confidence bands for the residual plot and non of the averaged residuals falls outside the grey 95% confidence bands for the binned residuals. This is different for *NB.PFE.COL*. Although still most of the points are inside the 95% confidence bands, we see a clear relationship between the fitted values and the residuals for the small values and the binned residuals for these small values all fall outside the confidence bands. However, this phenomenon can be explained if we inspect the variable further. Most establishments don't have any participants in further training with college

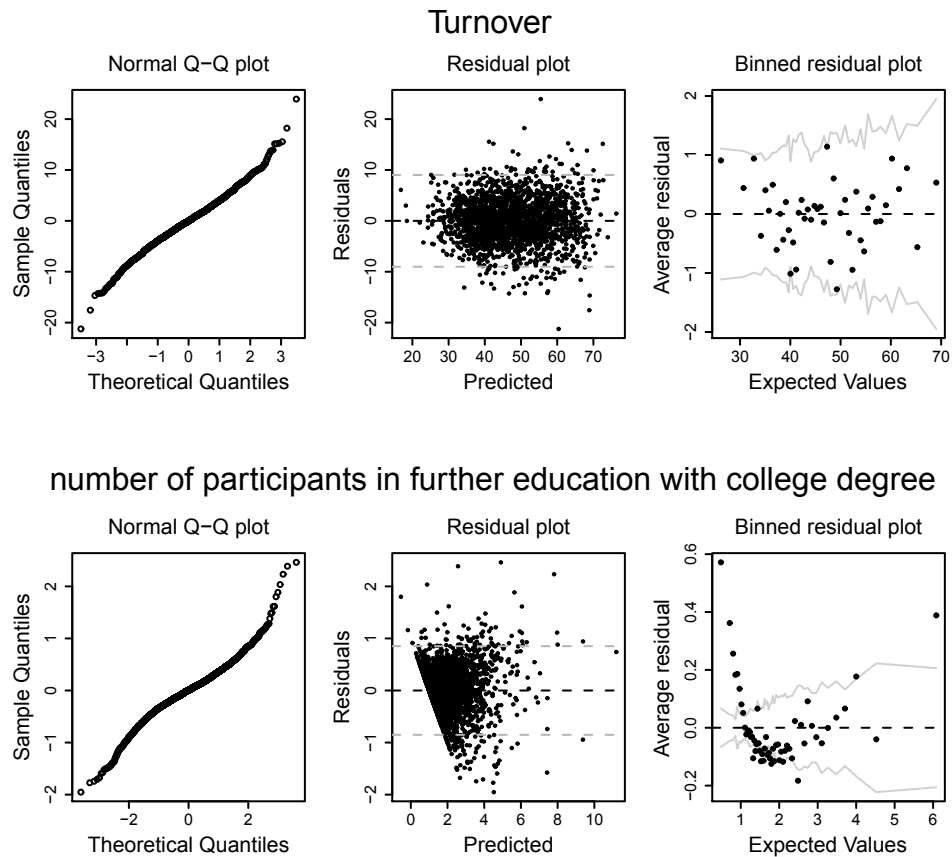


Figure 7.2: Model checks for *turnover* and *number of participants in further education with college degree*.

degree and we fitted the model only to the 3,426 units that reported to have at least one participant. 648 of these units report that they had only 1 participant, leading to a spike at 1 in the original data. Since we simply fit a linear model to the observed data, the almost vertical line in the residual plot is not surprising. It contains all the residuals for all the units with only 1 participant in the original data. The binned residual plot indicates that the small fitted values sometimes severely underestimate the original values. The reason for this again is the fact that the original data is truncated at 1 whereas the fitted values are predictions from a standard linear model that would even allow negative fitted values, since we computed the fitted values before the adjustments for non-negativity described in Section 3.3.3. The consequence is a slight overestimation for the larger fitted values.

We found similar patterns in some other variables that had huge spikes at 1.

We could have tried to model the data with a truncated distribution or we could have applied the semi-continuous approach described in Section 3.3.2 to model the spike at 1 separately, but since we expect that the non-negativity adjustments reduce this effect, we decided to avoid making the already complex modeling task even more difficult.

Missing rates are substantially lower for the categorical variables. Only 59 out of the close to 200 categorical variables in the dataset have missing rates above 1% and we limited our evaluation to these variables. We compared the relative number of responses in each category for the observed and the imputed values and flagged a variable for closer inspection, if the relative number of responses in one imputed category differed more than 20% from the relative number in the observed category. We further limited our search to categories that contained at least 25 units, since small changes in categories with less units would lead to significant changes in the relative differences for these categories. All 15 variables that were flagged by this procedure had a missing rate below 5% and the differences between the imputed and original response rates could be explained by the missingness pattern for all of them. We select one variable here to illustrate the significant differences between observed and imputed values that can arise from a missingness pattern that is definitely not missing completely at random. The variable under consideration asks for the expectations about the investment in 2007 compared to 2006. Table 7.2 provides some summary statistics for this variable. We find a substantial difference for the second and the third category, if we simply compare the observed response rates (column 1) with the imputed response rates (column 2). But the missing rate is only 0.2% for this variable for units with investments in 2006 but soars to 10.5% for units without investments in 2006. Thus, the response rates across categories for the imputed values will be influenced by the expectations for those units that had no investments in 2006 (column 4) even though only 12.9% of the participants who planned investments for 2007 reported no investments in 2006. These response rates differ completely from the response rates for units that reported investments in 2006 (column 3). Thus, the percentage of establishments that expect an increase in investments is significantly larger in the imputed data than it is in the original data.

For categorical data the Normal Q-Q plot is not appropriate as an internal diagnostic tool and the residual plot is difficult to interpret if the outcome is discrete. Therefore, we only examined the binned residual plots for the 59

Table 7.2: Expectations for the investments in 2007 (response rates in % for each category).

category	obs. data	imp. data	obs. units with investment 2006	obs. units without investment 2006
will stay the same	36.57	37.96	41.33	0.59
increase expected	38.79	57.66	30.74	99.41
decrease expected	20.33	0.73	23.05	0.00
don't know yet	4.31	3.65	4.88	0.00

categorical variables with missing rates above 1%. All plots indicate a good model fit. We move all graphics to the Appendix A.2 for brevity.

To check for possible problems with the iterative imputation procedure, we stored the mean for several continuous variables after every imputation round. We did not find any inherent trend for the imputed means for any of the variables. Of course, this is no guarantee for convergence. A possible strategy to measure the convergence of the algorithm is implemented in the new `mi` package by Su *et al.* (2009) following the ideas in Gelman *et al.* (2004). If different imputation chains are run to generate the m imputations, convergence can be monitored by calculating the variance of a given estimate of interest ψ (Su *et al.* (2009) use the mean and the standard deviation of each variable) within and between different imputation chains. Let ψ_{ij} denote the estimate obtained at iteration i , $i = 1, \dots, T$ in chain j , $j = 1, \dots, m$. The between-sequence variance B and the average within-sequence variance W can be calculated as:

$$B = \frac{T}{m-1} \sum_{j=1}^m (\bar{\psi}_{.j} - \bar{\psi}_{..})^2, \quad \text{where} \quad \bar{\psi}_{.j} = \frac{1}{T} \sum_{i=1}^T \psi_{ij}, \quad \bar{\psi}_{..} = \frac{1}{m} \sum_{j=1}^m \bar{\psi}_{.j}$$

$$W = \frac{1}{m} \sum_{j=1}^m s_j^2, \quad \text{where} \quad s_j^2 = \frac{1}{T-1} \sum_{i=1}^T (\psi_{ij} - \bar{\psi}_{.j})^2.$$

Gelman *et al.* (2004), p.297 suggest that convergence can be assumed if

$$\hat{R} = \sqrt{\frac{(1-1/T)W + 1/T * B}{W}} \quad (7.8)$$

is less than 1.1. We did not monitor this measure in our imputation routines.

7.4 Generating synthetic datasets from the multiply imputed IAB Establishment Panel

Once all missing values in the original data have been imputed, we can begin with the actual synthesis. The first and crucial step in the synthesis process is to decide which variables need to be synthesized and whether it is necessary to synthesize all records in the dataset. In general agencies can decide if they only want to select key variables for synthesis or if they also want to synthesize some of the sensitive variables. Key variables are those variables that could be used for re-identification purposes, i.e. variables for which the intruder knows the true values for some target records from external databases like business or credit information databases. Sensitive variables are all those variables that contain information that a survey respondent would not be willing to provide to the general public. In theory there is often no need to synthesize sensitive variables that are not considered key variables. If all key variables are sufficiently protected it will not be possible to link any record in the dataset to a specific respondent. Synthesizing sensitive variables is a conservative approach that might be justified since the amount of data available in external databases might increase over time, so records that are considered safe now might be at risk later. It also helps to convince survey respondents that their information is sufficiently protected.

In our project we decided to synthesize a combination of both variable types. Obviously key variables like *establishment size*, *region* and *industry code* need to be protected, since a combination of the three variables would enable the intruder to identify most of the larger establishments, but we also synthesized the most sensitive variables in the dataset like *turnover* or *amount of subsidies received from the government*. Almost all numerical and some of the categorical variables are synthesized.

In many datasets it is sufficient to alter only the subset of records that are actually at risk. These records can be found by cross tabulating the key variables. Only those records in cross tabulation cells with cell counts below an agency defined threshold might need protection. The selective multiple imputation of keys (SMIKE, Liu and Little (2002)) approach aims in that direction. In our application it might have been sufficient to synthesize values only for

the larger establishments since the sampling uncertainty and the similarities of the small establishments will make re-identification very difficult. Besides, arguably intruders will only be interested to identify some larger establishments. However, we decided to synthesize all records since, given the large amount of information contained in the dataset (close to 300 variables), all records are sampling uniques arguably even population uniques. Of course only a few variables in the dataset can be considered key variables, but once the dataset is released, a survey respondent might try to identify herself in the released dataset. Since the respondent knows all the answers she provided, it will be easy for her to find herself in the dataset. If she realizes that her record is included completely unchanged, she will feel that her privacy is at risk, even if an intruder that will not have the same background information will never be able to identify this respondent. To drive down this perceived risk we decided to synthesize all 15,644 records in the dataset.

7.4.1 The synthesis task

For the synthesis we use the sequential regression multivariate imputation approach (SRMI, Raghunathan *et al.* (2001)) with linear regression models for the continuous variables and logit models for the binary variables (See Section 3.3 for details on how to adjust these methods for skip patterns and logical constraints). Since we always replace all records with synthetic values and leave some of the variables unchanged, we do not have to iterate between the imputations like in standard SRMI for missing values. For illustration let Y_1, \dots, Y_3 be some sensitive variables in a dataset selected for replacement and let X be all variables that remain unchanged in the released dataset. To generate valid synthetic datasets, we need to draw replacement values from the joint distribution $f(Y_1, Y_2, Y_3|X)$. Note that we can write this distribution as $f(Y_1, Y_2, Y_3|X) = f(Y_1|X)f(Y_2|Y_1, X)f(Y_3|Y_1, Y_2, X)$.

Thus, we start our synthesis by drawing new values for Y_1 from an imputation model that only conditions on the unchanged variables X . Next, we built a model for Y_2 conditioning on the originally observed values of X and Y_1 . However, we use the imputed values of Y_1 when drawing new values for Y_2 . Finally, we built a model for Y_3 conditioning on all variables in the original data. New values for Y_3 are drawn using the imputed values of Y_1 and Y_2 . This approach, originally proposed in the missing data context for so called monotone

missingness pattern (Rubin (1987), Chapter 5.4), speeds up the imputation process, because we do not need to iterate before and between the imputations to guarantee convergence to the joint distribution and independence of the draws respectively.

Since all records are replaced with imputed values in our synthesis, developing good models is essential. All variables that don't contain any structural missings are used as predictors in the imputation models in hopes of reducing problems from uncongeniality (Meng, 1994). For the synthesis we use several imputation models for every variable whenever possible. Different models are defined for West and East Germany and for different establishment size classes defined by quantiles. Depending on the number of observations that could be used for the modeling, we define up to 8 different regression models. We do not use the multinomial logit model for the synthesis of the polytomous variables since we already experienced problems with this approach when imputing the missing values in the dataset. For the synthesis we do not want to limit our imputation models to some 30 explanatory variables. Furthermore, we also have to synthesize variables with a large number of categories like region (16 categories) and industry code (41 categories). The multinomial model would hardly ever converge for these variables.

The standard approach for a model based imputation of categorical variables with many categories is the multinomial/Dirichlet approach (see for example Abowd *et al.* (2006)). The disadvantage of this approach is that covariates can not be incorporated in the model directly. In general, a different model is fit for a large number of subcategories of the data defined by cross-classifying some of the covariates to preserve the conditional distributions in the defined classes. This approach is impractical if the number of observations in a survey is low, because the number of observations will be too low to define suitable models in every subclass for which the marginal distribution should be preserved. For this reason we follow a different strategy when synthesizing the categorical variables in our dataset. We generate synthetic values using CART models as suggested by Reiter (2005d).

CART models are a flexible tool for estimating the conditional distribution of a univariate outcome given multivariate predictors. Essentially, the CART model partitions the predictor space so that subsets of units formed by the partitions have relatively homogeneous outcomes. The partitions are found by recursive binary splits of the predictors. The series of splits can be effectively

represented by a tree structure, with leaves corresponding to the subsets of units.

CART models also can be used to generate partially synthetic data (Reiter, 2005d). To illustrate the approach, let us assume that we only want to synthesize three categorical variables: region, industry code, and legal form. To generate synthetic datasets for these three variables we proceed as follows. Using the original data D_{obs} we fit a tree of region on all other variables that don't contain any structural missings except industry code and legal form.³ Label this tree $\mathcal{Y}_{(R)}$. We require a minimum of five records in each leaf of the tree and do not prune it; see Reiter (2005d) for discussion of pruning and minimum leaf size. Let L_{Rw} be the w th leaf in $\mathcal{Y}_{(R)}$, and let $Y_{(R)}^{L_{Rw}}$ be the $n_{L_{Rw}}$ values of $Y_{(R)}$ in leaf L_{Rw} . In each L_{Rw} in the tree, we generate a new set of values by drawing from $Y_{(R)}^{L_{Rw}}$ using the Bayesian bootstrap (Rubin, 1981). These sampled values are the replacement imputations for the $n_{L_{Rw}}$ units that belong to L_{Rw} . Repeating the Bayesian bootstrap in each leaf of the region tree results in the i th set of synthetic regions, $Y_{(R)\text{rep},i}$.

Imputations are next made for the industry code. Using D_{obs} , we fit the tree, $\mathcal{Y}_{(I)}$, with all variables except legal form as predictors. To maintain consistency with $Y_{(R)\text{rep},i}$, units' leaves in $\mathcal{Y}_{(I)}$ are located using $Y_{(R)\text{rep},i}$. Occasionally, some units may have combinations of values that do not belong to one of the leaves of $\mathcal{Y}_{(I)}$. For these units, we search up the tree until we find a node that contains the combination, then treat that node as if it were the unit's leaf. Once each unit's leaf is located, values of $Y_{(I)\text{rep},i}$ are generated using the Bayesian bootstrap. Imputing legal form follows the same process: we fit the tree $\mathcal{Y}_{(L)}$ using all variables that don't contain any structural missings as predictors, place each unit in the leaves of $\mathcal{Y}_{(L)}$ based on their synthesized values of region and industry code, and sample new legal forms using the Bayesian bootstrap. We generate $r = 5$ datasets for every imputed dataset, i.e. $m * r = 25$ synthetic datasets will be released. Reiter (2008b) elaborates on the number of imputations on stage one and two when using multiple imputation for nonresponse and disclosure control simultaneously. He suggests to set $m > r$, especially if the fraction of missing information is large, to reduce variance from estimating

³To improve the data quality we actually grow several trees for different subsets of the data. The subsets are defined by West and East Germany and by up to 25 different establishment size classes defined by quantiles. To simplify the notation, we illustrate the approach assuming that only one tree is fit for each variable.

missing values. But this approach will increase the risk of negative variance estimates since \bar{b}_M will increase relative to B_M .

In our dataset only 12 variables (out of more than 300) have missing rates above 5%. On the other hand, we always synthesize 100% of the records. In his simulations Reiter (2008b) does not find a significant reduction in variance with increasing m compared to r for 100% synthesis paired with low missing rates. On the other hand, the risk of negative variance estimates increases significantly. From these results, we conclude that it would be better to set $m = r$ in our case.

7.4.2 Measuring the data utility

We evaluate the data utility of the generated datasets by comparing analytic results achieved with the original (fully imputed) data⁴ with results from the synthetic data. To provide realistic analyses, we use two regressions suggested by colleagues at the IAB, who regularly use the panel for applied analyses. The probit regression displayed in Tables 7.3 and 7.4 is adapted from a regression originally based on a different wave of the establishment panel. The dependent variable indicates if an establishment employs part-time employees. The 19 explanatory variables include among others dummies for the establishment size, whether the establishment expects changes in the number of employees, and information on the personnel structure. Since there are still differences within Germany, the results are computed for West Germany (Table 7.3) and East Germany (Table 7.4) separately.

Both regressions clearly demonstrate the good data quality. All point estimates from the synthetic data are close to the point estimates from the original data and the confidence interval overlap (See Section 5.2) is higher than 90% for most estimates with an average of 90% for West Germany and 93% for East Germany. We also report the z-scores for all regressions, because some researchers are concerned that synthetic datasets will provide valid results for the significant variables, but might provide less accurate results for variables with lower z-scores. From the results it is obvious that this is not true. We also note that the z-scores from the synthetic data are very close to the z-scores from the original data. This is an important result, since model selections are

⁴For convenience, we will refer to the dataset with all missing values multiply imputed as the original data from here on.

Table 7.3: Regression results from a probit regression of *part time-employees (yes/no)* on 19 explanatory variables in West Germany. For the CI length ratio the CI length of the original datasets is in the denominator.

	original data	synth. data	CI overlap	z-score org.	z-score syn	CI length ratio
Intercept	-0.809	-0.752	0.87	-7.23	-6.85	0.99
5-10 employees	0.443	0.437	0.97	8.52	7.99	1.06
10-20 employees	0.658	0.636	0.90	11.03	10.88	0.98
20-50 employees	0.797	0.785	0.95	13.02	12.36	1.04
100-200 employees	0.892	0.908	0.96	9.23	9.48	0.99
200-500 employees	1.131	1.125	0.99	9.99	9.87	1.01
>500 employees	1.668	1.641	0.97	8.22	8.33	0.97
growth in employment exp.	0.010	0.006	0.98	0.18	0.12	0.99
decrease in emp. expected	0.087	0.100	0.96	1.11	1.27	1.00
share of female workers	1.449	1.366	0.73	17.63	18.71	0.89
sh. of emp. with uni. degree	0.319	0.368	0.91	2.18	2.59	0.97
sh. of low qualified workers	1.123	1.148	0.93	12.17	11.87	1.05
sh. of temporary employees	-0.327	-0.138	0.75	-1.74	-0.71	1.05
share of agency workers	-0.746	-0.856	0.88	-3.09	-4.24	0.84
empl. in the last 6 mths	0.394	0.369	0.87	8.33	7.82	1.00
dismissal in the last 6 mths	0.294	0.279	0.92	6.38	6.03	1.00
foreign ownership	-0.113	-0.117	0.99	-1.33	-1.38	0.99
good/very good profitability	0.029	0.033	0.98	0.72	0.82	0.99
salary above coll. wage agr.	0.020	0.031	0.95	0.35	0.54	0.99
collective wage agreement	0.016	0.007	0.95	0.31	0.13	0.97

often based on significance levels. The last column reports the 95% confidence interval length ratio with the confidence interval length of the original data in the denominator. Since the multiple imputation procedure for generating synthetic datasets correctly reflects the uncertainty in the imputation models, it can happen that the confidence intervals from the synthetic datasets are much wider and thus less efficient than the confidence intervals from the original data. We find that only for the variable *share of low qualified workers* in Table 7.4 the confidence interval length is increased by 19%. For all other estimands the intervals are never increased more than 7%.

The second regression is an ordered probit regression with the expected employment trend in three categories (increase, no change, decrease) as the de-

Table 7.4: Regression results from a probit regression of *part time-employees (yes/no)* on 19 explanatory variables in East Germany. For the CI length ratio the CI length of the original datasets is in the denominator.

	original data	synth. data	CI over- lap	z- score org.	z- score syn	CI length ratio
Intercept	-0.712	-0.742	0.93	-6.42	-7.21	0.93
5-10 employees	0.266	0.257	0.96	4.81	4.53	1.03
10-20 employees	0.416	0.399	0.93	6.94	6.76	0.99
20-50 employees	0.542	0.532	0.96	9.18	8.72	1.04
100-200 employees	0.757	0.808	0.86	8.02	8.47	1.01
200-500 employees	0.971	1.013	0.91	8.25	8.57	1.00
>500 employees	1.401	1.422	0.98	5.69	5.66	1.02
growth in employment exp.	-0.041	-0.040	1.00	-0.73	-0.73	1.00
decrease in emp. expected	0.035	0.040	0.98	0.44	0.50	1.00
share of female workers	1.006	1.041	0.88	12.63	14.93	0.88
sh. of emp. with uni. degree	0.221	0.197	0.95	1.86	1.76	0.95
sh. of low qualified workers	0.976	1.042	0.87	8.44	7.84	1.19
sh. of temporary employees	-0.049	0.049	0.84	-0.31	0.34	0.91
share of agency workers	-0.176	-0.232	0.94	-0.73	-1.08	0.89
empl. in the last 6 mths	0.230	0.210	0.89	4.95	4.55	1.00
dismissal in the last 6 mths	0.301	0.295	0.97	6.43	6.35	0.99
foreign ownership	-0.176	-0.176	1.00	-1.83	-1.84	1.00
good/very good profitability	0.097	0.097	1.00	2.35	2.37	1.00
salary above coll. wage agr.	0.080	0.086	0.98	1.04	1.10	1.01
collective wage agreement	0.097	0.069	0.86	1.87	1.36	0.98

pendent variable. In the regression, we use 39 explanatory variables and the industry dummies as covariates. Again the analysis is computed for West Germany and East Germany separately. Figure 7.3 contains a plot of the original point estimates against the synthetic point estimates and a boxplot for the confidence interval overlap and the confidence interval length ratio. All graphs are based on the 78 estimates from the two regressions. Most of the point estimates in the first graph are close to the 45 degree line indicating that the point estimates from the synthetic data are very close to the point estimates from the original data. But even if the point estimates differ, we find that the data utility measured by the confidence interval overlap is high. The measure never drops below 61% and the median overlap is 92.7%. Thus, even though some

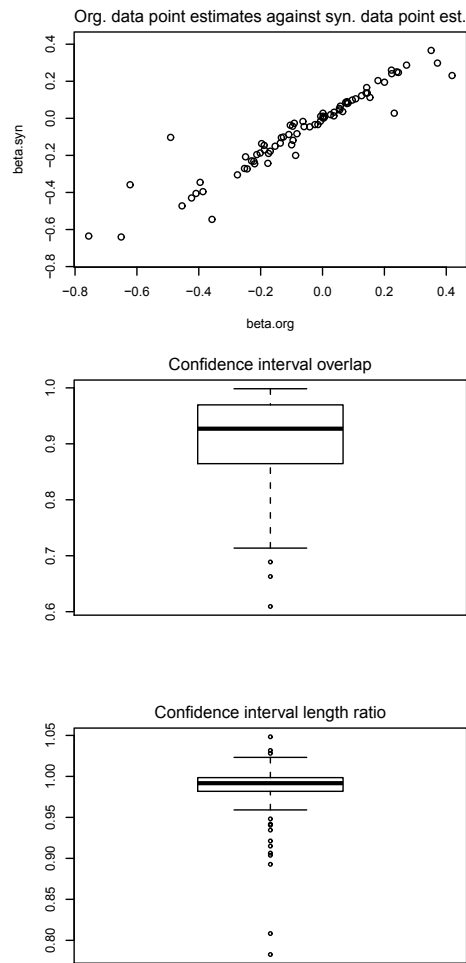


Figure 7.3: Ordered probit regression of *expected employment trend* on 39 explanatory variables and industry dummies.

estimates are a little off the 45 degree line, the results are close to the original results since these coefficients are estimated with a high standard error. The boxplot of the confidence interval length ratio indicates that we do not lose much efficiency by using the synthetic data instead of the original data. The confidence interval never increases by more than 5% compared to the original data.

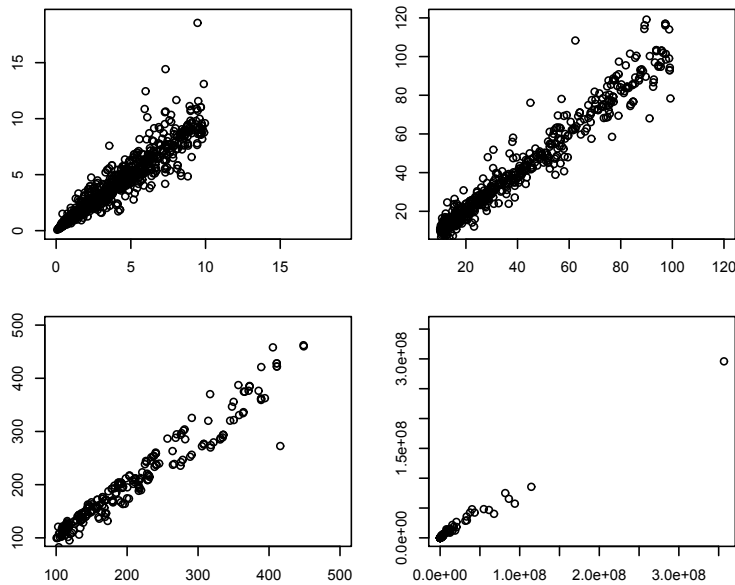


Figure 7.4: Original point estimates against synthetic point estimates for the overall mean and the means in subgroups defined by establishment size class, industry code and region.

Not all users of the data will be interested in multivariate regression analysis. For this reason we also included an evaluation of the data utility for descriptive statistics. For this, we compare the unweighted⁵ overall mean and the unweighted mean in different subgroups for all continuous variables in the dataset. The subgroups are defined by establishment size (10 categories, defined by quantiles), industry code (17 categories) and region (16 categories). We do not investigate any cross classifications since the cell sizes would be too small to obtain meaningful results. We also limit our evaluation to cells with at least 200 observations above zero for the same reason. This leads to a final number of 2,170 estimates. Figure 7.4 again presents the plots of the estimates from the original fully imputed datasets against the synthetic estimates. For readability the plots are divided in four parts depending on the original value of the mean ($[0; 10]$, $(10; 100]$, $(100; 1000]$, $(1000; \infty)$). We find that most of the synthetic estimates are close to their original counterparts. Only few estimates differ significantly from the original values. Figure 7.5 contains box plots for the confidence interval overlap. The results for each stratifying variable and

⁵We use the unweighted mean, because the weights were still under development when we performed this evaluation.

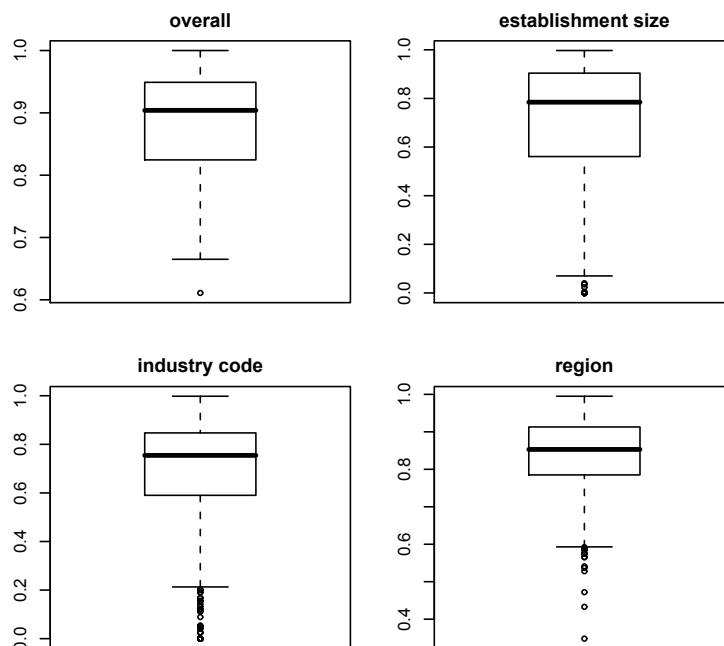


Figure 7.5: Box plots of CI overlaps for all continuous variables for the overall mean and the means in all subgroups defined by different stratifying variables.

the overall mean are reported separately. The median overlap is always higher than 75% indicating a very good overall quality of the data. Not surprisingly the overall mean (based on 92 estimates) provides the best results, with an overlap that never falls below 60% and a median overlap above 90%. The results for the means in different establishment size classes (based on 552 estimates) and the means for different industry codes (based on 720 estimates) are good for most of the estimates with a median overlap of 78.5% and 75.5% respectively, but for a small number of estimates (6.2% and 3.3% of the estimands respectively) the overlap is actually zero. The results are better for the region. The median overlap (based on 806 estimates) is 85.3%, the overlap never falls below 34% and only 3 estimates have an overlap below 50%.

7.4.3 Caveats in the use of synthetic datasets

Despite these very promising results, it would be overly optimistic to assume that synthetic datasets will provide results of similar quality for any potential analysis. It is crucial that the potential user of the data knows which analysis might provide valid results and for which analysis she might have to apply for direct access to the data at the research data center. To enable the user to make these decisions it is very important that additional information about the imputation models is released in combination with the synthetic data. For example the IAB could release information about which explanatory variables were used in the imputation models for each variable.

To give an example for which the synthetic data would not give valid results, we run a probit regression with the same explanatory variables as in Table 7.3 but we replace the dependent variable with an employment trend variable that equals 1 if the number of employees covered by social security increases between 2006 and 2007 and is 0 otherwise. We don't claim that this is a useful applied analysis, it only helps to illustrate that users should be careful when fitting models with dependent variables derived from two or more variables.

Table 7.5 provides the results for this regression and it is obvious that they are by no means close to the results given above in terms of data quality. 6 of the 20 estimates actually have no confidence interval overlap at all and the point estimates and z-scores often differ substantially from the original estimates. So the question arises, what is the reason for the poor performance of the synthetic datasets for this regression? To understand the problem, we first compare the number of employees covered by social security 2006 and 2007 between the original data and the synthetic data. Figure 7.6 presents QQ-plots of the original values against the synthetic values. The first two graphs present the plots for the two variables and the last plot depicts the QQ-plot for the difference in the number of employees between 2006 and 2007. We find that the synthesis model did a very good job in capturing the distribution of the variables for 2006 and 2007, the quantiles are more or less identical. The distribution of the difference between the number of employees covered by social security between 2006 and 2007 is also well preserved. If we would run a simple linear regression with the same covariates but with the difference in employment as the dependent variable, the average confidence interval overlap would be 75%, a significant improve compared to 42% for the results in Table 7.5.

Table 7.5: Regression results from a probit regression of *employment trend (increase/no increase)* on 19 explanatory variables in West Germany. For the CI length ratio the CI length of the original datasets is in the denominator.

	original data	synth. data	CI over- lap	z- score org.	z- score syn	CI length ratio
Intercept	-1.396	-0.978	0.05	-11.99	-9.28	0.92
5-10 employees	0.130	0.354	0.00	2.61	7.75	0.92
10-20 employees	0.316	0.495	0.05	6.19	11.19	0.87
20-50 employees	0.355	0.541	0.05	7.33	10.93	1.06
100-200 employees	0.366	0.351	0.94	5.69	6.09	0.91
200-500 employees	0.475	0.347	0.48	7.29	5.80	0.92
>500 employees	0.375	0.472	0.66	5.06	6.58	0.99
growth in employment exp.	0.374	0.148	0.00	9.29	3.59	1.05
decrease in emp. expected	-0.376	-0.020	0.00	-6.16	-0.38	0.86
share of female workers	-0.140	-0.054	0.67	-2.09	-0.84	1.00
sh. of emp. with uni. degree	0.229	0.199	0.91	1.94	2.05	0.83
sh. of low qualified workers	-0.043	-0.004	0.84	-0.68	-0.07	0.97
sh. of temporary employees	0.434	0.226	0.62	3.25	1.60	1.07
share of agency workers	0.058	0.013	0.69	0.94	0.08	2.61
empl. in the last 6 mths	0.948	0.368	0.00	24.94	11.60	0.84
dismissal in the last 6 mths	-0.172	-0.030	0.00	-4.42	-0.97	0.81
foreign ownership	-0.165	-0.113	0.79	-2.60	-1.90	0.98
good/very good profitability	0.248	0.100	0.00	7.69	3.35	0.93
salary above coll. wage agr.	0.039	0.033	0.96	0.87	0.81	0.91
collective wage agreement	0.003	0.063	0.62	0.06	1.72	0.85

The actual problem stems from the fact that there is not much variation between the employment numbers 2006 and 2007. In the original dataset 5,376 of the 15,644 establishments report no change in employment numbers and more than 90% of the establishments report change rates of $\pm 5\%$. It can easily happen that in the original data, an establishment reported an increase from 30 to 31 employees, but in the synthetic data this establishment might have imputed values of 30 in both years or maybe 29 in the second year. Thus, the actual number is estimated very well and even the predicted difference is very close, but this record will change from an establishment with positive employment trend to an establishment with no change or even negative employment trend. The opposite is likely to occur as well: A record with a small negative employ-

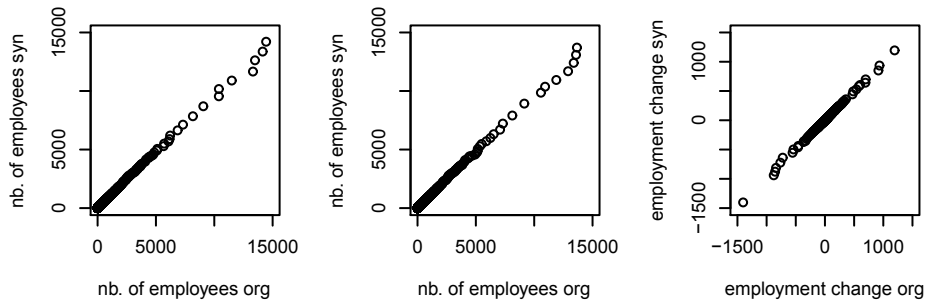


Figure 7.6: QQ-plots for the *number of employees covered by social security* 2006 and 2007 and the employment trend between the two years.

ment trend might end up with a positive employment trend. If this happens for many records, which is to be expected, since changes are very small for most records in the original dataset, the binary variable for employment trend will assign ones to a completely different subset of records. It is not surprising that results from the synthetic data will be different from the results in the original data in this case. It is important that users are made aware of this problem that is likely to occur, if the user derives her variable of interest from two or more variables in the dataset and small changes in the underlying variables can have huge impacts on the derived variable. On a side note, this problem is not limited to multiply imputed synthetic datasets. In fact, most if not all standard perturbative SDC methods like swapping, adding noise or micro aggregation will lead to similar problems.

7.4.4 Assessing the disclosure risk

It is unlikely that an intruder has detailed information about who participated in the survey, thus using the actual true data from the survey for the disclosure risk calculations is an unrealistic conservative scenario. For this reason we apply the disclosure risk measures described in Section 6.3.2 that account for the additional uncertainty from sampling.

To obtain a set of target records for which the intruder has some knowledge from external databases that she uses to identify units in the survey, we sample new records from the sampling frame of the survey, the German Social Security Data (GSSD). We sample from this frame, using the same sampling design as for the IAB Establishment Panel: Stratification by establishment size, region and industry code. Merging the stratification matrix from the panel to the stratification matrix of the GSSD reveals that there are 14 stratification cells with positive entries in the panel matrix that are empty in the GSSD matrix. This is a result of the fact that some establishments don't provide answers only for their own entity. They erroneously provide the numbers for the whole concern they belong to instead. By doing so, the establishment might move to another stratification cell that is empty in the original sampling frame. We remove these 14 entries from the stratification matrix of the survey. For the same reason it is possible that some panel cells contain more records than the corresponding GSSD cell. If this happens, we sample all records in this GSSD cell. Overall this leads to a reduction from 15,644 establishments in the original data to 15,624 records in the target sample.

Merging the GSSD and the IAB Establishment Panel using the establishment identification number, we find that 1,360 units from the panel are not included in the GSSD.⁶ As a consequence, these records will never appear in the target sample. Since more than 93% of these records are establishments with less than 100 employees, only 4 of them have between 1,000-5,000 employees and non has more than 5,000 employees, we are not concerned that we underestimate the disclosure risk by excluding these records from the target sample.

We find that 917 records from the target sample are also included in the original sample. Table 7.6 displays the percentage of records from the original dataset

⁶There are several possible reasons for this, e.g. re-organization of the firm leading to new establishment identification numbers, coding errors, or delays in the notifications for an establishment in the GSSD.

Table 7.6: Probabilities to be included in the target sample and in the original sample depending on establishment size.

establishment size class	probability(%)
1-4 employees	0.91
5-9 employees	1.62
10-19 employees	2.87
20-49 employees	4.10
50-99 employees	6.55
100-199 employees	11.39
200-499 employees	16.69
500-999 employees	20.48
1000-4999 employees	31.89
≥ 5000 employees	39.39

that are also included in the target sample for different establishment size classes. As expected, this probability increases with the establishment size. For establishments with less than 100 employees the probability is always less than 10% whereas large establishments with more than 5,000 employees are included in both samples with a probability close to 40%.

For the disclosure scenario we assume, the intruder has information on region, industry code (in 17 categories) and establishment size (measured by the number of employees covered by social security) for her target records and uses this information to identify units in the survey. We further assume that she would consider any record in the synthetic datasets a potential match for a specific target record, if it fulfills two criteria: First, the record's synthetic industry code and region exactly matches the target's true industry code and region. Second, the record's synthetic number of employees lies within a defined interval around the target's number of employees. To define these intervals, we divide the number of employees by the 10 stratification classes for establishment size and calculate the standard deviation within each size class. The interval is $t_e \pm \sqrt{sd_s}$, where t_e is the target's true value and sd_s is standard deviation of the size class in which the true value falls. We investigated several other intervals, e.g. using the standard deviation directly or defining the intervals by 10-20 establishment size classes as we did in the example in Section 6.4.2 instead of using the stratification classes. However, we found that the criteria above led to the highest risk of disclosure.

7.4.4.1 Log-linear modeling to estimate the number of matches in the population

In general, the intruder will not know the number of records F_t that fulfill the matching criteria in the population to estimate the matching probabilities given in (6.7). One way to estimate the population counts from the released samples was suggested by Elamir and Skinner (2006). We apply this approach to our data assuming that the population counts follow an all-two-way-interactions log-linear model. To simplify the computation, we use the original sample to fit the log-linear model instead of fitting a log-linear model to each synthetic dataset separately. Arguably, this will slightly increase the estimated risk, but we don't expect much difference in the results.

To fit the log-linear model, we need to cross tabulate the three matching dimensions region, industry code and establishment size in the sample. To obtain the correct number of establishment size matches, we need to identify all records that fulfill the establishment size match criterium in the survey sample for each integer value of establishment size in the target sample. This leads to a 16x17x1102 dimensional table to which we fit an all-two-way-interactions log-linear model. To calculate \hat{F}_t , we need the sampling probabilities for each entry in this table. We obtain these probabilities by dividing the stratification matrix from the original sample by the stratification matrix from the GSSD. We assign the same probability to all establishment size values that fall into the same stratification cell. Again, an intruder will not know the exact sampling probabilities because she can only estimate the stratification matrix of the original sample from the synthetic samples, but arguably it is possible to obtain information about the number of establishments in Germany by region times industry times establishment size class. Since the stratification matrix from the synthetic samples will not differ very much from the matrix of the original sample, the estimated sampling probabilities might be reasonably close to the true sampling estimates. In any case, using the true sampling probabilities provides an upper bound for the disclosure risk.

Since we can actually compute the true F_t from the GSSD, we are able to evaluate, how well we can estimate the true population counts with the log-linear modeling approach. In Table 7.7 and Figure 7.7 we compare the estimated \hat{F}_t with the true F_t . In Table 7.7 we compute the average \hat{F}_t and F_t for the target records in the 10 establishment size stratification classes. The average

Table 7.7: Average F_t and \hat{F}_t for different establishment size classes.

establishment size class	$mean(\hat{F}_t)$	$mean(F_t)$
1-4 employees	6467.66	6685.90
5-9 employees	1661.49	1737.89
10-19 employees	408.78	440.85
20-49 employees	161.98	179.01
50-99 employees	47.07	52.60
100-199 employees	17.91	22.89
200-499 employees	8.06	9.23
500-999 employees	2.17	2.88
1000-4999 employees	1.51	2.03
≥ 5000 employees	1.00	1.11

estimated population count slightly underestimates the true counts but nevertheless is always very close to the average true population count. In Figure 7.7 we plot \hat{F}_t against F_t for each record in the target sample. The left graph presents the results for all establishments, the right graph is limited to establishments with more than 100 employees. We find that the log-linear modeling approach performs very well even on the record level.

7.4.4.2 Results from the disclosure risk evaluations

To estimate the actual risk of disclosure, we use the summary statistics presented in Section 6.3.1 accounting for the uncertainty from sampling as described in Section 6.3.2. These statistics are presented in Table 7.8. Notice that using \hat{F}_t instead of F_t gives almost similar results. In both cases, we find that the disclosure risk is very low. Overall only about 150 of the 15,624 records in the target sample are matched correctly and the false match rate is 98.8%. We evaluate the disclosure risk in different establishment size classes and find that the percentage of true matches increase with the establishment size, but never exceeds 7%. We also investigate, if the risks increase, if the intruder only matches, when the average match probability exceeds a predefined threshold γ . Table 7.9 lists the false match rate and the number of true matches for different threshold values using F_t (there is almost no difference in the results if we use \hat{F}_t instead). The false match rates continually decrease to almost 80% at $\gamma \leq 0.5$. Further reducing γ leads to no improvements in

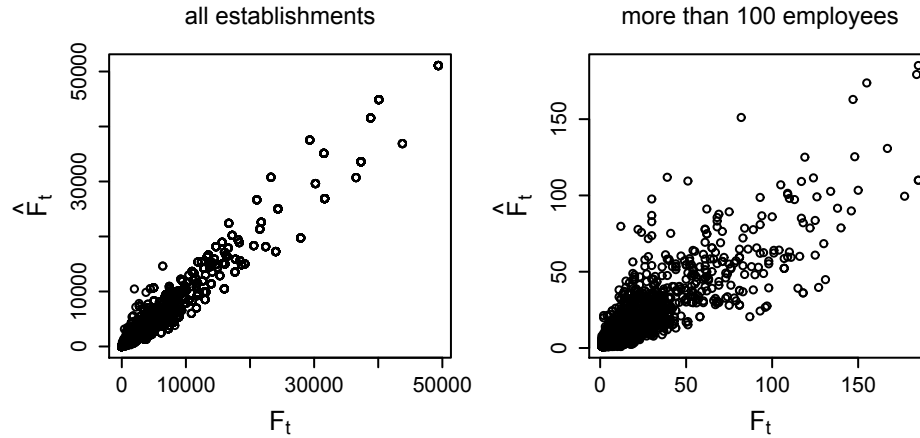


Figure 7.7: Plots of F_t against \hat{F}_t for all establishments and for establishments with more than 100 employees.

terms of the false match rate. Only for $\gamma \leq 0.1$ the rate drops to 66.7%. On the same time, the number of true matches continuously decreases until no true match is found at a threshold of $\gamma = 0$. Since the intruder never knows, which matches actually are true matches, these results indicate that the data seems to be well protected at least under the given assumptions about the information an intruder can gather in her target data.

Table 7.8: Disclosure risk summaries for the synthetic establishment panel wave 2007.

	$mean(\hat{F}_t)$	$mean(F_t)$
expected match risk	162.34	160.92
true match risk	152	150
false match rate (%)	98.75	98.76

7.4.4.3 Disclosure risk for large establishments

Even though the results in the last section indicate a low risk of disclosure, large establishments might still be at risk because these establishments might be identifiable by matching on establishment size alone. Since a potential intruder will know that region and industry code have been synthesized, she might match only on establishment size for large establishments and ignore

Table 7.9: False match rate and true match risk for different levels of γ .

γ	false match rate	true match risk
1	98.76	150
0.9	94.42	97
0.8	91.47	59
0.7	88.72	38
0.6	84.57	27
0.5	81.91	17
0.4	84.62	8
0.3	82.14	5
0.2	85.71	2
0.1	66.67	1
0.0	-	0

that region and industry code are different between the target record and the found match in the synthetic data.

To quantify the risk from this approach, we evaluate two disclosure risk scenarios. In the first scenario, the intruder ranks all synthetic datasets by establishment size and considers the mode of the ranks for one unit across the synthetic datasets as the true rank of this unit. She then links that unit to the unit with the same rank in her target dataset. The second scenario assumes that the intruder performs a simple nearest neighbor matching between the records in her target data and the records in the synthetic samples using the establishment size variable.

Since the largest establishments are sampled with probability close to 1, we treat the original sample as the target sample from which the intruder knows the true reported establishment size. This is still conservative, since the reported establishment size might differ from the size reported in other databases, but it is not unlikely that the intruder will get reasonable close estimates of the true establishment size for large establishments in Germany.

Table 7.10 provides the results for the largest 25 establishments. The average match rate in column three is the percentage of times the declared match from the nearest neighbor matching approach actually is the true match across the 25 synthetic datasets. Obviously the largest establishments face a very high risk of disclosure in both scenarios. The mode of the ranks in the synthetic datasets is almost always the same as the rank in the original sample and the

nearest neighbor matching approach will lead to correct matches for most of the datasets. If the intruder would also pick the mode of declared matches as

Table 7.10: Mode of the establishment size rank and average match rate for large establishments.

original rank	mode of synthetic ranks	Average match rate
1	1	0.96
2	2	0.72
3	3	1.00
4	4	1.00
5	5	1.00
6	6	0.88
7	7	0.64
8	8	0.56
9	9	0.44
10	10	0.32
11	11	0.84
12	12	0.56
13	13	0.56
14	14	0.68
15	15	0.76
16	17	0.56
17	18	0.48
18	16	0.00
19	19	0.56
20	20	0.04
21	22	0.44
22	23	0.72
23	21	0.00
24	24	0.40
25	25	0.28

the correct match, she would be right for 21 of the 25 establishments. Clearly, there is a need to further protect the largest establishments in the dataset.

7.4.4.4 Additional protection for the largest establishments in the survey

A simple strategy to better protect large establishments would be, to reduce the quality of the imputation model for establishment size for example by dropping explanatory variables from the imputation model until a predefined criterium of variability between the imputations is met. However, since we would have to drop the variables with the highest explanatory power to significantly increase the variability, important relationships between the variables would not be reflected in the released data leading to uncongeniality problems if the analyst's model differs from the imputation model. It is also not an option to use other SDL techniques since methods like noise addition would have to be applied on a very high level and other methods like data swapping and micro aggregation are well known to have severe negative consequences for data quality in the upper tail of the distribution. We therefore decided to inflate the variance of the beta coefficients in the imputation model instead. Remember that the imputation process always consists of two steps: In the first step, new parameters for the imputation model are drawn from their posterior distribution given the observed data. In the second step, new values for the variable to be imputed are drawn from the posterior predictive distribution given the parameters drawn in step one. For the standard linear model this means that step one consists of drawing new values of σ^2 and β from their posterior distributions. We decided to protect records at risk by inflating the variance of β in the underlying imputation models. We inflate the variance by drawing new values of β from:

$$\beta|\sigma^2 \sim N(\hat{\beta}, \alpha\sigma^2(X'X)^{-1}) \quad (7.9)$$

where α is the variance inflation factor, $\hat{\beta}$ and X are the regression coefficients and the explanatory variables from the underlying imputation model, and σ^2 is the new value of the variance drawn from its posterior distribution. Of course, imputation under this variance inflated model is not proper in Rubin's sense (see Rubin (1987), pp. 118–119), so we conducted a small simulation study to evaluate the impact of different levels of α on the validity of the results from a frequentist perspective. In our simulation, reported in the Appendix A.3, we found almost no impact on coverage rates. Even when synthesizing all records with $\alpha = 1,000$, the coverage rate for the 95% confidence interval never dropped below 90% and was close to the nominal 95% for most of the estimates of interest. The most notable consequence is that we loose efficiency

since the between imputation variance increases linearly with α . But since we only want to replace some records at risk, we are not concerned that this will have huge impacts on data utility.

To apply the variance inflation approach, we need to define which records we consider to be at risk. We define a record to be at risk, if one of the two following conditions is fulfilled:

1. The standard deviation of the establishment size rank across the synthetic datasets for the record is less than 2.
2. The mode of the declared matches in the nearest neighbor matching scenario is the correct match.

The threshold value for the standard deviation of the ranks is chosen somewhat arbitrarily. Defining justifiable threshold rules is an area for future research. To keep the negative impacts of this procedure at a minimum, we developed an iterative replacement algorithm. For a given level of α , all records that fulfill one of the above criteria are replaced by new draws from the variance inflated imputation model. Records that still are at risk after 10 rounds of repeatedly drawing from this model are replaced by draws from a model with the next higher level of α . In our application, we set the levels arbitrarily to $\alpha = (10; 100; 1,000)$. Developing methods to derive useful levels of α is an area of future research. Overall we replace 79 records in our dataset by this procedure. Less than 10 are replaced by draws from imputation models with $\alpha \geq 100$. Evaluating the disclosure risk for large establishments again, we find that the mode of the establishment size rank in the synthetic datasets is equal to the rank in the original data for only 12 of the largest 100 establishments. Since the intruder never knows, if her match is correct and it is also unlikely that the intruder will know the original rank beyond the largest 20 establishments in the survey, the data is well protected for these kind of attacks. For the nearest neighbor matching scenario, we guaranteed that the mode of the declared matches is never the correct match. We also find that no record is identified correctly in more than 5 of the 25 datasets. These results together with the results in Section 7.4.4 and the promising results on data utility in Section 7.4.2 demonstrate that our dataset is ready for release.

Chapter 8

A Two Stage Imputation Procedure to Balance the Risk-Utility-Trade-Off¹

There has been little discussion in the literature on how many multiply-imputed datasets an agency should release. From the perspective of the secondary data analyst, a large number of datasets is desirable. The additional variance introduced by the imputation decreases with the number of released datasets. For example, Reiter (2003) finds nearly a 100% increase in variance of regression coefficients when going from fifty to two partially synthetic datasets. From the perspective of the agency, a small number of datasets is desirable. The information available to ill-intentioned users seeking to identify individuals in the released datasets increases with the number of released datasets. Thus, agencies considering the release of partially synthetic data generally are confronted with a trade off between disclosure risk and data utility.

The empirical investigations presented in Section 8.3 indicate that increasing m results in both higher data utility and higher risk of disclosures. In this chapter, we present an alternative synthesis approach that can maintain high utility while reducing disclosure risks. The basic idea behind this approach is to impute variables that drive the disclosure risk only a few times and other variables many times. This can be accomplished by generating data in two stages, as described by Reiter and Drechsler (2010). In general, two

¹Most of this chapter is taken from Drechsler and Reiter (2009) and Reiter and Drechsler (2010).

stage and one stage approaches require similar amounts of modeling efforts; however, in some settings, the two stage approach can reduce computational burdens associated with generating synthetic data and thereby speed up the process; see (Reiter and Drechsler, 2010) for further discussion of this point. The two stage imputation procedure is applicable to both, partially and fully synthetic datasets. In the next sections, we present the combining rules to obtain valid inferences for both approaches and provide an application of the two stage partially synthetic approach to illustrate the potential benefits of this procedure.

8.1 Inference for synthetic datasets generated in two stages

For a finite population of size N , let $I_l = 1$ if unit l is included in the survey, and $I_l = 0$ otherwise, where $l = 1, \dots, N$. Let $I = (I_1, \dots, I_N)$, and let the sample size $s = \sum I_l$. Let X be the $N \times d$ matrix of sampling design variables, e.g. stratum or cluster indicators or size measures. We assume that X is known approximately for the entire population, for example from census records or the sampling frame(s). Let Y be the $N \times p$ matrix of survey data for the population. Let $Y_{inc} = (Y_{obs}, Y_{mis})$ be the $s \times p$ sub-matrix of Y for all units with $I_l = 1$, where Y_{obs} is the portion of Y_{inc} that is observed, and Y_{mis} is the portion of Y_{inc} that is missing due to nonresponse. Let R be an $N \times p$ matrix of indicators such that $R_{lk} = 1$ if the response for unit l to item k is recorded, and $R_{lk} = 0$ otherwise. The observed data is thus $D_{obs} = (X, Y_{obs}, I, R)$.

8.1.1 Fully synthetic data

Let Y_a be the values simulated in stage 1, and let Y_b be the values simulated in stage 2. The agency seeks to release fewer replications of Y_a than of Y_b , yet do so in a way that enables the analyst of the data to obtain valid inferences with standard complete data methods. To do so, the agency generates synthetic datasets in a three-step process. First, the agency fills in the unobserved values of Y_a by drawing values from $f(Y_a | D_{obs})$, creating a partially completed population. This is repeated independently m times to obtain $Y_a^{(i)}$, for $i = 1, \dots, m$. Second, in each partially completed population defined by nest i , the agency

generates the unobserved values of Y_b by drawing from $f(Y_b | D_{obs}, Y_a^{(i)})$, thus completing the rest of the population values. This is repeated independently r times for each nest to obtain $Y_b^{(i,j)}$ for $i = 1, \dots, m$ and $j = 1, \dots, r$. The result is $M = mr$ completed populations, $P^{(i,j)} = (D_{obs}, Y_a^{(i)}, Y_b^{(i,j)})$, where $i = 1, \dots, m$ and $j = 1, \dots, r$. Third, the agency takes a simple random sample of size n_{syn} from each completed population $P^{(i,j)}$ to obtain $D^{(i,j)}$. These M samples, $D_{syn} = \{D^{(i,j)} : i = 1, \dots, m; j = 1, \dots, r\}$, are released to the public. Each released $D^{(i,j)}$ includes a label indicating its value of i , i.e. an indicator for its nest.

The agency can sample from each $P^{(i,j)}$ using designs other than simple random samples, such as the stratified sampling in the IAB Establishment Panel synthesis. A complex design can improve efficiency and ensure adequate representation of important sub-populations for analyses. When synthetic data are generated using complex samples, analysts should account for the design in inferences, for example with survey-weighted estimates. One advantage of simple random samples is that analysts can make inferences with techniques appropriate for simple random samples.

The agency could simulate Y for all N units, thereby avoiding the release of actual values of Y . In practice, it is not necessary to generate completed-data populations for constructing the $D^{(i,j)}$; the agency need only generate values of Y for units in the synthetic samples. The formulation of completing the population, then sampling from it, aids in deriving inferential methods.

Let Q be the estimand of interest, such as a population mean or a regression coefficient. For all (i, j) , let $q^{(i,j)}$ be the estimate of Q , and let $u^{(i,j)}$ be the estimate of the variance associated with $q^{(i,j)}$. The $q^{(i,j)}$ and $u^{(i,j)}$ are computed based on the design used to sample from $P^{(i,j)}$. Note that when $n_{syn} = N$, the $u^{(i,j)} = 0$.

The following quantities are necessary for inferences

$$\bar{q}_r^{(i)} = \sum_{j=1}^r q^{(i,j)} / r \quad (8.1)$$

$$\bar{q}_M = \sum_{i=1}^m \bar{q}_r^{(i)} / m = \sum_{j=1}^r \sum_{i=1}^m q^{(i,j)} / mr \quad (8.2)$$

$$b_M = \sum_{i=1}^m (\bar{q}_r^{(i)} - \bar{q}_M)^2 / (m - 1) \quad (8.3)$$

$$w_r^{(i)} = \sum_{j=1}^r (q^{(i,j)} - \bar{q}_r^{(i)})^2 / (r - 1) \quad (8.4)$$

$$\bar{u}_M = \sum_{j=1}^r \sum_{i=1}^m u^{(i,j)} / mr \quad (8.5)$$

The analyst then can use \bar{q}_M to estimate Q and

$$T_{2st,f} = (1 + m^{-1})b_M + (1 - 1/r)\bar{w}_M - \bar{u}_M \quad (8.6)$$

to estimate the variance of \bar{q}_M , where $\bar{w}_M = \sum_{i=1}^m w_r^{(i)} / m$. Inferences can be based on a t -distribution with degrees of freedom

$$\nu_{2st,f} = \left(\frac{((1 + 1/m)b_M)^2}{(m - 1)T_{2st,f}^2} + \frac{((1 - 1/r)\bar{w}_M)^2}{(m(r - 1))T_{2st,f}^2} \right)^{-1}.$$

Derivations of these methods are presented in Reiter and Drechsler (2010). It is possible that $T_{2st,f} < 0$, particularly for small values of m and r . Generally, negative values of $T_{2st,f}$ can be avoided by making n_{syn} or m and r large. To adjust for negative variances, one approach is to use the always positive variance estimator, $T_{2st,f}^* = T_{2st,f} + \lambda \bar{u}_M$, where $\lambda = 1$ when $T_{2st,f} \leq 0$ and $\lambda = 0$ when $T_{2st,f} > 0$. When $T_{2st,f} < 0$, using $\nu_{2st,f}$ is overly conservative, since $T_{2st,f}^*$ tends to be conservative when $\lambda = 1$. To avoid excessively wide intervals, analysts can base inferences on t -distributions with degrees of freedom $\nu_{2st,f}^* = \nu_{2st,f} + \lambda \infty$.

8.1.2 Partially synthetic data

We assume that $Y_{inc} = Y_{obs}$, i.e., there is no missing data. Methods for handling missing data and one stage of partial synthesis simultaneously are presented in

Chapter 7. Developing two stage imputation methods for data that is subject to nonresponse is an area for future research.

The agency generates the partially synthetic data in two stages. Let $Y_a^{(i)}$ be the values imputed in the first stage in nest i , where $i = 1, \dots, m$. Let $Y_b^{(i,j)}$ be the values imputed in the second stage in dataset j in nest i , where $j = 1, \dots, r$. Let Y_{nrep} be the values of Y_{obs} that are not replaced with synthetic data and hence are released as is. Let $Z_{a,l} = 1$ if unit l , for $l = 1, \dots, s$, is selected to have any of its first-stage data replaced, and let $Z_{a,l} = 0$ otherwise. Let $Z_{b,l}$ be defined similarly for the second-stage values. Let $Z = (Z_{a,1}, \dots, Z_{a,s}, Z_{b,1}, \dots, Z_{b,s})$.

To create $Y_a^{(i)}$ for those records with $Z_{a,l} = 1$, first the agency draws from $f(Y_a | D_{obs}, Z)$, conditioning only on values not in Y_b . Second, in each nest, the agency generates $Y_b^{(i,j)}$ for those records with $Z_{b,l} = 1$ by drawing from $f(Y_b^{(i,j)} | D_{obs}, Z, Y_a^{(i)})$. Each synthetic data set $D^{(i,j)} = (X, Y_a^{(i)}, Y_b^{(i,j)}, Y_{nrep}, I, Z)$. The entire collection of $M = mr$ datasets, $D_{syn} = \{D^{(i,j)}, i = 1, \dots, m; j = 1, \dots, r\}$, with labels indicating the nests, is released to the public.

To obtain inferences from nested partially synthetic data, we assume the analyst acts as if each $D^{(i,j)}$ is a sample according to the original design. Unlike in fully synthetic data, there is no intermediate step of completing populations. The analyst again can use \bar{q}_M to estimate Q and

$$T_{2st,p} = \bar{u}_M + b_M/m. \quad (8.7)$$

to estimate the variance of \bar{q}_M . Inferences can be based on a t -distribution with degrees of freedom $\nu_{2st,p} = (m-1)(1 + m\bar{u}_M/b_M)^2$. Derivations of these methods are presented in Reiter and Drechsler (2010). We note that $T_{2st,p} > 0$ always holds, so that negative variance estimates do not arise in two-stage partial synthesis.

8.2 Data utility and disclosure risk

To evaluate the data utility and disclosure risk, we can apply the same methods as with standard one stage synthesis. We refer to Section 5.2 for possible data utility measures and to Section 5.3 and Section 6.3 for possible disclosure risk evaluations.

8.3 Application of the two stage approach to the IAB Establishment Panel

To assess the impact of different numbers of imputations, we first evaluate the trade-off between risk and utility as a function of m for standard one stage imputation. We then compare the results with results achievable with the proposed two stage imputation approach.

For this simulation study, we synthesize two variables in the IAB Establishment Panel for 1997: the number of employees and the industry coded in 16 categories. For both variables, all 7,332 observations are replaced by imputed values. Employment size and industry code are high risk variables since (i) they are easily available in other databases and (ii) the distribution for the number of employees is heavily skewed. Imputations are based on linear models with more than 100 explanatory variables for the number of employees and on a multinomial logit model with more than 80 explanatory variables for the industry. We use large numbers of predictors in hopes of reducing problems from uncongeniality (Meng, 1994). Some variables for the multinomial logit model are dropped for multicollinearity reasons.

8.3.1 Data utility for the panel from one stage synthesis

We investigate data utility for some descriptive statistics and a probit regression. The descriptive statistics are the (unweighted) average number of employees by industry; they are based solely on the two variables we synthesized. The probit regression, which originally appeared in an article by Zwick (2005), is used in various places throughout the book, see Section 5.4.2 for a detailed description.

Tables 8.1 – 8.4 display point estimates and the interval overlap measures for different values of m . For most parameters, increasing m moves point estimates closer to their values in the original data and increases the overlaps in the confidence intervals. Increasing $m = 3$ to $m = 10$ results in the largest increase in data utility, as the average confidence interval overlap over all 31 parameters in Table 8.3 and Table 8.4 increases from 0.828 to 0.867. Increasing $m = 50$ to $m = 100$ does not have much impact on data utility.

Each entry in Table 8.1 – 8.4 results from one replication of a partially synthetic data release strategy. To evaluate the variability across different replications,

Table 8.1: Average number of employees by industry for one stage synthesis.

	original data	m=3	m=10	m=50	m=100
Industry 1	71.5	84.2	84.2	82.6	82.4
Industry 2	839.1	919.4	851.2	870.2	852.9
Industry 3	681.1	557.7	574.5	594.4	593.1
Industry 4	642.9	639.9	644.8	643.5	649.6
Industry 5	174.5	179.8	176.0	183.5	187.4
Industry 6	108.9	132.4	121.8	120.8	120.7
Industry 7	117.1	111.6	112.9	117.1	119.6
Industry 8	548.7	455.3	504.3	514.2	513.0
Industry 9	700.7	676.9	689.4	711.8	713.4
Industry 10	547.0	402.4	490.3	499.3	487.7
Industry 11	118.6	142.7	130.2	132.1	131.0
Industry 12	424.3	405.6	414.9	424.5	425.2
Industry 13	516.7	526.1	549.1	550.2	551.9
Industry 14	128.1	185.8	167.1	160.0	159.0
Industry 15	162.0	292.8	233.4	221.9	238.1
Industry 16	510.8	452.8	449.9	441.5	439.3

we repeated each simulation ten times. Table 8.5 presents the average confidence interval overlap over all 31 estimands for the ten simulations. The variation in the overlap measures decreases with m . This is because the variability in \bar{q}_m and T_m decreases with m , so that results stabilize as m gets large. We believe most analysts would prefer to have stable results across different realizations of the synthesis and hence favor large values of m .

8.3.2 Disclosure risk for the panel from one stage synthesis

To assess the disclosure risk, we assume that the intruder knows which establishments are included in the survey and the true values for the number of employees and industry, i.e. we assume the intruder scenario described in Section 6.3.1. This is a conservative scenario but gives, in some sense, an upper bound on the risk for this level of intruder knowledge. Intruders might also know other variables on the file, in which case the agency may need to synthesize them as well.

The intruder computes probabilities using the approach outlined in Section

Table 8.2: Results from the vocational training regression for one stage partial synthesis revisited.

	original data	m=3	m=10	m=50	m=100
Intercept	-1.319	-1.323	-1.322	-1.323	-1.324
Redundancies expected	0.253	0.268	0.262	0.264	0.264
Many emp. exp. on mat. leave	0.262	0.334	0.316	0.312	0.314
High qualification need exp.	0.646	0.636	0.640	0.640	0.639
Appr. tr. react. on skill short.	0.113	0.098	0.106	0.110	0.112
Training react. on skill short.	0.540	0.529	0.538	0.542	0.543
Establishment size 20-199	0.684	0.718	0.709	0.705	0.701
Establishment size 200-499	1.352	1.363	1.333	1.339	1.343
Establishment size 500-999	1.346	1.315	1.386	1.377	1.367
Establishment size 1000 +	1.955	1.782	1.800	1.778	1.776
Share of qualified employees	0.787	0.787	0.788	0.784	0.785
State-of-the-art tech. equipment	0.171	0.183	0.178	0.174	0.174
Collective wage agreement	0.255	0.268	0.264	0.267	0.268
Apprenticeship training	0.490	0.501	0.510	0.507	0.507
East Germany	0.058	0.038	0.033	0.042	0.044

6.3.1. We assume that the agency does not reveal the synthesis model to the public, so that the only information in M is that employee size and industry were synthesized. For a given target \mathbf{t} , records from each $D^{(i)}$ must meet two criteria to be possible matches. First, the record's synthetic industry code exactly matches the target's true industry code. Second, the record's synthetic number of employees lies within an agency-defined interval around the target's true number of employees. Acting as the agency, we define the interval as follows. We divide the cubic root of the true number of employees into twenty quantiles and calculate the standard deviation of the number of employees within each quantile. The interval is $t_e \pm sd_s$, where t_e is the target's true value and sd_s is the standard deviation of the quantile in which the true value falls. When there are no synthetic records that fulfill both matching criteria, the intruder matches only on the industry code.

We use 20 quantiles because this is the largest number of groups that guarantees some variation within each group. Using more than 20 quantiles results in groups with only one value of employment, which forces exact matching for targets in those quantiles. On the other hand, using a small number of

Table 8.3: Confidence interval overlap for the average number of employees for one stage synthesis.

	m=3	m=10	m=50	m=100
Industry 1	0.778	0.770	0.777	0.782
Industry 2	0.844	0.893	0.853	0.874
Industry 3	0.730	0.776	0.797	0.800
Industry 4	0.983	0.992	0.995	0.971
Industry 5	0.920	0.935	0.863	0.817
Industry 6	0.605	0.749	0.764	0.767
Industry 7	0.809	0.820	0.863	0.876
Industry 8	0.692	0.862	0.894	0.890
Industry 9	0.926	0.966	0.968	0.963
Industry 10	0.660	0.876	0.897	0.871
Industry 11	0.609	0.804	0.773	0.792
Industry 12	0.903	0.912	0.916	0.918
Industry 13	0.946	0.814	0.809	0.799
Industry 14	0.408	0.589	0.655	0.664
Industry 15	0.586	0.639	0.654	0.638
Industry 16	0.666	0.645	0.583	0.566
Average	0.754	0.815	0.816	0.812

quantiles does not differentiate adequately between small and large establishments. For small establishments, we want the potential matches to deviate only slightly from the original values. For large establishments, we accept higher deviations.

We studied the impact of using different numbers of groups for $m = 50$. We found a substantial increase in the risks of identifications, especially for the small establishments, when going from exact matching to five quantiles. Between five and twenty quantiles, the disclosure risk doesn't change dramatically. For more than twenty quantiles, the number of identifications starts to decline again.

Table 8.6 displays the average true matching risk and expected matching risk over the ten simulation runs used in Table 8.5. Since the largest establishments are usually considered as the records most at risk of identification, we also include the risk measures for the largest 25 establishments in parentheses. There is clear evidence that a higher number of imputations leads to a higher risk of disclosure, especially for the largest establishments. This is because,

Table 8.4: Confidence interval overlap for the vocational training regression for one stage synthesis.

	m=3	m=10	m=50	m=100
Intercept	0.987	0.989	0.986	0.984
Redundancies expected	0.931	0.958	0.946	0.948
Many emp. exp. on maternity leave	0.808	0.856	0.867	0.861
High qualification need exp.	0.965	0.977	0.978	0.976
Appr. tr. react. on skill shortages	0.928	0.964	0.984	0.996
Training react. on skill shortages	0.946	0.989	0.989	0.982
Establishment size 20-199	0.802	0.856	0.879	0.902
Establishment size 200-499	0.934	0.939	0.935	0.933
Establishment size 500-999	0.926	0.907	0.928	0.953
Establishment size 1000 +	0.731	0.763	0.727	0.723
Share of qualified employees	0.995	0.997	0.989	0.993
State-of-the-art tech. equipment	0.919	0.953	0.976	0.977
Collective wage agreement	0.926	0.952	0.934	0.927
Apprenticeship training	0.937	0.883	0.899	0.899
East Germany	0.872	0.840	0.899	0.912
Average	0.907	0.922	0.928	0.931

as m increases, the intruder has more information to estimate the distribution that generated the synthetic data. It is arguable that the gains in utility, at least for these estimands, are not worth the increases in disclosure risks.

The establishments that are correctly identified vary across the 10 replicates. For example, for $m = 50$, the total number of identified records over all 10 replicates is 614. Of these records, 319 are identified in only one simulation, 45 are identified in more than five simulations, and only 10 records are identified in all 10 replications. For $m = 10$, no records are identified more than seven times.

The risks are not large on an absolute scale. For example, with $m = 10$, we anticipate that the intruder could identify only 83 establishments out of 7,332. This assumes that the intruder already knows the establishment size and industry classification code and also has response knowledge, i.e. he knows which establishments participated in the survey. Furthermore, the intruder will not know how many of the unique matches (i.e. $c_j = 1$) actually are true matches.

We also investigated the disclosure risk for different subdomains for $m = 50$.

Table 8.5: Average confidence interval overlap for all 31 estimands for ten independent simulations of one stage synthesis.

	m=3	m=10	m=50	m=100
Simulation 1	0.828	0.867	0.870	0.870
Simulation 2	0.864	0.869	0.869	0.874
Simulation 3	0.858	0.866	0.873	0.868
Simulation 4	0.881	0.861	0.874	0.871
Simulation 5	0.872	0.865	0.866	0.875
Simulation 6	0.845	0.862	0.869	0.865
Simulation 7	0.849	0.851	0.871	0.873
Simulation 8	0.841	0.862	0.871	0.873
Simulation 9	0.841	0.877	0.873	0.872
Simulation 10	0.861	0.865	0.874	0.867
Average	0.854	0.865	0.871	0.871

Four of the sixteen industry categories had less than 200 units in the survey. For these categories, the percentage of identified records ranged between 5% and almost 10%. For the remaining categories, the percentage of correct identifications never went beyond 2.3%. If these risks are too high, the agency could collapse some of the industry categories.

The percentage of identified establishments was close to 5% for the largest decile of establishment size and never went beyond 2.5% for all the other deciles. The identification risk is higher for the top 25 establishments, but still moderate. When $m = 3$ only two of these establishments are correctly identified; this increases to seven establishments with $m = 100$. The intruder also makes many errors when declaring matches for these establishments. In fact, the false match rate for these top establishments is 87% for $m = 3$, 77% for $m = 10$, and approximately 70% for $m = 50$ and $m = 100$. None of the top 10 establishments are identified in all ten simulations.

The largest establishment's size is reduced by at least 10% in all synthetic datasets. We note that this can be viewed as reduction in data utility, since the tail is not accurate at extreme values. It may be possible to improve tail behavior with more tailored synthesis models, such as CART approaches (Reiter, 2005d).

As noted previously, these risk computations are in some ways conservative. First, they presume that the intruder knows which records are in the survey.

Table 8.6: Averages of disclosure risk measures over ten simulations of one stage synthesis. Measures for the 25 largest establishments are reported in parentheses.

	m=3	m=10	m=50	m=100
Expected match risk	67.8 (3.2)	94.8 (5.2)	126.9 (6.9)	142.5 (7.1)
True match risk	35.2 (2.0)	82.5 (4.9)	126.1 (6.8)	142.4 (7.1)

This is not likely to be true, since most establishments are sampled with probability less than one. However, large establishments are sampled with certainty, so that the risk calculations presented here apply for those records. Drechsler and Reiter (2008) show how to adjust the identification disclosure probabilities for intruder uncertainty due to sampling. In their application, the true match rate is 6% when the intruder knows which records are in the sample, and only 1% when the intruder does not know which records are in the sample. Second, the risk measurements presume that the intruder has precise information on establishment size and industry code. In Germany, it is not likely that intruders will know the sizes of all establishments in the survey, because there is no public information on small establishments. However, intruders can obtain size and industry type for large companies from public databases. They also can purchase large private databases on establishments, although the quality of these databases for record linkage on employee size is uncertain. Thus, except for possibly the largest establishments, the risk measures here could overstate the probabilities of identification.

8.3.3 Results for the two stage imputation approach

For the two stage imputation, we impute the industry in stage one and the number of employees in stage two. Exchanging the order of the imputation does not materially impact the results. We consider different values of m and r . We run ten simulations for each setting and present the average estimates over these ten simulations.

Table 8.7 displays the average confidence interval overlap for all 31 parameters and the two disclosure risk measures for the different settings. As with one stage synthesis, there is not much difference in the data utility measures for different M , although there is a slight increase when going from $M = 9$ to

Table 8.7: Average CI overlap and match risk for two stage synthesis based on ten simulations. Match risk for largest 25 establishments is in parentheses.

m,r	Avg. overlap	Expected match risk	True match risk
m=3,r=3	0.867	83.1 (4.0)	67.6 (3.4)
m=3,r=16	0.868	98.0 (4.1)	91.8 (4.0)
m=3,r=33	0.870	99.8 (3.8)	96.3 (3.8)
m=5,r=10	0.869	106.1 (4.6)	101.2 (4.4)
m=10,r=5	0.875	113.8 (5.0)	109.4 (5.0)
m=16,r=3	0.874	119.9 (5.2)	116.4 (5.2)

$M \approx 50$. The two stage results with $M = 9$ (average overlap of .867) are slightly better than the one stage results with $m = 10$ (average overlap of .865). The two stage results with $M \approx 50$ are approximately on the same level or slightly above the one stage results for $m = 50$ (average overlap of .871).

The improvements in data utility when using the two stage approach are arguably minor, but the reduction in disclosure risks is more noticeable. The measures are always substantially lower for the two stage approach compared to the one stage approach with approximately the same number of synthetic datasets. For example, releasing two stage synthetic data with $M = 9$ carries an average true match risk of 67 (3.4 for the top 25 establishments), whereas releasing one stage synthetic data with $m = 10$ has a true match risk of 82 (4.9). Risks are lower for $M \approx 50$ as compared to one stage with $m = 50$ as well. Additionally, for the top 25 establishments, the percentage of unique matches that are true matches is lower for the two stage approach. When $M = 9$, this percentage is 17% for the two stage approach compared to around 23% for one stage synthetic data with $m = 10$. When $M \approx 50$, this percentage varied between 18% and 22%, whereas it is around 30% for one stage synthetic data with $m = 50$.

The two stage methods have lower disclosure risks at any given total number of released datasets because they provide fewer pieces of data about industry codes. This effect is evident in the two stage results with $M \approx 50$. The risks increase monotonically with the number of imputations dedicated to the first stage.

Chapter 9

Chances and Obstacles for Multiply Imputed Synthetic Datasets

The main focus of the first statistical disclosure limitation (SDL) techniques proposed in the literature was on providing sufficient disclosure protection. At that time, agencies paid only little attention to the negative impacts of these approaches on data utility. Over the years more and more sophisticated methods evolved. However, these methods also became more complicated to implement and often required correction methods difficult to apply for non standard analysis. For these reasons most agencies still tend to rely on standard, easy to implement SDL techniques like data swapping or noise addition although it has been repeatedly shown that these methods can have severe negative consequences on data utility and may even fail to fulfill their primary goal - to protect the data sufficiently (see for example Winkler (2007b)).

Generating multiply imputed synthetic datasets is a promising alternative. With this approach the user doesn't have to learn complicated adjustments that might differ depending on the kind of analysis the user wants to perform. Instead he can use the simple and straightforward to calculate combining rules presented in this book. With any synthetic data approach that is based on multiple imputation, the point estimate is simply the average of the point estimates calculated for every dataset and its variance is estimated by a simple combination of the estimated variance within each dataset and the variance of the point estimates between the dataset. Furthermore, it is possible with syn-

thetic datasets to account for many real data problems like skip patterns and logical constraints (see Section 3.3 for details). Most standard SDL techniques can not deal with these problems. Besides, it is very easy to address missing data problems and confidentiality problems at the same time when generating partially synthetic datasets. Since both problems can be handled by multiple imputation, it is reasonable to impute missing values first and then generate synthetic datasets from the multiply imputed datasets as described in Chapter 7. This will actually increase the value of the generated datasets since the fully imputed, not synthesized datasets could be used by other researchers inside the agency that otherwise might not be able to adjust their analyses to account for the missing values properly.

However, most research on generating synthetic data, especially with real data applications, dates back no more than 5 years, so it is not surprising that at the current stage there are some obstacles for this approach that still need to be addressed. First and foremost, many agencies complain that developing synthetic datasets for complex surveys is too labor intensive, takes too long and requires experts that are familiar with the data on the one hand, but also need detailed knowledge in Bayesian statistics and excellent modeling skills to generate synthetic data with a high level of data utility. Many small agencies cannot afford to fund research on synthetic data for several months or even years. Other agencies are reluctant to invest into a new data disseminating strategy before the usefulness of this strategy has been clearly demonstrated. This may change with the release of high quality synthetic data in the U.S. and in Germany. Besides, a new version of the multiple imputation software IVEware (Raghunathan *et al.*, 2002) for generating synthetic datasets is under development at the University of Michigan. This software will allow researchers without a sound background in modeling and Bayesian statistics to develop synthetic data. Another promising approach that might speed up the synthetic data generation is the use of non parametric imputation methods like CART (Reiter, 2005d). With this approach, the modeling is mostly automatic, the researcher only needs to define the minimum number of records in each leaf and a threshold value for the homogeneity criterion below which no splits should occur. This can significantly simplify the modeling task. Evaluating to what extent the synthesis can be automated and testing the feasibility of this approach for complex datasets with skip patterns and logical constraints is an area for future research.

But it is not only the agencies that are concerned about this new data disseminating strategy. Many potential users of the released data are skeptical about the approach, too. They insist that they would only work with the original data, ignoring the fact that unrestricted access to the original data is not an option in many cases. It is important that users understand that they should focus on the potential benefits of this approach relative to other SDC methods instead of comparing the approach with unrestricted access. They also tend to see the original data as the true data ignoring other sources of uncertainty and potential bias like nonresponse, undercoverage, reporting or coding errors, etc. that might dwarf the additional bias potentially introduced by the synthesis. Furthermore, a common misconception is that the synthetic data will only provide valid results, if the imputation model and the analysts model match exactly. This is not true. If the imputation model contains more information than the analysts model, the results will still be valid albeit with a reduced efficiency. But even if the imputation model does not contain all the variables that are included in the analysts model, this does not necessarily mean that the results will be biased. In fact, if one variable is omitted from the imputation, the model implicitly assumes conditional independence between the dependent variable and this variable. Now, if the imputation model is already based on hundreds of variables, the assumption of conditional independence given all the other variables might be appropriate. In this case, the analyst would obtain valid results with the released data, even if some of the information in her model was not included in the imputation model.

Still, it would be misleading to praise the synthetic data approach as the panacea for data dissemination. It is simply impossible to generate a dataset with any kind of statistical disclosure limitation technique that provides valid results for any potential analysis while at the same time guaranteeing 100% disclosure protection. The synthetic data reflect only those relationships included in the data generation models. When the models fail to reflect accurately certain relationships, analysts' inferences also will not reflect those relationships. Similarly, incorrect distributional assumptions built into the models will be passed on to the users' analyses. In practice, this dependence means that some analyses cannot be performed accurately, and that agencies need to release information that helps analysts decide whether or not the synthetic data are reliable for their analyses. For example, agencies might include summaries of the posterior distributions of parameters in the data generation models as

attachments to public releases of data. Or, they might include generic statements that describe the imputation models, such as “Main effects for age, sex, and race are included in the imputation models for education.” This transparency also is a benefit of the synthetic data approach: analysts are given indications of which analyses can be reliably performed with the synthetic data. Analysts who desire finer detail than afforded by the imputations may have to apply for special access to the observed data.

To overcome the scepticism against synthetic data, agencies can also offer some incentives to work with the synthetic data. For example, both, the research team at Cornell and the research team at the IAB independently decided to offer the guarantee that for an initial phase any research that is performed on the synthetic data will also be run on the original data and the results from the original data will be sent back to the researcher after checks for potential confidentiality violations. This is a very strong incentive since researchers do not have to apply for access to the research data center but still can be sure that they will finally get the results from the original data. At the same time, they can compare the results from the original data with the results from the synthetic data and if they repeatedly find out that the results actually do not differ very much, they hopefully give up some of their reservations against the use of synthetic data over time.

Finally, researchers tend to be reluctant to use new methods until they are implemented in standard statistical software and results are easily obtainable using standard commands. For example, the use of multiple imputation has significantly increased since routines to multiply impute missing values and to analyze the imputed data are readily available in all major statistical software packages like Stata, SAS or R. We suggest that agencies work with academic researchers and software developers to write software routines that implement the combining rules necessary to obtain valid results for the different synthetic data approaches.

The interest in synthetic data is ever growing and many seemingly insurmountable obstacles have been overcome in the last few years. There are still some efforts necessary to make the concept a universal, widely accepted, and easy to implement approach, but the first releases of partially synthetic datasets in the US and in Germany demonstrate that the approach successfully managed the critical step from a pure theoretical concept to practical implementation. Nevertheless, plenty of room remains for future research in this area that will

further improve the feasibility of this approach. With the continuous proliferation of publicly available databases and improvements in record linkage technologies releasing synthetic datasets might soon be the only reasonable strategy to balance the trade-off between disclosure risk and data utility when disseminating data collected under the pledge of privacy to the public.

Appendix

A.1 Bill Winkler's Microdata Confidentiality References (01. August 2009)

- Abowd, J. M., and Vilhuber, L. (2008). "How Protective are Synthetic Data?" in (J. Domingo-Ferrer and V. Yucel, eds.) *Privacy in Statistical Databases*, New York, N.Y.: Springer, 239-246.
- Abowd, J. M., and Woodcock, S. D. (2002), "Disclosure Limitation in Longitudinal Linked Data", in (P. Doyle et al, eds.) *Confidentiality, Disclosure, and Data Access*, Amsterdam, The Netherlands: North Holland.
- Abowd, J. M., and Woodcock, S. D. (2004), "Multiply-Imputing Confidential Characteristics and File Links in Longitudinal Linked Data, in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, New York: Springer .
- Adams, N. R., and Wortmann, J. C., (1989), "Security-control Methods for Statistical Databases, A Comparative Study", *ACM Computing Surveys*, 21, 515-556.
- Aggarwal, C. C., (2005), "On k-Anonymity and the Curse of Dimensionality", *Proceedings of the 31st VLDB Conference*, Trondheim, Norway, <http://www.vldb2005.org/program/paper/fri/p901-aggarwal.pdf>.
- Aggarwal, C. C., and Parthasarathy, S. (2001), "Mining Massively Incomplete Data Sets through Conceptual Reconstruction", *Proceedings of the ACM KDD Conference*, 227-232.
- Aggarwal, C. C., and Yu, P. (2004), "A Condensation Approach to Privacy Preserving Data Mining", *Proceedings of the EBDT Conference*, 183-199.
- Aggarwal, G., Feder, T., Kenthapadi, K., Motwani, R., Panigrahy, R., Thomas, D., and Zhu, A. (2005), "Anonymizing Tables", *International*

Conference on Database Theory.

- Aggarwal, G., Feder, T., Kenthapadi, K., Khuller, S. Panigrahy, R., Thomas, D., and Zhu, A. (2006), "Achieving Anonymity via Clustering", *ACM PODS '06*.
- Agrawal, D., and Aggarwal, C. C. (2001), "On the Design and Quantification of Privacy Preserving Data Mining Algorithms", Association of Computing Machinery, *Proceedings of PODS 2001*, 247-255.
- Agrawal, R., and Srikant, R. (2000), Privacy Preserving Data Mining, *Proceedings of the ACM SIGMOD 2000*, 439-450.
- Agrawal, R., Srikant, R., and Thomas, D. (2005), "Privacy Preserving OLAP", *ACM SIGMOD Conference*.
- Agrawal, R., Kiernan, J., Srikant, R., and Xu, Y. (2002), "Hippocratic Databases", *Very Large Databases 2002*.
- Bacher, J., Bender, S., and Brand, R. (2001), "Re-identifying Register Data by Survey Data: An Empirical Study", presented at the *UNECE Workshop On Statistical Data Editing*, Skopje, Macedonia, May 2001.
- Bacher, J., Brand, R., and Bender, S. (2002) Re-identifying Register Data by Survey Data using Cluster Analysis: An Empirical Study, *International Journal of Uncertainty, Fuzziness, Knowledge-Based Systems*, 10 (5) 589-608.
- Barak, B., Chaudhuri, K., Dwork, C., Kale, S., McSherry, F., and Talwar, K. (2007), "Privacy, Accuracy, and Consistency Too: A Holistic Solution to Contingency Table Release", *PODS '07*, Beijing, China.
- Bayardo, R. J., and Agrawal, R. (2005), "Data Privacy Through Optimal K-Anonymization", *IEEE 2005 International Conference on Data Engineering*.
- Bethlehem, J. A., Keller, W. J., and Pannekoek, J., (1990), "Disclosure Control of Microdata", *Journal of the American Statistical Association*, 85, 38-45.
- Blien, U., Wirth, U., and Muller, M. (1992), "Disclosure Risk for Microdata Stemming from Official Statistics", *Statistica Neerlandica*, 46, 69-82.
- Blum, A., Dwork, C., McSherry, F., and Nissim, K. (2005), "Practical Privacy: The SuLQ Framework", *ACM SIGMOD Conference* (also <http://research.microsoft.com/research/sv/DatabasePrivacy/bdmn.pdf>).

- Brand, R. (2002), "Microdata Protection Through Noise Addition", in (J. Domingo-Ferrer, ed.) *Inference Control in Statistical Databases*, New York: Springer, 97-116.
- Castro, J. (2004), "Computational Experience with Minimum-Distance Controlled Perturbation Methods", in (J. Domingo-Ferrer, ed.), *Privacy in Statistical Databases 2004*, Springer: New York.
- Chawla, S., Dwork, C., McSherry, F., Smith, A., and Wee, H. (2004), "Toward Privacy in Public Databases", Microsoft Research Technical Report, Theory of Cryptography Conference.
- Chawla, S., Dwork, C., McSherry, F., and Talwar, K. (2005), "On the Utility of Privacy-Preserving Histograms", <http://research.microsoft.com/research/sv/DatabasePrivacy/cdmt.pdf>.
- Dalenius, T., and Reiss, S.P. (1982), "Data-swapping: A Technique for Disclosure Control", *Journal of Statistical Planning and Inference*, 6, 73-85.
- Dandekar, R. A. (2004), Maximum Utility Minimum Information Loss Table Server Design of Statistical Disclosure Control of Tabular Data, in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, New York: Springer, 121-135.
- Dandekar, R. A., Domingo-Ferrer, J., and Sebe, F. (2002), "LHS-Based Hybrid Microdata vs Rank Swapping and Microaggregation for Numeric Microdata Protection", in (J. Domingo-Ferrer, ed.) *Inference Control in Statistical Databases*, New York: Springer, 175-186.
- Dandekar, R., Cohen, M., and Kirkendal, N. (2002), "Sensitive Microdata Protection Using Latin Hypercube Sampling Technique", in (J. Domingo-Ferrer, ed.) *Inference Control in Statistical Databases*, New York: Springer, 117-125.
- Defays, D., and Anwar, M. N. (1998), "Masking Microdata Using Microaggregation", *Journal of Official Statistics*, 14, 449-461.
- Defays, D., and Nanopolis, P. (1993), "Panels of Enterprises and Confidentiality: the Small Aggregates Method", in *Proceedings of the 1992 Symposium on Design and Analysis of Longitudinal Surveys*, 195-204.
- De Waal, A. G., and Willenborg, L.C.R.J. (1995), "Global Recodings and Local Suppressions in Microdata Sets", *Proceedings of Statistics Canada Symposium 95*, 121-132.

- De Waal, A. G., and Willenborg, L.C.R.J. (1996), "A View of Statistical Disclosure Control for Microdata", *Survey Methodology*, 22, 95-103.
- De Wolf, P.-P. (2007), "Risk, Utility and PRAM: A Comparison of Proximity Swap and Data Shuffle for Numeric Data", in (J. Domingo-Ferrer, ed.) *Statistical Data Protection 2006*, Springer: New York, N.Y.
- Dinur, I., and Nissim, K. (2003), "Revealing Information while Preserving Privacy", *ACM PODS Conference*, 202-210.
- Domingo-Ferrer, J. (2001), "On the Complexity of Microaggregation", presented at the *UNECE Workshop On Statistical Data Editing*, Skopje, Macedonia, May 2001.
- Domingo-Ferrer, J. (ed.) (2002) *Inference Control in Statistical Databases*, New York: Springer
- Domingo-Ferrer, J., and Mateo-Sanz, J. M. (2001), "An Empirical Comparison of SDC Methods for Continuous Microdata in Terms of Information Loss And Re-Identification Risk", presented at the *UNECE Workshop On Statistical Data Editing*, Skopje, Macedonia, May 2001.
- Domingo-Ferrer, J., and Mateo-Sanz, J. M. (2002), "Practical Data-Oriented Microaggregation for Statistical Disclosure Control", *IEEE Transactions on Knowledge and Data Engineering*, 14 (1), 189-201.
- Domingo-Ferrer, J., Mateo-Sanz, J., Oganian, A., and Torres, A. (2002), "On the Security of Microaggregation with Individual Ranking: Analytic Attacks", *International Journal of Uncertainty, Fuzziness, and Knowledge-Based Systems*, 477-492.
- Domingo-Ferrer, J., Sebé, F., and Castellà-Roca, J. (2004), "On the Security of Noise Addition for Privacy In Statistical Databases, in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, New York: Springer, 149-161.
- Domingo-Ferrer, J., and Torra, V. (2001) A Quantitative Comparison of Disclosure Control Methods for Microdata in (P. Doyle, J. Lane, J. Theeuwes, and L. Zayatz, eds.) *Confidentiality, Disclosure Control and Data Access: Theory and Practical Applications*, Amsterdam, The Netherlands: North Holland, 111-134.
- Domingo-Ferrer, J., and Torra, V. (2003), "Statistical Data Protection in Statistical Microdata Protection Via Advanced Record Linkage", *Statistics and Computing*, 13 (4), 343-354.

- Du, W., Han, Y. S., and Chen, S. (2004), "Privacy Preserving Multivariate Statistical Analysis: Linear Regression and Classification", *SIAM International Conference on Data Mining 2004*.
- Du, W., and Zhan, Z. (2003), "Using Randomized Response Techniques for Privacy Preserving Data Mining", *ACM Knowledge Discovery and Data Mining Conference 2003*, 505-510.
- DuMouchel, W., Volinsky, C., Johnson, T., Cortes, C., and Pregibon, D. (1999), "Squashing Flat Files Flatter", *Proceedings of the ACM Knowledge Discovery and Data Mining Conference*, 6-15.
- Duncan, G. T., Keller-McNulty, S. A., and Stokes, S. L. (2001), "Disclosure Risk vs. Data Utility: The R-U Confidentiality Map", Los Alamos National Laboratory Technical Report LA-UR-01-6428.
- Dwork, C. (2006), "Differential Privacy", *33rd International Colloquium on Automata, Languages and Programming - ICALP 2006*, Part II, 1-12.
- Dwork, C. (2008), "Differential Privacy: A Survey of Results", in (M. Agrawal et al., eds.) *TAMC 2008*, LNCS 4978, 1-19.
- Dwork, C., and Lei, J. (2009), "Differential Privacy and Robust Statistics", *STOC '09*.
- Dwork, C., McSherry, F., Nissim, K., and Smith, A. (2006), "Calibrating Noise to Sensitivity in Private Data Analysis", *3rd Conference on Cryptography - TCC 2006*, 365-384.
- Dwork, C., McSherry, F., and Talwar, K. (2007a), "The Price of Privacy and the Limits of LP Decoding", *STOC '07*, San Diego, CA.
- Dwork, C., McSherry, F., and Talwar, K. (2007b), "Differentially Private Marginals Release with Mutual Consistency and Error Independent Sample Size", *UNECE Worksession on Statistical Data Confidentiality*, Manchester, UK, at <http://www.unece.org/stats/documents/2007/12/confidentiality/wp.19.e.pdf>.
- Dwork, C. and Yekhanin, S. (2008), "New Efficient Attacks on Statistical Disclosure Control Mechanisms", *Advances in Cryptology-CRYPTO 2008*, to appear, also at <http://research.microsoft.com/research/sv/DatabasePrivacy/dy08.pdf>.
- Dwork, C. and Nissim, K. (2004), "Privacy-Preserving Datamining on Vertically Partitioned Databases", Microsoft Research Technical Report.

- Elamir, E. A. H. (2004), "Analysis of Re-identification Risk based on Log-Linear Model", in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, Springer: New York, 273-281.
- Elamir, E. A. H. and Skinner, C. J. (2006), "Record Level Measures of Disclosure Risk for Survey Microdata", *Journal of Official Statistics*, 22, 525-539.
- Elliott, M. A., Manning, A. M., and Ford, R. W. (2002), "A Computational Algorithm for Handling the Special Uniques Problem", *International Journal of Uncertainty, Fuzziness, and Knowledge-Based Systems*, 10 (5), 493-510.
- Elliot, M.J., Skinner, C. A., and Dale, A. (1998), "Special Uniques, Random Uniques, and Sticky Populations: Some Counterintuitive Effects of Geographic Detail in Disclosure Risk", *Research in Official Statistics*, 1, 53-68.
- Evfimievski, A. (2004), "Privacy Preserving Information Sharing", Ph.D. Dissertation, Cornell University, <http://www.cs.cornell.edu/aevf/>.
- Evfimievski, A., Gehrke, J., and Srikant, R. (2003), "Limiting Privacy Breaches in Privacy Preserving Data Mining", *ACM PODS Conference*, 211-222.
- Evfimievski, A., Srikant, R., Agrawal, R., and Gehrke, J. (2002), "Privacy Preserving Mining of Association Rules", Association of Computing Machinery, Special Interest Group on Knowledge Discovery and Datamining 2002.
- Fellegi, I. P., and Sunter, A. B. (1969), "A Theory for Record Linkage", *Journal of the American Statistical Association*, 64, 1183-1210.
- Fellegi, I. P. (1972), "On the Question of Statistical Confidentiality", *Journal of the American Statistical Association*, 67, 7-18.
- Fellegi, I. P. (1999), "Record Linkage and Public Policy - A Dynamic Evolution", *Proceedings of the Record Linkage Workshop 1997*, Washington, DC: National Academy Press, 3-12.
- Fienberg, S. E. (1997), "Confidentiality and Disclosure Limitation Methodology: Challenges for National Statistics and Statistical Research, commissioned by Committee on National Statistics of the National Academy of Sciences.

- Fienberg, S. E. and MacIntyre, J. (2005), "Data Swapping: Variations on a Theme of Dalenius and Reiss", *Journal of Official Statistics*, 21 (2), 309-323.
- Fienberg, S. E., Makov, E. U., and Sanil, A. P., (1997), "A Bayesian Approach to Data Disclosure: Optimal Intruder Behavior for Continuous Data", *Journal of Official Statistics*, 14, 75-89.
- Fienberg, S. E., Makov, E. U., and Steel, R. J. (1998), "Disclosure Limitation using Perturbation and Related Methods for Categorical Data", *Journal of Official Statistics*, 14, 485-502.
- Frakes, W., and Baeza-Yates, R. (1992), *Information Retrieval - Data Structures and Algorithms*, Upper Saddle River, NJ: Prentice-Hall.
- Franconi, L., Capobianchi, A., Polletini, S., and Seri, G. (2001), "Experiences in Model-Based Disclosure Protection", presented at the *UNECE Workshop on Statistical Data Confidentiality*, Skopje, Macedonia, May 2001.
- Fuller, W. A. (1993), "Masking Procedures for Microdata Disclosure Limitation", *Journal of Official Statistics*, 9, 383-406 (<http://www.jos.nu/Articles/abstract.asp?article=92383>).
- Fung, B. C. M., Wang, K., and Yu, P. S. (2005), "Top-Down Specialization for Information and Privacy Preservation", *IEEE International Conference on Data Engineering*, 205-216.
- Ganta, S., Prasad, S., and Smith, A. (2008), "Compositional Attacks and Auxiliary Information in Data Privacy", *ACM KDD '08*, 265-273.
- Gopal, R., Goes, P. and Garfinkel, R. (1998) "Confidentiality Via Camouflage: The CVC Approach to Database Query Management", in *Statistical Data Protection '98*, Eurostat, Brussels, Belgium, 1-8. (also (2002) *Operations Research*, 50 (3)).
- Gilburd, B., Schuster, A., and Wolff, R. (2004b), "k-TTP: A New Privacy Model for Large-Scale Distributed Environments", *ACM Knowledge Discovery and Data Mining Conference 2004*, 599-604.
- Gill, L. (1999), "OX-LINK: The Oxford Medical Record Linkage System", in *Record Linkage Techniques 1997*, Washington, DC: National Academy Press, 15-33.
- Gomatam, S. V., and Karr, A. (2003), "On Data Swapping of Categorical Data", American Statistical Association, *Proceedings of the Section on*

Survey Research Methods, CD-ROM.

- Gouweleeuw, J.M., Kooiman, P., Willenborg, L. C. R. J., and De Wolf, P.-P. (1998), "Post Randomisation For Statistical Disclosure Control: Theory and Implementation", *Journal of Official Statistics*, 14, 463-478.
- Graham, P., Young, J., and Penny, R. (2009), "Multiply Imputed Synthetic Data: Evaluation of Hierarchical Bayesian Imputation Models", *Journal of Official Statistics*, 25 (2), 245-268.
- Grim, J., Bocek, P., and Pudil, P. (2001), "Safe Dissemination of Census Results by Means of Interactive Probabilistic Models", *Proceedings of 2001 NTTS and ETK*, Luxembourg: Eurostat, 849-856.
- Huang, Z., Du, W., and Chen, B. (2005), "Deriving Private Information from Randomized Data", *ACM SIGMOD 2005 Conference*, 37-48.
- Huckett, J. C. (2008), "Synthetic Data Methods for Disclosure Limitation", Ph.D. Thesis, Department of Statistics, Iowa State University.
- Huckett, J. C., and Larsen, M. D. (2007), "Microdata Simulation for Confidentiality Protection Using Regression Quantiles and Hot Deck", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM, 3053-3060.
- Hwang, J. T. (1986), "Multiplicative Error-in-Variables Models with Applications to Recent Data Released by the U.S. Department of Energy", *Journal of the American Statistical Association*, 81 (395), 680-688.
- Iyengar, V. (2002), "Transforming Data to Satisfy Privacy Constraints", Association of Computing Machinery, *Knowledge Discovery and Datamining Conference 2002*, 279-288.
- Kantarcioglu, M., and Clifton, C. (2004a), "Privacy-Preserving Distributed Mining of Association Rules on Horizontally Partitioned Data", *IEEE Transactions on Knowledge and Data Engineering*, 16 (9), 1026-1037.
- Kantarcioglu, M., and Clifton, C. (2004b), "When Do Data Mining Results Violate Privacy?" Association of Computing Machinery, *Knowledge Discovery and Data Mining Conference 2004*, 599-604.
- Kargupta, H., Datta, S., Wang, Q., and Ravikumar, K. (2003) "Random Data Perturbation Techniques and Privacy Preserving Data Mining", Expanded version of best paper awarded paper from the *IEEE International Conference on Data Mining*, November, 2003, Orlando, FL, (also

- version to appear in Knowledge and Information Systems Journal, http://www.cs.umbc.edu/~hillol/PUBS/kargupta_privacy03a.pdf).
- Kaufmann, S., Seastrom, M., and Roey, S. (2005), "Do Disclosure Controls to Protect Confidentiality Degrade the Quality of the Data? Data from the 2003 Trends in Mathematics and Science Study", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.
- Keller-McNulty, S., and Unger, E. (2003), "Database Systems: Inferential Security", *Journal of Official Statistics*, 9 (2), 475-499.
- Kennickell, A. B. (1999), "Multiple Imputation and Disclosure Control: The Case of the 1995 Survey of Consumer Finances", in *Record Linkage Techniques 1997*, Washington, DC: National Academy Press, 248-267 (available at <http://www.fcs.m.govunderMethodologyreports>).
- Kifer, D. (2009), "Attacks on Privacy and deFinetti's Theorem", *ACM SIGMOD Conference*.
- Kifer, D. and Gehrke, J. (2006), "Injecting Utility into Anonymized Data Sets", *ACM SIGMOD*.
- Kim, J. J. (1986), "A Method for Limiting Disclosure in Microdata Based on Random Noise and Transformation", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 370-374 (http://www.amstat.org/sections/SRMS/Proceedings/papers/1986_069.pdf).
- Kim, J. J. (1990), "Subdomain Estimation for the Masked Data", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 456-461 (http://www.amstat.org/sections/SRMS/Proceedings/papers/1990_075.pdf).
- Kim, J. J., and Winkler, W. E. (1995), "Masking Microdata Files", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 114-119 (http://www.amstat.org/sections/SRMS/Proceedings/papers/1995_017.pdf), longer report <http://www.census.gov/srd/papers/pdf/rr97-3.pdf>).
- Kim, J. J., and Winkler, W. E. (2001), "Multiplicative Noise for Masking Continuous Data", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.

- Lambert, D. (1993), "Measures of Disclosure Risk and Harm", *Journal of Official Statistics*, 9, 313-331 (<http://www.jos.nu/Articles/abstract.asp?article=92313>).
- Lane, J. (2007), "Optimizing the Use of Microdata: An Overview of the Issues", *Journal of Official Statistics*, 23 (3), 299-317 (<http://www.jos.nu/Articles/abstract.asp?article=233299>).
- Lawrence, C., Zhou, J.L., and Tits, A. L. (1997), "User's Guide for CFSZP Version 2.5: A C Code for Solving (Large Scale) Constrained Nonlinear Inequality Constraints", Unpublished, Electrical Engineering Dept. and Institute for Systems Research, University of Maryland.
- Lakshmanan, L., Ng, R., and Ramesh, G. (2005), "To Do or Not To Do - The Dilemma of Disclosing Anonymized Data", *ACM SIGMOD Conference*.
- Lakshmanan, L. K. S., Ng, R., Ramesh, G. (2008), "On Disclosure Risk Analysis of Anonymized Itemsets in the Presences of Prior Knowledge", *ACM Transactions on Knowledge Discovery from Data*, 2 (3), 13.1-13.44.
- LeFevre, K., DeWitt, D. and Ramakrishnan, R. (2005), "Incognito: Efficient Full-Domain K-Anonymity", *ACM SIGMOD Conference*.
- Li, X.-B., and Sarkar, S. (2007), "A Tree-Based Data Perturbation Technique for Privacy-Preserving Data Mining", *IEEE Transactions on Knowledge and Data Engineering*, 18 (9), 1278-1283.
- Liew, C. K., Choi, U. J., and Liew, C. J. (1991), "A Data Distortion by Probability Distribution", *ACM Transactions on Database Systems*, 10, 395-411.
- Little, R. J. A. (1993), "Statistical Analysis of Masked Data", *Journal of Official Statistics*, 9, 407-426 (<http://www.jos.nu/Articles/abstract.asp?article=92407>).
- Little, R. J. A., and Liu, F. (2002), "Selective Multiple Imputation of Keys for Statistical Disclosure Control in Microdata", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.
- Little, R. J. A., and Liu, F. (2003), "Comparison of SMiKe with Data-Swapping and PRAM for Statistical Disclosure Control of Simulated Microdata", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.

- Lindell, Y. and Pinkas, B. (2002), "Privacy Preserving Data Mining", *Proceedings of Crypto 2000*, Springer LNCS 1880, 20-24.
- Liu, H., Kargupta, H., and Ryan, J. (2007), "Random Projection-Based Multiplicative Data Perturbation for Privacy Preserving Distributed Data Mining", *IEEE Transactions of Knowledge and Data Engineering*, 18 (1), 92-106.
- Machanavajjhala, A., Gehrke, and M., Goetz. (2009), "Data Publishing against Realistic Adversaries", *VLDB 2009*.
- Machanavajjhala, A., Gehrke, J., Kifer, D., Venkitasubramaniam, M. (2005), "l-Diversity: Privacy Beyond k-Anonymity", Cornell CS Dept. technical report, <http://www.cs.cornell.edu/johannes/papers/2005/publishing-icde-final.pdf>.
- Machanavajjhala, A., Kifer, D., Abowd, J. Gehrke, J., and Vilhuber, L. (2008), "Privacy: Theory meets Practice on the Map", *ICDE 2008*.
- Malin, B. Sweeney, L., and Newton, E. (2003), "Trail Re-identification: Learning Who You are from Where You have Been", *Workshop on Privacy in Data*, Carnegie-Mellon University, March 2003.
- McSherry, F. (2009), "Privacy Integrated Queries", *SIGMOD 2009*.
- McSherry, F., and Talwar, K. (2007), "Mechanism design via differential privacy", *Proceedings of the 48th Symposium of the Foundations of Computer Science*.
- Mera, R. (1998), "Matrix Masking Methods That Preserve Moments", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 445-450.
- Moore, R. (1995), "Controlled Data Swapping Techniques For Masking Public Use Data Sets", U.S. Bureau of the Census, Statistical Research Division Report rr96/04, (available at <http://www.census.gov/srd/www/byyear.html>).
- Moore, A. W., and Lee, M. S. (1998), "Cached Sufficient Statistics for Efficient Machine Learning with Large Datasets", *Journal of Artificial Intelligence Research*, 8, 67-91.
- Motwani, R., and Xu, Y. (2007), "Efficient Algorithms for Masking and Finding Quasi-Identifiers", *VLDB '07*.

- Müller, W., Blien, U., and Wirth, H. (1995), "Identification Risks of Micro Data", Evidence from Experimental Studies. *Sociological Methods and Research*, 24, 131-157.
- Muralidhar, K., Batrah, D., and Kirs, P.J. (1995), "Accessibility, Security, and Accuracy in Statistical Databases : The Case for the Multiplicative Fixed Data Perturbation Approach", *Management Science*, 41 (9), 1549-1584
- Muralidhar, K., Parsa, R., and Sarathy, R. (1999), "A General Additive Data Perturbation Method for Database Security", *Management Science*, 45 (10), 1399-1415.
- Muralidhar, K., and Sarathy, R. (1999) "Security of Random Data Perturbation Methods", *ACM Transactions on Database Systems*, 24 (4), 487-493.
- Muralidhar, K., and Sarathy, R. (2003), "A Theoretical Basis for Perturbation Methods", *Statistics and Computing*, 13 (4), 329-335.
- Muralidhar, K., and Sarathy, R. (2006a), "Data Shuffling - A New Masking Approach to Numerical Data", *Management Science*, 52 (5), 658-670.
- Muralidhar, K., and Sarathy, R. (2006b), "A Theoretical Basis for Perturbation Methods", *Journal of Statistics*, 22 (3), 507-524.
- Muralidhar, K. and Sarathy, R. (2007), " 'Easy to Implement' is Putting the Cart before the Horse: Effective Techniques for Masking Numerical Data", Federal Committee on Statistical Methodology Research Conference, to appear on CD-ROM.
- Muralidhar, K., Sarathy, R., and Dandekar, R. (2006), "Why Swap When You Can Shuffle? A Comparison of Proximity Swap and Data Shuffle for Numeric Data", in (J. Domingo-Ferrer, ed.) *Statistical Data Protection 2006*, Springer: New York, N.Y.
- Muralidhar, K., Sarathy, R., and Parsa, R. (2001) "An Improved Security Requirement for Data Perturbation with Implications for E-Commerce", *Decision Sciences*, 32 (4), 683-698.
- Nissim, K., Raskhodnikova, S., and Smith, A. (2007), "Smooth Sensitivity and Sampling in Private Data Analysis", *STOC'07*, June 11-13, 2007, San Diego, California, USA.
- Onn, S. (2007), "Entry Uniqueness in Margined Tables", in (J. Domingo-Ferrer, ed.) *Statistical Data Protection 2006*, Springer: New York, N.Y.

- Owen, A. (2003), "Data Squashing by Empirical Likelihood", *Data Mining and Knowledge Discovery*, 7 (1), 101-113.
- Paas, G. (1988), "Disclosure Risk and Disclosure Avoidance for Microdata", *Journal of Business and Economic Statistics*, 6, 487-500.
- Palley, M. A., and Simonoff, J. S. (1987), "The Use of Regression Methodology for the Compromise of Confidential Information in Statistical Databases", *ACM Transactions on Database Systems*, 12 (4), 593-608.
- Polletini, S. (2003), "Maximum Entropy Simulation for Microdata Protection", *Statistics and Computing*, 13 (4), 307-320.
- Polletini, S., Franconi, L., and Stander, J. (2002), "Model Based Disclosure Protection", in (J. Domingo-Ferrer, ed.) *Inference Control in Statistical Databases*: New York: Springer.
- Polletini, S., and Stander, J. (2004), "A Bayesian Hierarchical Model Approach to Risk Estimation in Statistical Disclosure Limitation", in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, New York: Springer, 247-261
- Raghunathan, T. E. (2003), Evaluation of Inferences from Multiple Synthetic Data Sets Created Using Semiparametric Approach, Panel on Confidential Data Access for Research Purposes, Committee On National Statistics, October 2003.
- Raghunathan, T. E., Lepkowski, J. M., Hoewyk, J. V., and Sollenberger, P. (1998), "A Multivariate Technique for Multiply Imputing Missing Values Using a Series of Regression Models", Survey Research Center, University of Michigan.
- Raghunathan, T. E., and Reiter, J. P. (2007), "The Multiple Adaptations of Multiple Imputation", *Journal of the American Statistical Association*, 102, 1462-1471.
- Raghunathan, T.E., Reiter, J. P., and Rubin, D.R. (2003), "Multiple Imputation for Statistical Disclosure Limitation", *Journal of Official Statistics*, 19, 1-16.
- Raghunathan, T.E., and Rubin, D.R. (2000), "Multiple Imputation for Disclosure Limitation" University of Michigan, Department of Biostatistics technical report
- Reiss, J.P. (1984), "Practical Data Swapping: The First Steps", *ACM Transactions on Database Systems*, 9, 20-37.

- Reiter, J.P. (2002), "Satisfying Disclosure Restrictions with Synthetic Data Sets", *Journal of Official Statistics*, 18, 531-543.
- Reiter, J.P. (2003a), "Inference for Partially Synthetic, Public Use Data Sets", *Survey Methodology*, 181-189.
- Reiter, J.P. (2003b), Estimating Probabilities of Identification for Microdata, Panel on Confidential Data Access for Research Purposes, Committee On National Statistics, October 2003.
- Reiter, J.P. (2005a), "Releasing Multiply Imputed, Synthetic Public-Use Microdata: An Illustration and Empirical Study", *Journal of the Royal Statistical Society*, A, 168 (1), 185-205.
- Reiter, J. P. (2005b), "Estimating Risk of Identify Disclosure in Microdata", *Journal of the American Statistical Association*, 100, 1103-1112.
- Reiter, J. P. (2008), "Selecting the number of imputed datasets when using multiple imputation for missing data and disclosure protection", *Statistics and Probability Letters*, 78, 15-20.
- Reiter, J. P. and Mitra, R. (2009), "Estimating risks of identification disclosure in partially synthetic data", *Journal of Privacy and Confidentiality*, 01 (01), 99-110.
- Roque, G. M. (2000), "Masking Microdata Files with Mixtures of Multivariate Normal Distributions", Ph.D.Dissertation, Department of Statistics, University of California at Riverside.
- Rubin, D. B. (1993), "Satisfying Confidentiality Constraints through the Use of Synthetic Multiply-imputed Microdata", *Journal of Official Statistics*, 91, 461-468.
- Samarati, P. (2001), "Protecting Respondents' Identity in Microdata Release", *IEEE Transactions on Knowledge and Data Engineering*, 13 (6), 1010-1027.
- Samarati, P., and Sweeney, L. (1998), "Protecting Privacy when Disclosing Information: k-anonymity and Its Enforcement through Generalization and Cell Suppression", Technical Report, SRI International.
- Sarathy, R., and Muralidhar, K. (2002), "The Security of Confidential Numerical Data in Databases", *Information Systems Research*, 48 (12), 1613-1627.

- Sarathy, R., Muralidhar, K., and Parsa, R. (2002), Perturbing Non-Normal Attributes: The Copula Approach, *Management Science*, 48 (12), 1613-1627
- Scheuren, F., and Winkler, W. E. (1993), "Recursive Merging and Analysis of Administration Lists", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 123-128 (presently available on <http://www.amstat.org/intheSectiononGovernmentStatistics>).
- Scheuren, F., and Winkler, W. E. (1997), "Regression Analysis of Data Files that are Computer Matched - Part II", *Survey Methodology*, 157-165).
- Schlörér, J. (1981), "Security of Statistical Databases: Multidimensional Transformation", *ACM Transactions on Database Systems*, 6, 91-112.
- Skinner, C. J., and Elliot, M. A. (2001), "A Measure of Disclosure Risk for Microdata", *Journal of the Royal Statistical Society*, 64 (4), 855-867.
- Skinner, C. J., and Holmes, D. J. Estimating the Re-identification Risk per Record in Microdata, *Journal of Official Statistics*, 14 (1998) 361-372.
- Skinner, C. J. and Shlomo, N. (2007), "Assessing Identification Risk in Survey Data Using Loglinear Models", *Journal of the American Statistical Association*, 103 (483), 989-1001, also at <http://eprints.soton.ac.uk/48103/>.
- Stander, J., and Franconi, L. (2001), "A Model-Based Disclosure Limitation Method for Business Microdata", presented at the *UNECE Workshop On Statistical Data Editing*, Skopje, Macedonia, May 2001.
- Sullivan, G., and Fuller, W. A. (1989), "The Use of Measurement Error to Avoid Disclosure", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 802-807.
- Sullivan, G., and Fuller, W. A. (1990), "Construction of Masking Error for Categorical Variables", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 435-439.
- Sweeney, L. (1999), "Computational Disclosure Control for Medical Microdata: The Datafly System" in *Record Linkage Techniques 1997*, Washington, DC: National Academy Press, 442-453.
- Sweeney, L. (2002), "Achieving k-Anonymity Privacy Protection Using Generalization and Suppression", *International Journal of Uncertainty, Fuzziness, and Knowledge-Based Systems*, 10 (5), 571-588.

- Sweeney, L. (2004), "Optimal Anonymity using K-similar, a New Clustering Algorithm", manuscript.
- Takemura, A. (2002), "Local Recoding and Record Swapping by Maximum Weight Matching for Disclosure Control of Microdata Sets", *Journal of Official Statistics*, 18 (2), 275-289.
- Tendick, P., and Matloff, N. (1994), "A Modified Random Perturbation Method for Database Security", *ACM Transactions on Database Systems*, 19, 47-63.
- Thibaudeau, Y. (2004), "An Algorithm for Computing Full Rank Minimal Sufficient Statistics with Application to Confidentiality Protection, *UN-ECE Statistical Journal*, to appear.
- Thibaudeau, Y., and Winkler, W.E. (2002), "Bayesian Networks Representations, Generalized Imputation, and Synthetic Microdata Satisfying Analytic Restraints", Statistical Research Division Report RR 2002/09 at <http://www.census.gov/srd/www/byyear.html>.
- Thibaudeau, Y., and Winkler, W.E. (2004), "Full Rank Minimal Statistics for Disclosure Limitation and Variance Estimation: A Practical Way to Release Count Information", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.
- Torra, V. (2004), "OWA Operators in Data Modeling and Re-identification", *IEEE Transactions on Fuzzy Systems*, 12 (5), 652-660.
- Torra, V., Domingo-Ferrer, J., and Abowd, J. (2007), "Using Mahalanobis-Distance Record Linkage for Disclosure Risk Assessment", in (J. Domingo-Ferrer, ed.) *Statistical Data Protection 2006*, Springer: New York, N.Y.
- Torra, V., and Miyamoto, S. (2004), "Evaluating Fuzzy Clustering Algorithms for Microdata Protection", in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, New York: Springer, 175-186.
- Trottini, M., and Fienberg, S. E. (2002), "Modelling User Uncertainty for Disclosure Risk and Data Utility", *International Journal of Uncertainty, Fuzziness, and Knowledge-Based Systems*, 511-528.
- Van Den Hout, A., and Van Der Heijden, P. G. M. (2002), "Randomized Response, Statistical Disclosure Control, and Misclassification: A Review", *International Statistical Review*, 70 (2), 269-288.

- Van Gewerden, L., Wessels, A., and Hundepol, A. (1997), "Mu-Argus Users Manual, Version 2", Statistics Netherlands, Document TM-1/D.
- Willenborg, L., and De Waal, T. (1996), *Statistical Disclosure Control in Practice*, Vol. 111, Lecture Notes in Statistics, New York: Springer.
- Willenborg, L., and De Waal, T. (2000), *Elements of Statistical Disclosure Control*, Vol. 155, Lecture Notes in Statistics, New York: Springer.
- Willenborg, L. and Van Den Hout, A. (2006), "Peruco: A Method for Producing Safe and Consistent Microdata", *International Statistical Review*, 74 (2), 271-284.
- Winglee, M., Valliant, R., Clark, J., Lim, Y., Weber, M., and Strudler, M. (2002), "Assessing Disclosure Protection for the SOI Public Use File", American Statistical Association, *Proceedings of the Section on Survey Research Methods*, CD-ROM.
- Winkler, W. E. (1994), "Advanced Methods for Record Linkage, American Statistical Association, *Proceedings of the Section on Survey Research Methods*, 467-472 (<http://www.census.gov/srd/papers/pdf/rr94-5.pdf>).
- Winkler, W. E. (1995), "Matching and Record Linkage", in (B. G. Cox et al, ed.) *Business Survey Methods*, New York: J. Wiley, 355-384 (also <http://www.fcs.m.gov/working-papers/wwinkler.pdf>).
- Winkler, W. E. (1997), "Views on the Production and Use of Confidential Microdata", Statistical Research Division report RR 97/01 at <http://www.census.gov/srd/www/byyear.html>.
- Winkler, W. E. (1998), "Re-identification Methods for Evaluating the Confidentiality of Analytically Valid Microdata", *Research in Official Statistics*, 1, 87-104, <http://www.census.gov/srd/papers/pdf/rrs2005-09.pdf>.
- Winkler, W. E. (2002a), "Using Simulated Annealing for k-anonymity", Statistical Research Division report RR 2002/07 at <http://www.census.gov/srd/www/byyear.html>.
- Winkler, W. E. (2002b), "Single Ranking Micro-aggregation and Re-identification", Statistical Research Division report RR 2002/08 at <http://www.census.gov/srd/www/byyear.html>.
- Winkler, W. E. (2004a), Masking and Re-identification Methods for Public-Use Microdata: Overview and Research Problems, in (J. Domingo-

- Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, Springer: New York, 231-247, also <http://www.census.gov/srd/papers/pdf/rrs2004-06.pdf>.
- Winkler, W. E. (2004b), Re-identification Methods for Masked Microdata, in (J. Domingo-Ferrer and V. Torra, eds.), *Privacy in Statistical Databases*, Springer: New York, 216-230, also <http://www.census.gov/srd/papers/pdf/rrs2004-03.pdf>.
- Winkler, W. E. (2005), "Modeling and Quality of Masked Microdata", American Statistical Association, *Proceedings of the Section on Survey Research Method*, CD-ROM, also <http://www.census.gov/srd/papers/pdf/rrs2006-01.pdf>.
- Winkler, W. E. (2007a), "Analytically Valid Discrete Microdata and Re-identification", <http://www.census.gov/srd/papers/pdf/rrs2007-19.pdf>.
- Winkler, W. E. (2007b), "Examples of Easy-to-implement, Widely Used Masking Methods for which Analytic Properties are not Justified", <http://www.census.gov/srd/papers/pdf/rrs2007-21.pdf>.
- Winkler, W. E. (2008), "General Discrete-data Modeling Methods for Producing Synthetic Data with Reduced Re-identification Risk that Preserve Analytic Properties", *IAB Workshop on Confidentiality and Disclosure*, http://fdz.iab.de/en/FDZ_Events/SDC-Workshop.aspx, Nuremberg, Germany, November 20-21, 2008.
- Winkler, W. E. (2009), Should Social Security numbers be replaced by modern, more secure identifiers?" *Proceedings of the National Academy of Science*.
- Woo, M., Reiter, J. P., Oganian, A., Karr, A. F. (2009) "Global measures of data utility for microdata masked for disclosure limitation", *Journal of Privacy and Confidentiality*, 01 (01), 111-124.
- Xiao, X., and Tao, Y. (2007a), "m-Invariance: Towards Privacy Preserving Re-publication of Dynamic Datasets", *ACM SIGMOD*, 689-700.
- Xiao, X., and Tao, Y. (2007b), "Anatomy: Simple and Effective Privacy Preservation", *VLDB*, 139-150.
- Xiao, X., and Tao, Y. (2008a), "Dynamic Anonymization: Accurate Statistical Analysis with Privacy Preservation", *ACM SIGMOD*, 107-120.

- Xiao, X., and Tao, Y. (2008b), "Output Perturbation with Query Relaxation", *VLDB*, 857-869.
- Yancey, W.E., Winkler, W.E., and Creecy, R. H. (2002) "Disclosure Risk Assessment in Perturbative Microdata Protection", in (J. Domingo-Ferrer, ed.) *Inference Control in Statistical Databases*, New York: Springer, 135-151, (also <http://www.census.gov/srd/papers/pdf/rrs2002-01.pdf>).
- Yang, Z., Zhong, S., and Wright, R. (2005), "Anonymity Preserving Data Collection", *ACM KDD Conference*, 334-343.
- Zhang, N., Wang, S., and Zhao, W. (2005), "A New Scheme on Privacy-Preserving Data Classification", *ACM KDD Conference*, 374-383.
- Zhong, S., Yang, Z., and Wright, R. (2005), "Privacy-Enhancing k-Anonymization of Customer Data", *ACM Principals of Database Systems Conference 2005*, 139-147.
- Zhu, Y., and Liu, L. (2004), "Optimal Randomization for Privacy Preserving Data Mining", *ACM Knowledge Discovery and Data Mining Conference 2004*, 761-765.

A.2 Binned residual plots to evaluate the imputations for the categorical variables

Figures A.1-A.7 present the binned residual plots for all 59 categorical variables with missing rates $\geq 1\%$. For variables with more than two categories, we present a graph for each category (the first category is always defined as the reference category in the multinomial imputation model).¹

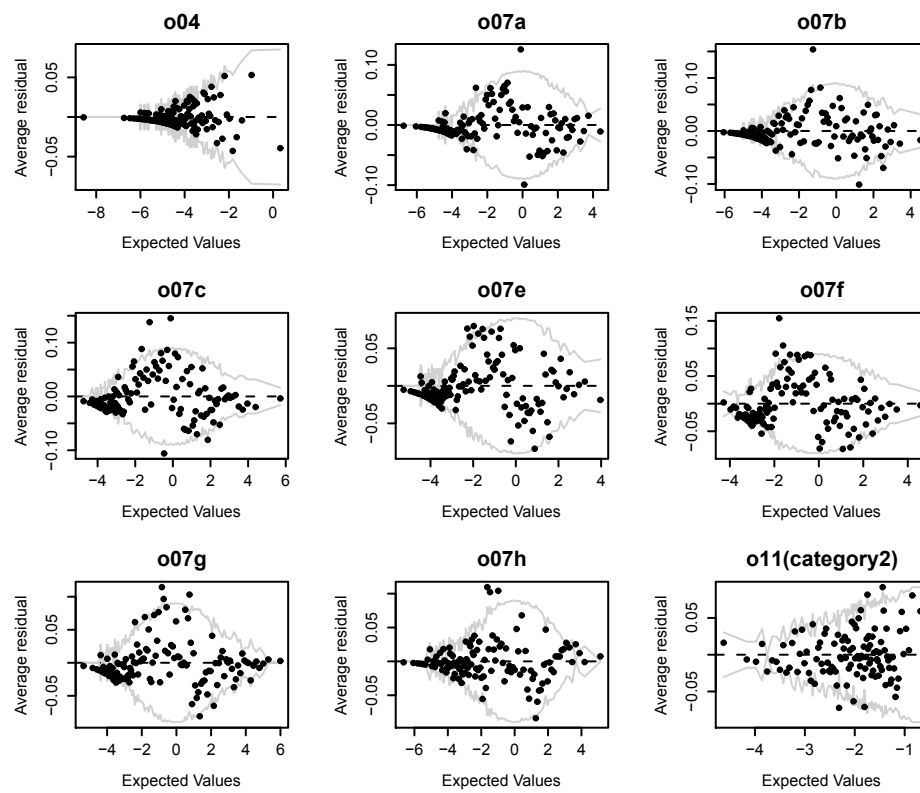


Figure A.1: Binned residual plots for the categorical variables with missing rates above 1%.

¹For readability, we use the internal labeling for the variables. A detailed description of all variables can be obtained from the author upon request.

A.2. BINNED RESIDUAL PLOTS FOR CATEGORICAL VARIABLES 133

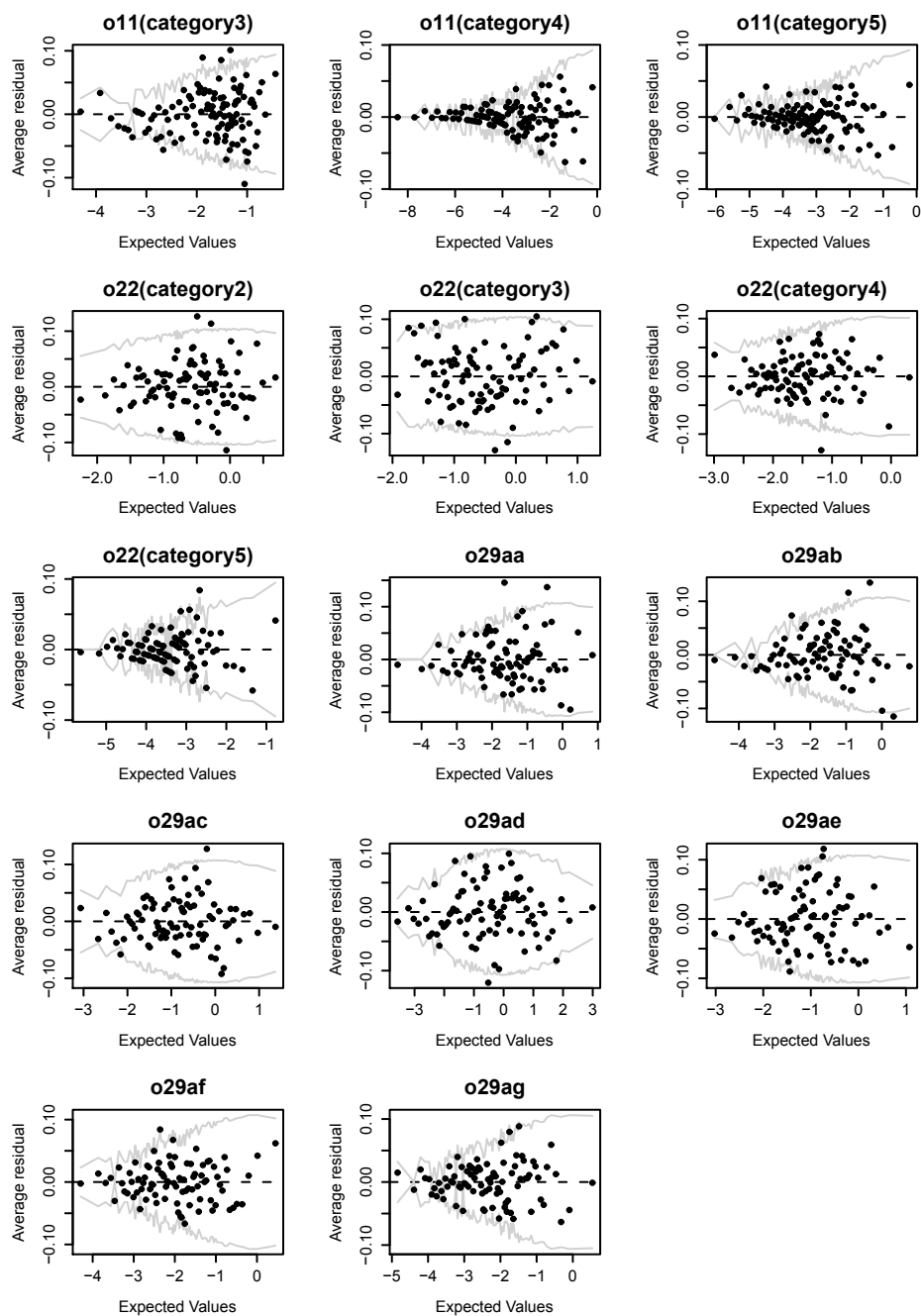


Figure A.2: Binned residual plots for the categorical variables with missing rates above 1%.

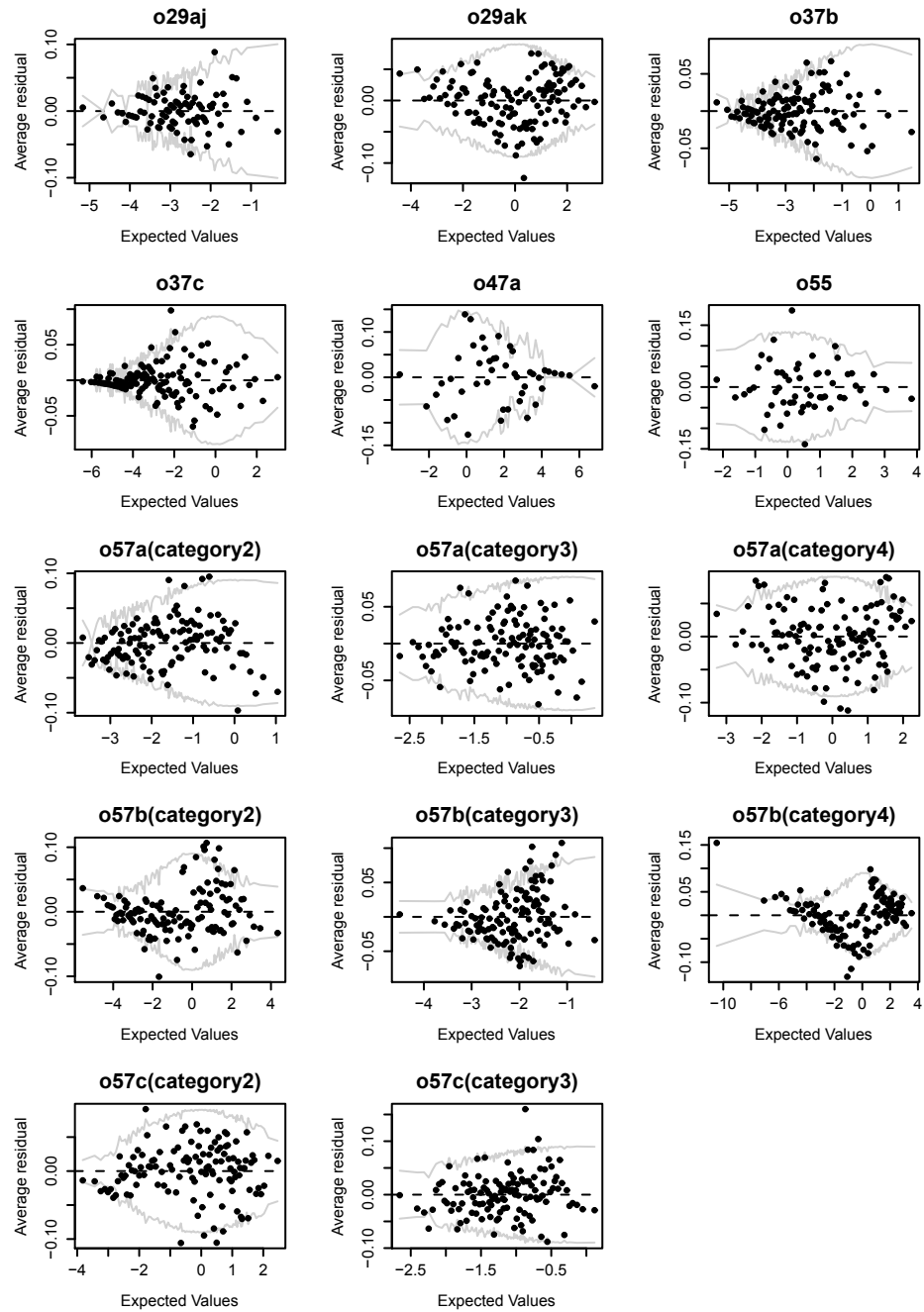


Figure A.3: Binned residual plots for the categorical variables with missing rates above 1%.

A.2. BINNED RESIDUAL PLOTS FOR CATEGORICAL VARIABLES 135

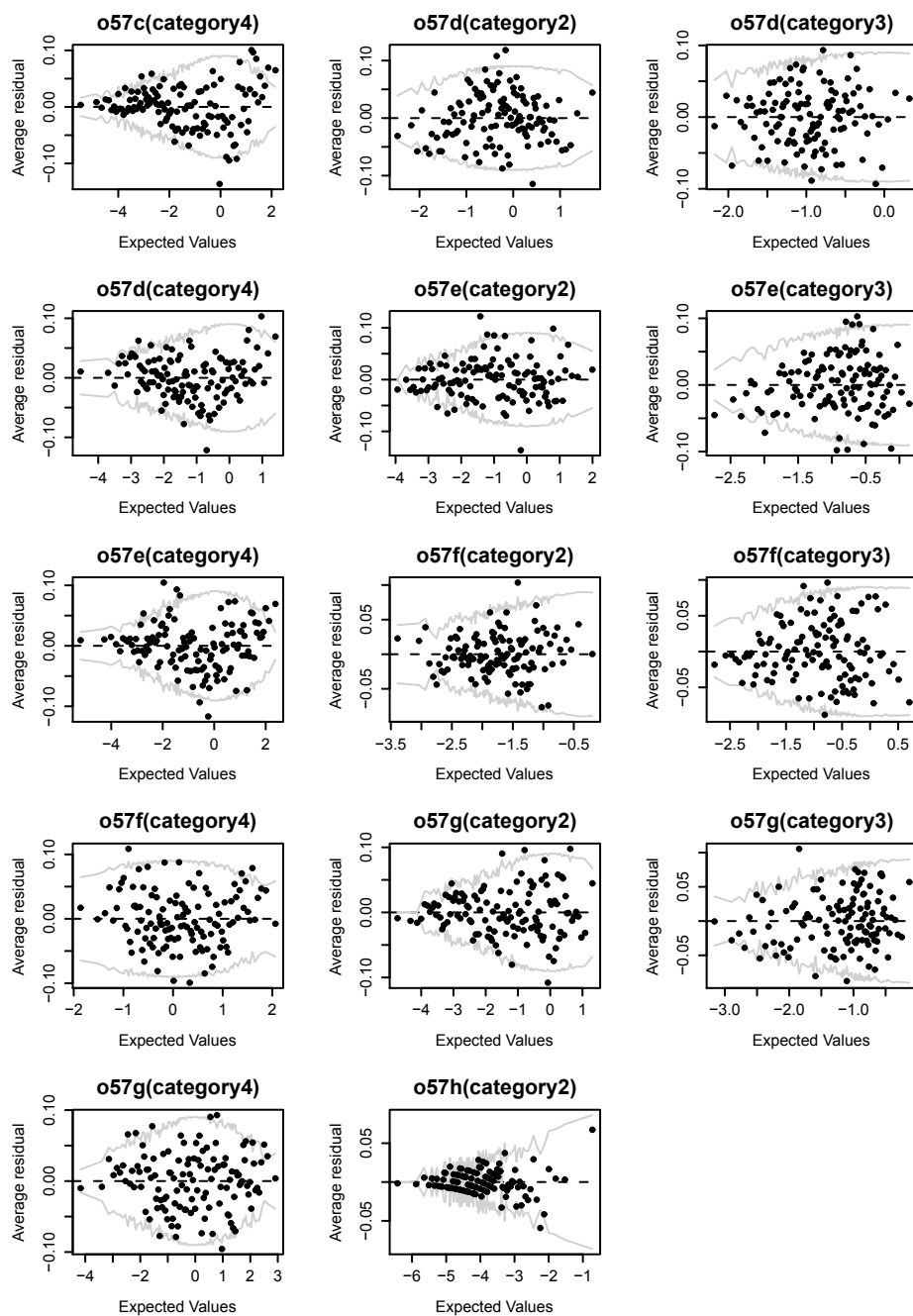


Figure A.4: Binned residual plots for the categorical variables with missing rates above 1%.

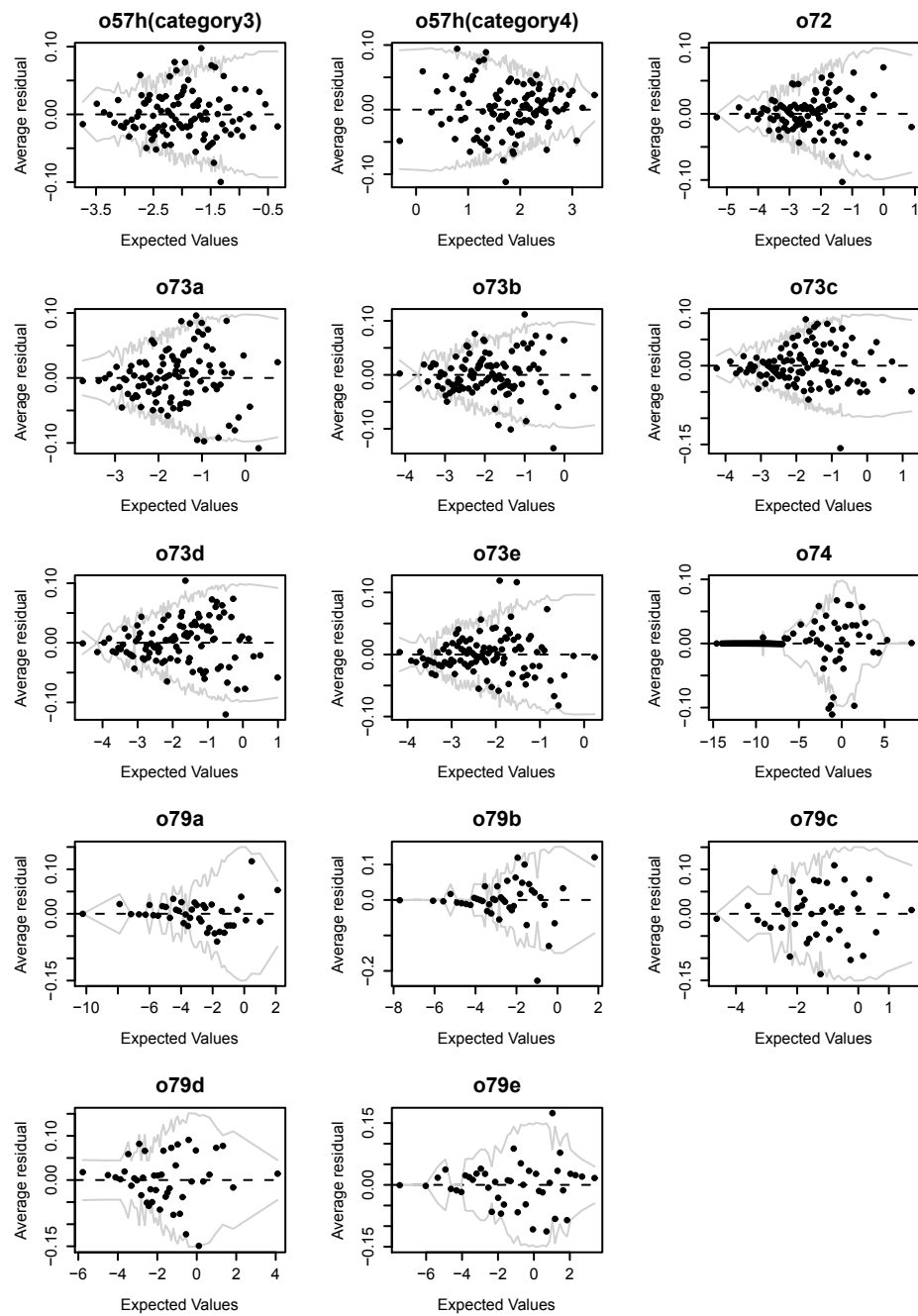


Figure A.5: Binned residual plots for the categorical variables with missing rates above 1%.

A.2. BINNED RESIDUAL PLOTS FOR CATEGORICAL VARIABLES 137

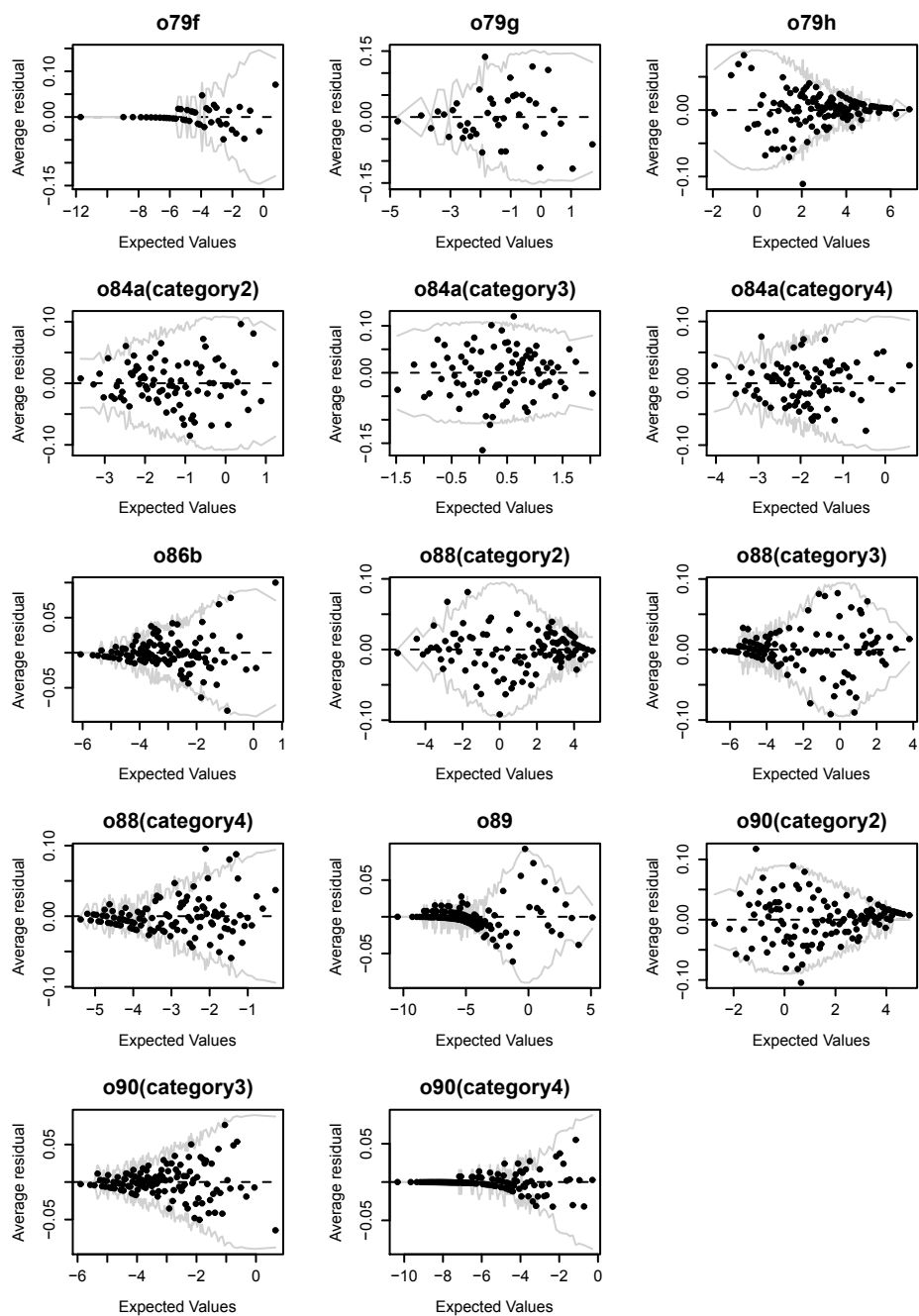


Figure A.6: Binned residual plots for the categorical variables with missing rates above 1%.

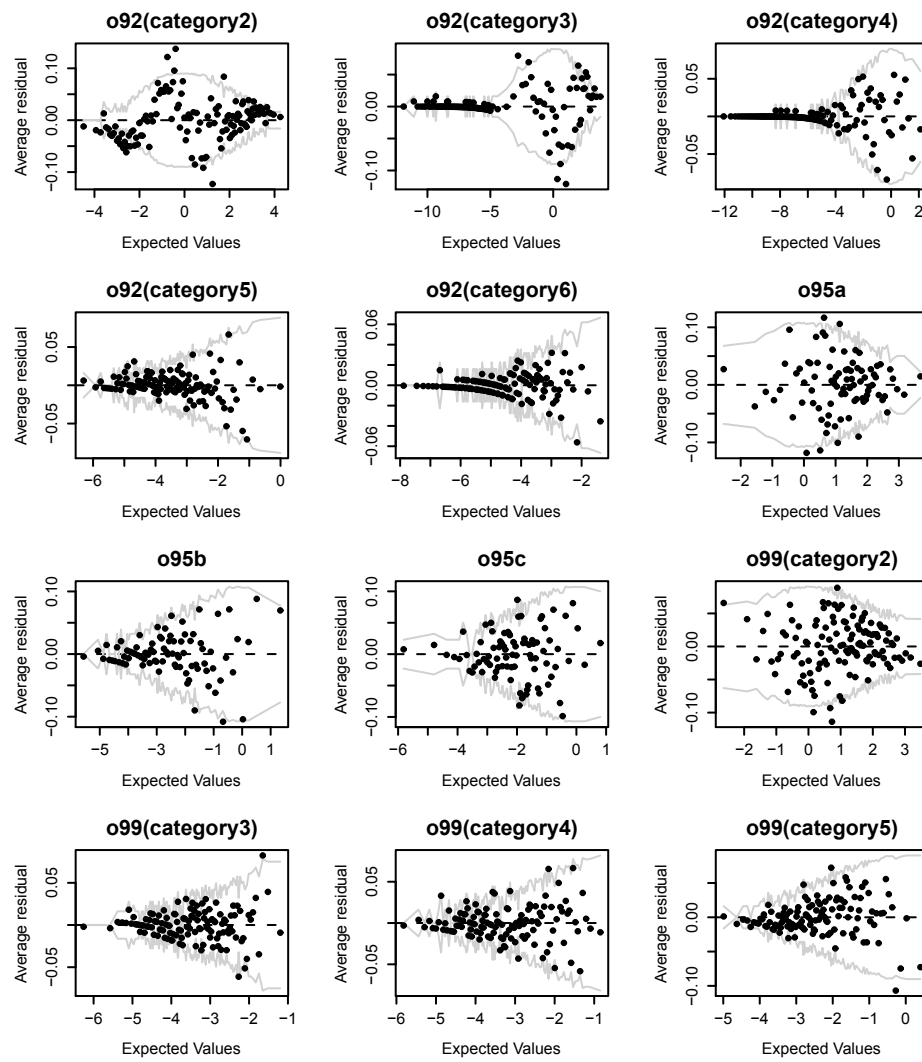


Figure A.7: Binned residual plots for the categorical variables with missing rates above 1%.

A.3 Simulation study for the variance inflated imputation model

Here we present results from a small simulation study that we conducted to evaluate the impact on data quality for the variance inflated imputation model described in Section 7.4.4.4. For the simulation, we generate a population of $N = 1,000,000$ records comprising three variables, Y_1, \dots, Y_3 , drawn from $N(0, \Sigma)$, where Σ has variances equal to one and correlations ranging from 0.3 to 0.7. From this population we repeatedly draw simple random samples of size $s = 10,000$ and treat these samples as the originally observed samples D_{obs} . For the synthesis we replace values of Y_3 for all records in D_{obs} . We generate replacement values by sampling from the posterior predictive distribution, $f(Y_3|D_{obs})$, using parameter values drawn from the variance inflated posterior distribution given in (7.9) with different levels of the variance inflation factor α . For comparison, we also generate synthetic datasets with Y_1 omitted from the imputation model to illustrate the negative consequences of dropping explanatory variables from the models to obtain a higher level of data protection. In analogy with our real data application, we generate $m = 5$ synthetic datasets for any one iteration of the simulation design. We obtain inferences for 4 quantities in each simulation run, including the population mean and the intercept and regression coefficients of Y_2 (β_1) and Y_3 (β_2) in a regression of Y_1 on Y_2 and Y_3 . We repeat each simulation 5,000 times.

Table A.1 displays key results from the simulations. The average \bar{q}_m across the 5,000 simulation runs is always very close to the average q_{obs} for $\alpha \leq 100$. For $\alpha = 1,000$ we find small biases for all point estimates. The variance estimator T_p (column four) correctly estimates the true variance of \bar{q}_m (column three) for any given level of α . Columns six and seven summarize the percentages of the 5,000 synthetic 95% confidence intervals that cover their corresponding Q for the original sample and the synthetic samples respectively. We find that the coverage rates from the synthetic samples are always close to the expected nominal coverage of 95% for $\alpha \leq 100$. Only for $\alpha = 100$ we notice a slight undercoverage for the regression coefficient β_2 compared to the coverage rate of β_2 in the original sample. The undercoverage increases for $\alpha = 1,000$. All estimates slightly undercover and for β_2 the coverage rate actually drops to 90.8%. The ninth column reports the ratio of the confidence interval length from the synthetic datasets over the confidence interval length from the original

Table A.1: Simulation results for the inflated variance imputation. The denominators of the confidence interval length ratios and the RMSE ratios are based on the point estimates from the sample without synthesis.

α	Q	q_{obs}	\bar{q}_m	$var(q_m)$	T_p	CI cov. org.	CI cov. syn.	CI length ratio	RMSE ratio	
1	$mean(Y_3)$	$-3.27 * 10^{-03}$	$-3.41 * 10^{-03}$	$3.65 * 10^{-03}$	$3.58 * 10^{-03}$	94.92	95.04	1.11	1.10	
	Intercept	$3.84 * 10^{-03}$	$3.93 * 10^{-03}$	$1.74 * 10^{-03}$	$1.76 * 10^{-03}$	95.54	95.26	1.10	1.10	
	β_1	$9.93 * 10^{-02}$	$9.93 * 10^{-02}$	$6.50 * 10^{-05}$	$6.40 * 10^{-05}$	94.78	94.96	1.10	1.08	
	β_2	$6.70 * 10^{-01}$	$6.70 * 10^{-01}$	$6.60 * 10^{-05}$	$6.50 * 10^{-05}$	95.26	94.96	1.11	1.11	
	10	$mean(Y_3)$	$-2.91 * 10^{-03}$	$-2.69 * 10^{-03}$	$6.23 * 10^{-03}$	$6.31 * 10^{-03}$	94.90	95.18	1.57	1.45
		Intercept	$3.79 * 10^{-03}$	$3.64 * 10^{-03}$	$2.97 * 10^{-03}$	$2.99 * 10^{-03}$	95.42	94.64	1.52	1.43
β_1		$9.93 * 10^{-02}$	$9.94 * 10^{-02}$	$1.08 * 10^{-04}$	$1.05 * 10^{-04}$	94.78	94.30	1.48	1.39	
β_2		$6.70 * 10^{-01}$	$6.70 * 10^{-01}$	$6.80 * 10^{-05}$	$6.70 * 10^{-05}$	94.74	94.82	1.12	1.11	
100		$mean(Y_3)$	$-1.18 * 10^{-03}$	$-1.20 * 10^{-03}$	$3.29 * 10^{-02}$	$3.34 * 10^{-02}$	94.32	94.58	4.19	3.29
		Intercept	$3.39 * 10^{-03}$	$3.41 * 10^{-03}$	$1.50 * 10^{-02}$	$1.50 * 10^{-02}$	95.08	94.44	3.94	3.17
	β_1	$9.93 * 10^{-02}$	$1.01 * 10^{-01}$	$4.87 * 10^{-04}$	$5.07 * 10^{-04}$	94.90	94.44	3.78	2.98	
	β_2	$6.70 * 10^{-01}$	$6.65 * 10^{-01}$	$8.90 * 10^{-05}$	$9.10 * 10^{-05}$	95.58	94.52	1.33	1.45	
	1000	$mean(Y_3)$	$-8.98 * 10^{-04}$	$9.97 * 10^{-03}$	$3.11 * 10^{-01}$	$3.02 * 10^{-01}$	95.00	94.40	13.29	10.22
		Intercept	$3.14 * 10^{-03}$	$-3.20 * 10^{-03}$	$1.23 * 10^{-01}$	$1.20 * 10^{-01}$	95.18	94.62	11.78	9.21
β_1		$9.93 * 10^{-02}$	$1.14 * 10^{-01}$	$4.16 * 10^{-03}$	$4.20 * 10^{-03}$	95.00	94.08	11.53	8.92	
β_2		$6.70 * 10^{-01}$	$6.24 * 10^{-01}$	$1.66 * 10^{-03}$	$1.68 * 10^{-03}$	95.06	90.82	5.65	8.35	

samples. Not surprisingly, the ratio increases with increasing α , since the variance inflated imputation model increases the between imputation variance b_m and thus the variance of \bar{q}_m . Comparing the confidence interval length ratio with the root mean squared error (RMSE) ratio in the last column, we notice that the RMSE ratio is always smaller than or equal to the confidence interval length ratio indicating that the increased RMSE in the synthetic datasets is likely due to the increased variance from the variance inflated imputation model. Only for the regression coefficient β_2 and $\alpha \geq 100$ we find an increased RMSE ratio compared to the confidence interval length ratio. Overall we find that at least for this simulation levels of $\alpha \leq 100$ only lead to reduced efficiency in the estimation, but not to any noticeable bias. For $\alpha = 1,000$ we find a small bias that leads to slight undercoverage, but note that we replaced all records with variance inflated imputations in these simulations. In practice agencies will only replace some records that are specifically at risk with draws from the variance inflated imputation model. We expect that the bias will be small in this context.

The results for the data generation that drops Y_1 from the imputation model to obtain a higher level of data protection are presented in Table A.2. \bar{Y}_3 and the intercept from the regression are not affected, but the two regression coefficients are completely biased leading to a 0% coverage rate for both estimates and a significantly increased RMSE ratio. It is obvious that the variance inflated imputation model provides far better results in terms of data validity. Dropping variables from the imputation models should only be considered an option, if the data disseminating agency knows that the data user will never evaluate the relationship between the dropped variable and the variable to be imputed.

Table A.2: Simulation results if Y_1 is excluded from the imputation model. The denominators of the confidence interval length ratios and the RMSE ratios are based on the point estimates from the sample without synthesis.

Q	q_{obs}	\bar{q}_m	$var(q_m)$	T_p	CI cov. org.	CI cov. syn.	CI length ratio	RMSE ratio
$mean(Y_3)$	$-2.07 * 10^{-03}$	$-1.87 * 10^{-03}$	$4.12 * 10^{-03}$	$4.07 * 10^{-03}$	95.40	94.92	1.20	1.18
Intercept	$3.86 * 10^{-03}$	$2.51 * 10^{-03}$	$2.68 * 10^{-03}$	$2.70 * 10^{-03}$	94.88	95.08	1.35	1.34
β_1	$9.93 * 10^{-02}$	$3.00 * 10^{-01}$	$8.90 * 10^{-05}$	$1.01 * 10^{-04}$	95.10	0.00	1.36	27.15
β_2	$6.70 * 10^{-01}$	$-7.60 * 10^{-05}$	$2.00 * 10^{-05}$	$1.19 * 10^{-04}$	94.72	0.00	1.49	90.70

Bibliography

- Abayomi, K., Gelman, A., and Levy, M. (2008). Diagnostics for multivariate imputations. *Journal of the Royal Statistical Society, Series C* **57**, 273–291.
- Abowd, J. and Vilhuber, L. (2008). How protective are synthetic data? In J. Domingo-Ferrer and Y. Saygin, eds., *Privacy in Statistical Databases*, 239–246. New York: Springer-Verlag.
- Abowd, J. M. and Lane, J. I. (2004). New approaches to confidentiality protection: Synthetic data, remote access and research data centers. In J. Domingo-Ferrer and V. Torra, eds., *Privacy in Statistical Databases*, 282–289. New York: Springer-Verlag.
- Abowd, J. M., Stinson, M., and Benedetto, G. (2006). Final report to the social security administration on the SIPP/SSA/IRS public use file project. Tech. rep., U.S. Census Bureau Longitudinal Employer-Household Dynamics Program.
- Abowd, J. M. and Woodcock, S. D. (2001). Disclosure limitation in longitudinal linked data. In P. Doyle, J. Lane, L. Zayatz, and J. Theeuwes, eds., *Confidentiality, Disclosure, and Data Access: Theory and Practical Applications for Statistical Agencies*, 215–277. Amsterdam: North-Holland.
- Abowd, J. M. and Woodcock, S. D. (2004). Multiply-imputing confidential characteristics and file links in longitudinal linked data. In J. Domingo-Ferrer and V. Torra, eds., *Privacy in Statistical Databases*, 290–297. New York: Springer-Verlag.
- An, D. and Little, R. J. A. (2007). Multiple imputation: an alternative to top coding for statistical disclosure control. *Journal of the Royal Statistical Society, Series A* **170**, 923–940.

- Arnold, B. C., Castillo, E., and Sarabia, J. M. (1999). *Conditional Specification of Statistical Models*. Springer.
- Barnard, J. and Rubin, D. B. (1999). Small-sample degrees of freedom with multiple-imputation. *Biometrika* **86**, 948–955.
- Center of Excellence for Statistical Disclosure Control (2009). Handbook on statistical disclosure control. Available at: http://neon.vb.cbs.nl/casc/SDC_Handbook.pdf.
- Domingo-Ferrer, J., Drechsler, J., and Polettini, S. (2009). Report on synthetic data files. Tech. rep., Eurostat.
- Drechsler, J. (2009). Far from normal – multiple imputation of missing values in a German establishment survey. In *Proceedings of the UN/ECE Work Session on Statistical Data Editing and Imputation*, Available at <http://www.unece.org/stats/documents/ece/ces/ge.44/2009/wp.21.e.pdf>.
- Drechsler, J., Bender, S., and Rässler, S. (2008a). Comparing fully and partially synthetic data sets for statistical disclosure control in the German IAB Establishment Panel. *Transactions on Data Privacy* **1**, 105 – 130.
- Drechsler, J., Dundler, A., Bender, S., Rässler, S., and Zwick, T. (2008b). A new approach for disclosure control in the IAB Establishment Panel–Multiple imputation for a better data access. *Advances in Statistical Analysis* **92**, 439 – 458.
- Drechsler, J. and Rässler, S. (2008). Does convergence really matter? In Shalabh and C. Heumann, eds., *Recent Advances in Linear Models and Related Areas*, 341–355. Heidelberg: Physica.
- Drechsler, J. and Reiter, J. P. (2008). Accounting for intruder uncertainty due to sampling when estimating identification disclosure risks in partially synthetic data. In J. Domingo-Ferrer and Y. Saygin, eds., *Privacy in Statistical Databases*, 227–238. New York: Springer-Verlag.
- Drechsler, J. and Reiter, J. P. (2009). Disclosure risk and data utility for partially synthetic data: An empirical study using the German IAB Establishment Survey. *Journal of Official Statistics* **25**, 589–603.

- Duncan, G. T. and Lambert, D. (1989). The risk of disclosure for microdata. *Journal of Business and Economic Statistics* **7**, 207–217.
- Dwork, C. (2006). Differential privacy. In M. Bugliesi, B. Preneel, V. Sassone, and I. Wegener, eds., *ICALP 2006*, 1–12. New York: Springer-Verlag.
- Elamir, E. and Skinner, C. J. (2006). Record level measures of disclosure risk for survey microdata. *Journal of Official Statistics* **22**, 525–539.
- Fienberg, S. E. (1994). A radical proposal for the provision of micro-data samples and the preservation of confidentiality. Tech. rep., Department of Statistics, Carnegie-Mellon University.
- Fienberg, S. E. and Makov, U. E. (1998). Confidentiality, uniqueness, and disclosure limitation for categorical data. *Journal of Official Statistics* **14**, 361–372.
- Fienberg, S. E., Makov, U. E., and Sanil, A. P. (1997). A Bayesian approach to data disclosure: Optimal intruder behavior for continuous data. *Journal of Official Statistics* **13**, 75–89.
- Fienberg, S. E., Makov, U. E., and Steele, R. J. (1998). Disclosure limitation using perturbation and related methods for categorical data. *Journal of Official Statistics* **14**, 485–502.
- Fischer, G., Janik, F., Müller, D., and Schmucker, A. (2008). The IAB Establishment Panel - from sample to survey to projection. Tech. rep., FDZ-Methodenreport, No. 1 (2008).
- Franconi, L. and Stander, J. (2002). A model based method for disclosure limitation of business microdata. *The Statistician* **51**, 1–11.
- Franconi, L. and Stander, J. (2003). Spatial and non-spatial model-based protection procedures for the release of business microdata. *Statistics and Computing* **13**, 295–305.
- Gelman, A., Carlin, J. B., Stern, H. S., and Rubin, D. B. (2004). *Bayesian Data Analysis: Second Edition*. London: Chapman & Hall.
- Gelman, A. and Hill, J. (2006). *Data Analysis Using Regression and Multi-level/Hierarchical Models*. Cambridge University Press.

- Graham, P. and Penny, R. (2005). Multiply imputed synthetic data files. Tech. rep., University of Otago, <http://www.uoc.otago.ac.nz/departments/pubhealth/pgrahpub.htm>.
- Graham, P., Young, J., and Penny, R. (2009). Multiply imputed synthetic data: Evaluation of hierarchical bayesian imputation models. *Journal of Official Statistics* **25**, 407–426.
- Karr, A. F., Kohnen, C. N., Oganian, A., Reiter, J. P., and Sanil, A. P. (2006). A framework for evaluating the utility of data altered to protect confidentiality. *The American Statistician* **60**, 224–232.
- Kennickell, A. B. (1997). Multiple imputation and disclosure protection: The case of the 1995 Survey of Consumer Finances. In W. Alvey and B. Jamer-son, eds., *Record Linkage Techniques, 1997*, 248–267. Washington, D.C.: National Academy Press.
- Kinney, S. K. and Reiter, J. P. (2010). Significance testing when using multiple imputation for missing data and disclosure limitation. *Journal of Official Statistics* forthcoming.
- Knoche, P. (1993). Factual anonymity of microdata from household and person-related surveys - the release of microdata files for scientific purposes. In *Proceedings of the International Symposium on Statistical Confidentiality*, 407–413.
- Kölling, A. (2000). The IAB-Establishment Panel. *Journal of Applied Social Science Studies* **120**, 291–300.
- Lane, J. I. (2007). Optimizing the use of microdata: An overview of the issues. *Journal of Official Statistics* **23**, 299–317.
- Lechner, S. and Pohlmeier, W. (2005). Data masking by noise addition and the estimation of nonparametric regression models. *Journal of Economics and Statistics* **225**, 517–528.
- Little, R. J. A. (1993). Statistical analysis of masked data. *Journal of Official Statistics* **9**, 407–426.

- Little, R. J. A., Liu, F., and Raghunathan, T. E. (2004). Statistical disclosure techniques based on multiple imputation. In A. Gelman and X. L. Meng, eds., *Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives*, 141–152. New York: John Wiley & Sons.
- Little, R. J. A. and Raghunathan, T. E. (1997). Should imputation of missing data condition on all observed variables? In *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, 617–622.
- Little, R. J. A. and Rubin, D. B. (2002). *Statistical Analysis with Missing Data: Second Edition*. New York: John Wiley & Sons.
- Liu, F. and Little, R. J. A. (2002). Selective multiple imputation of keys for statistical disclosure control in microdata. In *ASA Proceedings of the Joint Statistical Meetings*, 2133–2138.
- Machanavajjhala, A., Kifer, D., Abowd, J. M., Gehrke, J., and Vilhuber, L. (2008). Privacy: Theory meets practice on the map. In *ICDE*, 277–286.
- Meng, X.-L. (1994). Multiple-imputation inferences with uncongenial sources of input (disc: P558-573). *Statistical Science* **9**, 538–558.
- Mitra, R. and Reiter, J. P. (2006). Adjusting survey weights when altering identifying design variables via synthetic data. In J. Domingo-Ferrer and L. Franconi, eds., *Privacy in Statistical Databases*, 177–188. New York: Springer-Verlag.
- Polettini, S. (2003). Maximum entropy simulation for microdata protection. *Statistics and Computing* **13**, 308–320.
- Polettini, S., Franconi, L., and Stander, J. (2002). Model-based disclosure protection. In J. Domingo-Ferrer, ed., *Inference Control in Statistical Databases*, 83–96. Berlin: Springer-Verlag.
- Raghunathan, T. E., Lepkowski, J. M., van Hoewyk, J., and Solenberger, P. (2001). A multivariate technique for multiply imputing missing values using a series of regression models. *Survey Methodology* **27**, 85–96.
- Raghunathan, T. E., Reiter, J. P., and Rubin, D. B. (2003). Multiple imputation for statistical disclosure limitation. *Journal of Official Statistics* **19**, 1–16.

- Raghunathan, T. E., Solenberger, P., and van Hoewyk, J. (2002). IVEware: Imputation and variance estimation software. Available at: <http://www.isr.umich.edu/src/smp/ive/>.
- Reiter, J. P. (2002). Satisfying disclosure restrictions with synthetic data sets. *Journal of Official Statistics* **18**, 531–544.
- Reiter, J. P. (2003). Inference for partially synthetic, public use microdata sets. *Survey Methodology* **29**, 181–189.
- Reiter, J. P. (2004). Simultaneous use of multiple imputation for missing data and disclosure limitation. *Survey Methodology* **30**, 235–242.
- Reiter, J. P. (2005a). Estimating identification risks in microdata. *Journal of the American Statistical Association* **100**, 1103–1113.
- Reiter, J. P. (2005b). Releasing multiply-imputed, synthetic public use microdata: An illustration and empirical study. *Journal of the Royal Statistical Society, Series A* **168**, 185–205.
- Reiter, J. P. (2005c). Significance tests for multi-component estimands from multiply-imputed, synthetic microdata. *Journal of Statistical Planning and Inference* **131**, 365–377.
- Reiter, J. P. (2005d). Using CART to generate partially synthetic, public use microdata. *Journal of Official Statistics* **21**, 441–462.
- Reiter, J. P. (2008a). Letter to the editor. *Journal of Official Statistics* **24**, 319–321.
- Reiter, J. P. (2008b). Selecting the number of imputed datasets when using multiple imputation for missing data and disclosure limitation. *Statistics and Probability Letters* **78**, 15–20.
- Reiter, J. P. and Drechsler, J. (2010). Releasing multiply-imputed, synthetic data generated in two stages to protect confidentiality. *Statistica Sinica* **20**, 405–421.
- Reiter, J. P. and Mitra, R. (2009). Estimating risks of identification disclosure in partially synthetic data. *Journal of Privacy and Confidentiality* **1**, 99–110.

- Rosenbaum, P. R. and Rubin, D. B. (1983). The central role of the propensity score in observational studies for causal effects. *Biometrika* **70**, 41–55.
- Royston, P. (2005). Multiple imputation of missing values: Update of ice. *The Stata Journal* **5**, 527–536.
- Rubin, D. B. (1978). Multiple imputations in sample surveys. In *Proceedings of the Section on Survey Research Methods of the American Statistical Association*, 20–34.
- Rubin, D. B. (1981). The Bayesian bootstrap. *The Annals of Statistics* **9**, 130–134.
- Rubin, D. B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York: John Wiley & Sons.
- Rubin, D. B. (1993). Discussion: Statistical disclosure limitation. *Journal of Official Statistics* **9**, 462–468.
- Rubin, D. B. (2004). The design of a general and flexible system for handling nonresponse in sample surveys. *The American Statistician* **58**, 298–302.
- Rubin, D. B. and Schenker, N. (1986). Multiple imputation for interval estimation from simple random samples with ignorable nonresponse. *Journal of the American Statistical Association* **81**, 366–374.
- Schafer, J. L. (1997a). *Analysis of Incomplete Multivariate Data*. Chapman & Hall.
- Schafer, J. L. (1997b). *Analysis of Incomplete Multivariate Data*. London: Chapman & Hall.
- Schenker, N., Raghunathan, T. E., Chiu, P. L., Makuc, D. M., Zhang, G., and Cohen, A. J. (2006). Multiple imputation of missing income data in the National Health Interview Survey. *Journal of the American Statistical Association* **101**, 924–933.
- Statistisches Bundesamt (2005). Statistik und Wissenschaft Band 4: Handbuch zur Anonymisierung wirtschaftsstatistischer Mikrodaten.

- Su, Y., Gelman, A., Hill, J., and Yajima, M. (2009). Multiple imputation with diagnostics (mi) in R: Opening windows into the black box. *Journal of Statistical Software* (forthcoming).
- Van Buuren, S. and Oudshoorn, C. (2000). Mice v1.0 user’s manual. report pg/vgz/00.038. Tech. rep., TNO Prevention and Health, Leiden.
- Winkler, W. E. (2007a). Analytically valid discrete microdata files and re-identification. Tech. rep., Statistical Research Division, U.S. Bureau of the Census.
- Winkler, W. E. (2007b). Examples of easy-to-implement, widely used methods of masking for which analytic properties are not justified. Tech. rep., Statistical Research Division, U.S. Bureau of the Census.
- Woo, M. J., Reiter, J. P., Oganian, A., and Karr, A. F. (2009). Global measures of data utility for microdata masked for disclosure limitation. *Journal of Privacy and Confidentiality* **1**, 111–124.
- Woodcock, S. D. and Benedetto, G. (2009). Distribution-preserving statistical disclosure limitation. *Computational Statistics and Data Analysis* **53**, 4228–4242.
- Zwick, T. (2005). Continuing vocational training forms and establishment productivity in germany. *German Economic Review* **6**, 155 – 184.