# An Analysis of the Missing Data Methodology
# for Different Types of Data

A THESIS PRESENTED IN PARTIAL FULFILMENT
OF THE REQUIREMENTS FOR THE DEGREE
OF

MASTER OF APPLIED STATISTICS

AT MASSEY UNIVERSITY, ALBANY
NEW ZEALAND

Judith-Anne Scheffer

2000

# Abstract

Missing data is an eternal problem in data analysis. It is widely recognised that data is costly to collect, and the methods used to deal with missing data in the past relied on case deletion. There is no one overall best fix, but many different methodologies to use in different situations.

This study was motivated by the writer's time spent analysing data in the nutrition study, and realising how much data was wasted by case deletion, and subsequently how this could bias inferences formed from the results. A better method (or methods), of dealing with missing data (than case deletion) is required, to ensure valuable information is not lost.

What is being done: What is in the literature? The literature on this topic has exploded with new methods in recent times. Algorithms have been written and incorporated based on these methods into a number of statistical packages and add-on libraries.

Statistical packages are also reviewed for their practicality and application in this area. The nutrition data is then applied to different methodologies, and software packages to assess different types of imputation.

A set of questions are posed; based on type of data, type of missingness, extent of missingness, the required end use of the data, the size of the dataset, and how extensive that analysis needs to be. This can guide the investigator into using an appropriate form of imputation for the type of data at hand.

A comparison of imputation methods and results is given with the principal result that imputing missing data is a very worthwhile exercise to reduce bias in survey results, which can be achieved by any researcher analysing their own data.

Further to this, a conjecture is given for using Data Augmentation for ordinal data, particularly Likert scales. Previously this has been restricted to either person or item mean imputation, or hot deck methods. Using model based methods for imputation is far superior for other types of data. Model based methods for Likert data are achieved by means of inserting the linear by linear association model into standard missing data methodology.

# Acknowledgements

I wish to offer my sincerest thanks to my supervisor, Doctor Barry W. M$^c$Donald, for all his helpful advice, comments and efforts on my behalf, and also for his encouragement and mentoring throughout the course of this degree.

My thanks also go to Doctor Howard P. Edwards for his assistance in 'Matters Bayesian', Ms Katya Ruggiero for her ability to challenge practices and ideas, Mrs Kay Rowbottom for her assistance with the production of the flowcharts, and Synthia for her encouragement.

Thanks also go to Mrs Patsy E. Watson for providing via my supervisor, the nutrition dataset; and also to Ms Janet Norton for providing her dataset, via Professor Graham R. Wood.

Lastly but not least, I would like to thank my family (the thesis orphans) for putting up with my frequent absences for long periods to do this work.

Blessed is the man who perseveres under trial,
because when he has stood the test,
he will receive the crown of life that
God has promised to those who love him.
James 1:12

# Table of Contents

# 3   LITERATURE REVIEW OF METHODOLOGY FOR ANALYSING MISSING DATA                                                                    32

# 8    SOME APPROACHES TO ORDINAL  CATEGORICAL DATA IMPUTATION: LIKERT DATA IN PARTICULAR  (A CONJECTURE)

# 9    ANALYSIS AND IMPUTATION OF DATA    157

# List of Tables and Figures

# Notation and Abbreviations

| | |
|---|---|
| BLR | Binary Logistic Regression |
| CD | Case Deletion |
| EM | Expectation Maximisation (algorithm) |
| EM Imp | Imputation via the EM algorithm |
| GLM Imp | General Location Model Imputation |
| HD | Hotdeck (Imputation) |
| iid | Independent identically distributed |
| LUM | Look up methods |
| LVCF | Last Value Carried Forwards |
| MCAR | Missing Completely at Random |
| MAR | Missing at Random |
| Mean Imp | Mean family of Imputation |
| MI | Multiple Imputation |
| MI BB | Multiple Imputation Bayesian Bootstrap |
| MICE | Multiple Imputation by Chained Equations |
| MI DA | Multiple Imputation via Data Augmentation |
| MI EM | Multiple Imputation via the EM algorithm |
| N.Neighbour | Nearest Neighbour |
| N Nets | Neural Networks |
| NLR | Nominal Logistic Regression |
| NMAR | Not Missing at Random (Informatively Missing) |
| OLR | Ordinal Logistic Regression |
| PMM | Predictive Mean matching |
| Reg Imp | Regression Imputation |
| SHHD | Sequential and/or Hierarchical Hotdeck |
| SI | Single Imputation |
| St Reg | Stochastic regression Imputation |

| | |
|---|---|
| W | Indicator for Missingness |
| X | Co-variate in model |
| Y | Variable of interest |
| | |
| $\hat{\alpha}$ | Gamma Parameter (Ch 8) |
| $\hat{\beta}$ | Gamma Parameter (Ch 8) |
| $\hat{\beta}$ | Regression Coefficient Estimate (Ch 9) |
| $\theta$ | Distribution Parameter |
| $\hat{\theta}$ | Maximum Likelihood Estimate of the Parameter |
| $\psi$ | Missingness Parameter in Model |