

COMBINING LATENT CLASS ANALYSIS AND MULTIPLE IMPUTATION TO CORRECT FOR MISCLASSIFICATION IN COMBINED DATASETS

Laura Boeschoten

Tilburg University – Statistics Netherlands

18 oktober 2018

OUTLINE

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
2/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- 1 Introduction
- 2 MILC
- 3 Application
- 4 Extensions
- 5 Application
- 6 Summary & Discussion

Starting point: Combined dataset

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
3/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Multiple sources:

- Administrative sources
 - (Large part of) population
- Surveys
 - Sample of population
- Linked on person level

Different sources sometimes contain the same (categorical) variable

Official Statistics output

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
4/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Predefined large cross-tables that meet a number of requirements

Requirement 1

Account for impossible combinations of scores

Problem:

- Sometimes, a combination of scores is observed that is not possible in practice
- Caused by misclassification in one of the variables (De Waal, Pannekoek & Scholtus, 2011)

Solution:

- Incorporate edit restrictions into the LC model

Example:

$$P(\text{Gender} = \text{Male} | \text{Pregnant} = \text{Yes}) = 0$$

Requirement 2

Numerical consistence over different cross-tables

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
6/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Problem:

- Sometimes variables are only observed by means of sample surveys
- Weighting leads to inconsistent estimates of the different cross-tables (De Waal, 2014)

Solution:

- Mass imputation

Requirement 3

Incorporate uncertainty into the variance estimates

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
7/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Problem:

- Missing and conflicting values make us more uncertain about our estimates

Solution:

- Multiple imputation

Summary

Multiple Imputation of Latent Classes (MILC)

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
8/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Latent Class model

- Variables that measure the same construct as indicators of a latent variable that measures the ‘true scores’
- Edit restrictions

Multiple Imputation

- Consistent estimates
- Uncertainty due to missing and conflicting values

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
9/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

MULTIPELE IMPUTATION OF LATENT CLASSES

MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
10/34

- BigSurv -
17/10/2018

Introduction

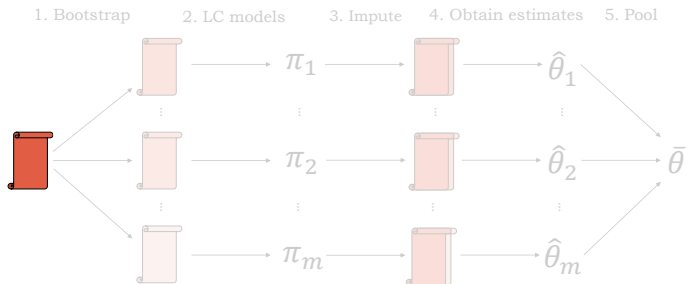
MILC

Application

Extensions

Application

Summary &
Discussion



MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
11/34

- BigSurv -
17/10/2018

Introduction

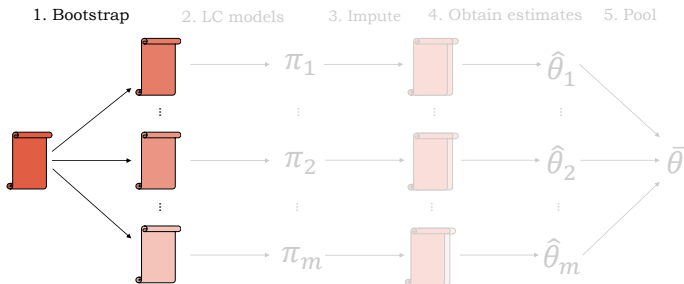
MILC

Application

Extensions

Application

Summary &
Discussion



MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
12/34

- BigSurv -
17/10/2018

Introduction

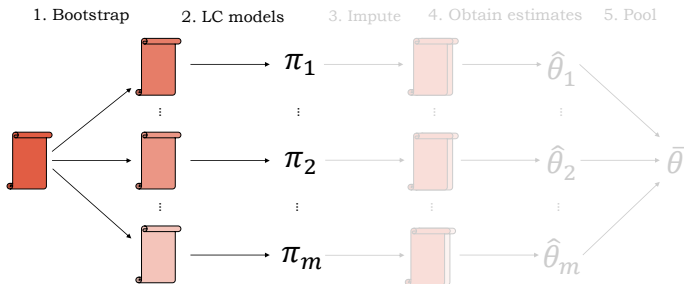
MILC

Application

Extensions

Application

Summary &
Discussion



LC MODEL

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
13/34

- BigSurv -
17/10/2018

Introduction

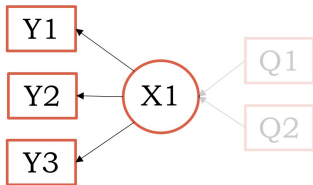
MILC

Application

Extensions

Application

Summary &
Discussion



- Variables measuring the same construct used as indicators of a latent variable
- Number of LC's = number of categories in indicators

LC MODEL

Assumptions

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
14/34

- BigSurv -
17/10/2018

Introduction

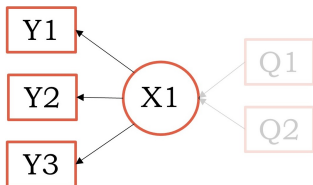
MILC

Application

Extensions

Application

Summary &
Discussion



- Mixture:

$$P(\mathbf{Y}) = \sum_x P(X)p(\mathbf{Y}|X)$$

- Local independence

$$P(\mathbf{Y}|X) = \prod_l P(Y_l|X)$$

LC MODEL

Covariates

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
15/34

- BigSurv -
17/10/2018

Introduction

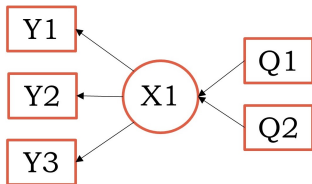
MILC

Application

Extensions

Application

Summary &
Discussion



- Misclassification independent of covariates
- Covariates are free of error
- Edit restrictions are ‘hard edits’:

$$P(X = 2|Q = 1) = 0$$

MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
16/34

- BigSurv -
17/10/2018

Introduction

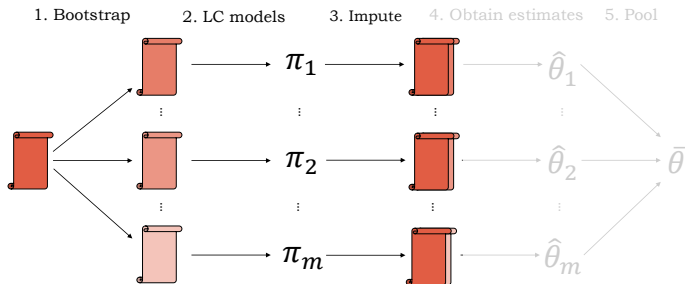
MILC

Application

Extensions

Application

Summary &
Discussion



MULTIPLE IMPUTATION

Using posteriors

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
17/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- For every case, one imputation is created using the posterior belonging to the corresponding profile

$$P(X|\mathbf{Y}) = \frac{P(X) \prod_l P(Y_l|X)}{\sum_x P(X) \prod_l P(Y_l|X)}$$

- This is done for every LC model (belonging to every bootstrap sample)
- Resulting in m imputations

MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
18/34

- BigSurv -
17/10/2018

Introduction

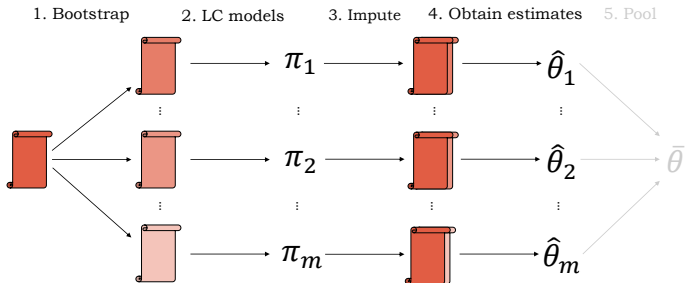
MILC

Application

Extensions

Application

Summary &
Discussion



MILC

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
19/34

- BigSurv -
17/10/2018

Introduction

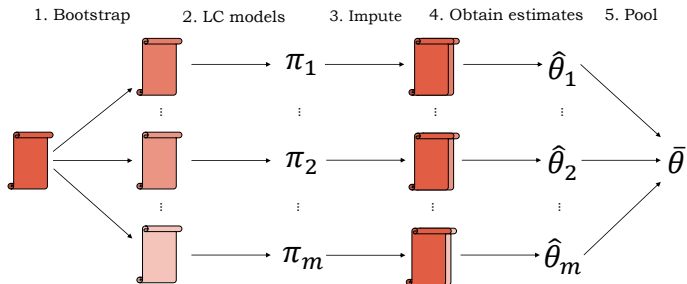
MILC

Application

Extensions

Application

Summary &
Discussion



POOLING OF THE RESULTS

Using the pooling rules defines by Rubin

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
20/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

$$T = \bar{U} + \left(1 + \frac{1}{m}\right)B$$

- T = Total variance
- \bar{U} = 'Within' variance (uncertainty about the assignend score)
- B = 'Between' variance (uncertainty about the model)

SIMULATION CONCLUSIONS

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
21/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- Quality of the imputations depends on the quality of the data
 - Entropy R^2 : Score between 0 and 1 indicating how well you can classify based on your observed data
- Required quality depends on desired output
- Low number of imputations seems sufficient
- (Boeschoten, Oberski & De Waal, 2017)

APPLICATION: NUMBER OF SERIOUS ROAD INJURIES PER VEHICLETYPE

Number of serious road injuries per vehicle type in 2013

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
23/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Police	Hospital
2,615 observed	189 missing
2,426 observed	In both sources
12,418 missing	14,844 observed

In total 15,033 serious road injuries

Number of serious road injuries per vehicle type in 2013

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
24/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- Nine categories for vehicle type
- Currently, if the police assigned a score, this is used. Otherwise, the score assigned by the hospital is used
- Vehicle type is classified differently for 30% of the cases

Number of serious road injuries per vehicle type in 2013

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
25/34

- BigSurv -
17/10/2018

Introduction

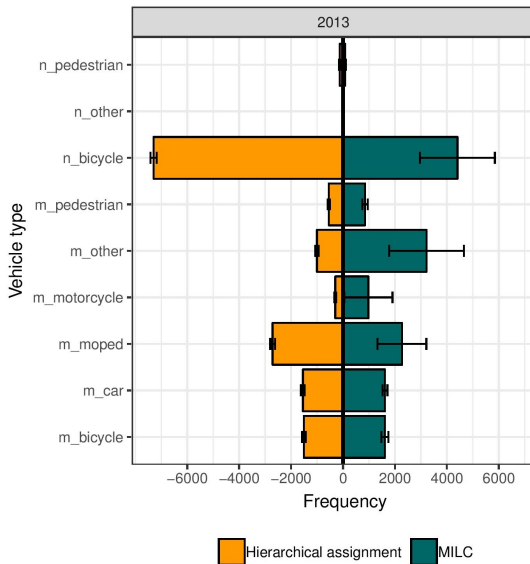
MILC

Application

Extensions

Application

Summary &
Discussion



Simultaneously impute missing values

By using a quasi-latent variable

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
26/34

- BigSurv -
17/10/2018

Introduction

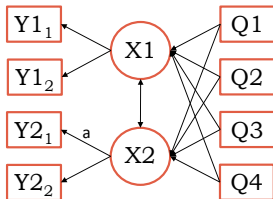
MILC

Application

Extensions

Application

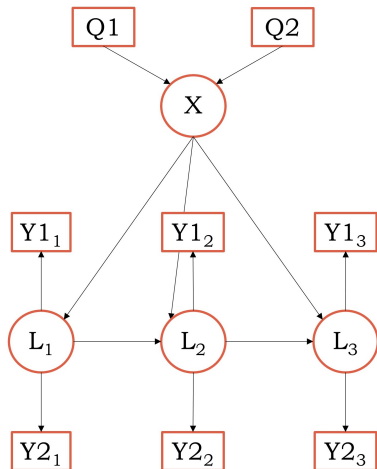
Summary &
Discussion



- Simultaneously impute the variable 'region of accident'
- Only measured by police
- restriction 'a': 'region of accident is a perfect indicator of X_2 '.
- If missing: LC model is used, with 'region of hospital' as another indicator

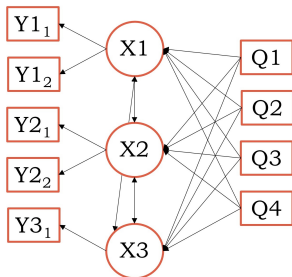
EXTENSIONS OF THE METHOD

Extension 1: Longitudinal



- Model to estimate monthly employment rates (Pavlopoulos & Vermunt, 2015)
- Create monthly imputations of employment status
- Investigated and applied in collaboration with ISTAT

Extension 2: Population census



- Comparable to model discussed in application
- Simulate from finite population
 - Different approach for bootstrap
- Many observations; many cells

Extension 3: External covariates

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
30/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

Using three-step methodology

- Estimate the relationship between an imputed latent variable and external covariates (not included in the LC model)
- Apply a correction procedure (ML or BCH)
- Update the imputations with the newly obtained posteriors
- (Boeschoten, Oberski, De Waal & Vermunt, 2018) & (Boeschoten, Croon & Oberski, 2018)

Summary

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
31/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- Quality of the output depends on the quality of the data
- Missing values can be imputed simultaneously
- MILC can easily be adjusted to specific situations
- Covariates can be added at a later time-point

Discussion

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
32/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- Is there a more flexible alternative for the bootstrap?
 - Gibbs sampler?
- How can we use MILC for other types of errors?
 - Or use it in combination with other types of correction methods?

Literature

MULTIPLE
IMPUTATION
OF LATENT
CLASSES
33/34

- BigSurv -
17/10/2018

Introduction

MILC

Application

Extensions

Application

Summary &
Discussion

- A.G. De Waal, J. Pannekoek, S. Scholtus (2011) Handbook of statistical data editing and Imputation *John Wiley & Sons Inc.*
- A.G. De Waal (2014) Consistent estimates for categorical data based on a mix of administrative data sources and surveys *CBS Discussion Paper*
- L. Boeschoten, D.L. Oberski & A.G. De Waal (2017) Estimating classification errors under edit restrictions in composite survey-register data using multiple imputation latent class modelling *Journal of Official Statistics* 33-4
- L. Boeschoten, D.L. Oberski, A.G. De Waal & J.K. Vermunt (2018) Updating latent class imputations with external auxiliary variables *Structural Equation Modeling: A Multidisciplinary Journal*
- L. Boeschoten, M. Croon & D.L. Oberski (2018) A note on applying the BCH method under linear equality and inequality constraints *Journal of Classification*
- D. Pavlopoulos & J.K. Vermunt (2015) Measuring temporary employment. Do survey or register data tell the truth? *Survey Methodology* 41-1, p.197-214

THANK YOU!!

`l.boeschoten@tilburguniversity.edu`