

# Imputation for the Treatment of Item Nonresponse in Surveys: An Overview

David Haziza

Department of mathematics and statistics  
Université de Montréal

ITACOS 2019  
Florence, Italy

June 5, 2019

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks

# Levels of nonresponse

- It is customary to distinguish **unit nonresponse** from **item nonresponse**:
  - 1 **Unit nonresponse**:
    - ▶ No usable information is collected on a sample unit;
    - ▶ Reasons: refusal, inability to contact the sample unit, etc.
  - 2 **Item nonresponse**:
    - ▶ Absence of information limited to some survey variables only;
    - ▶ Reasons: sensitive items (e.g., income), invalid or inconsistent responses.

# Methods of treatment

- **Treatment of unit nonresponse:** typically, handled through weight adjustment procedures methods. They consist of
  - eliminating the nonrespondents from the data file;
  - **increasing the sampling weight** of the respondents to compensate for the elimination of the nonrespondents;
- **Treatment of item nonresponse: Imputation!**
  - It consists of consists of constructing **one or more replacement values** to "fill in" for a missing value.
  - **Single imputation:** replace a missing value by a single imputed value → a single completed data file.
  - **Multiple imputation:** replace a missing value by  $M \geq 2$  replacement values →  $M$  completed data files.

# Effects of nonresponse

- Why is nonresponse an issue?
  - **Nonresponse bias**: due to the fact that respondents and nonrespondents do not have the same characteristics with respect to the survey variables.
  - Sample size is smaller than expected → **Variance of estimators is greater** than that of estimators that would have been used if there were no nonresponse. → **nonresponse variance**.
- Main objective of all the treatment methods: **reduce the nonresponse bias** and possibly control the variance due to nonresponse.
- **Some review papers on imputation**: Durrant (2005), Haziza (2009), Andridge and Little (2010) and **Chen and Haziza (2019)**.

# Questions of interest in the last two decades

- (1) How to obtain asymptotically unbiased and efficient point estimators?
- (2) How to obtain some protection against the misspecification of the underlying model(s)?
- (3) How to consistently estimate the variance of imputed estimators?
- (4) How to preserve relationships between survey variables?

# Finite population parameters

- Let  $\mathcal{P}$  be a finite population of size  $N$ .
- $y$ : a survey variable
- $y_i$ :  $y$ -value attached to unit  $i$ ,  $i = 1, \dots, N$ .
- **Goal**: estimate a finite population parameter  $\theta_N$  defined as the solution of **the census estimating equation**:

$$U_N(\theta_N) = \frac{1}{N} \sum_{i \in \mathcal{P}} u(y_i; \theta_N) = 0.$$

Parameter	$u(y_i; \theta_N)$	Explicit form of $\theta_N$
Total	$y_i - n^{-1} \pi_i \theta_N$	$t_y = \sum_{i \in \mathcal{P}} y_i$
Mean	$y_i - \theta_N$	$\bar{Y} = t_y / N$
Distribution function	$1(y_i \leq t) - \theta_N$	$F_N(t) = \frac{1}{N} \sum_{i \in \mathcal{P}} 1(y_i \leq t)$



# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks

# Estimation in ideal conditions

- Ideal conditions: no nonsampling errors
- $S$ : sample of size  $n$  selected according to a given sampling design  $P(S)$ .
- $\pi_i = P(i \in S) > 0$ : first-order inclusion probability for unit  $i$  (known prior to sampling)
- $w_i = 1/\pi_i$ : design weight attached to unit  $i$
- A full sample estimator of  $\theta_N$ , denoted by  $\hat{\theta}_F$ , is the solution of the sample estimating equation

$$\hat{U}_F(\theta_N) = \frac{1}{\hat{N}} \sum_{i \in S} w_i u(y_i; \theta_N) = 0.$$

# Estimation in ideal conditions

Parameter	$u(y_i; \theta_N)$	Full sample estimator $\hat{\theta}_F$
Total	$y_i - n^{-1}\pi_i\theta_N$	$\hat{t}_{HT} = \sum_{i \in S} w_i y_i$
Mean	$y_i - \theta_N$	$\hat{Y}_{HA} = \sum_{i \in S} w_i y_i / \sum_{i \in S} w_i$
Distribution function	$1(y_i \leq t) - \theta_N$	$\hat{F}_n(t) = \frac{\sum_{i \in S} w_i 1(y_i \leq t)}{\sum_{i \in S} w_i}$

- Under ideal conditions,  $\hat{\theta}_F$  is **design-consistent** for  $\theta_N$ ; that is,

$$\hat{\theta}_F - \theta_N = O_p(1/\sqrt{n}) \quad \text{or} \quad O_p(N/\sqrt{n});$$

e.g., Wang and Opsomer (2011) and Breidt and Opsomer (2017).

# Estimation after imputation

- In practice the  $y$ -variable is subject to missingness.
- Let  $r_i$  be a response indicator for unit  $i$  such that  $r_i = 1$  if  $y$  is observed and  $r_i = 0$ , otherwise.
- Let  $\tilde{y}$  denote the  $y$ -variable after imputation. We have

$$\tilde{y}_i = r_i y_i + (1 - r_i) y_i^*,$$

where  $y_i^*$  denotes the imputed value used to replace the missing  $y_i$ .

- An estimator of  $\theta_N$  after imputation, denoted by  $\hat{\theta}_I$ , is obtained by solving the estimating equation

$$\hat{U}_I(\theta_N) = \frac{1}{N} \sum_{i \in S} w_i u(\tilde{y}_i; \theta_N) = 0.$$

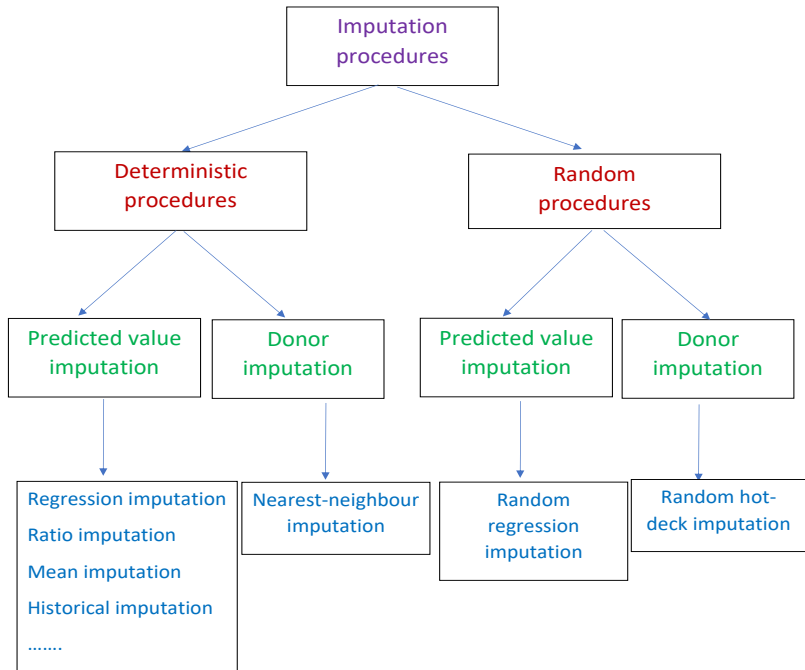
# Estimation after imputation

Parameter	$u(y_i; \theta_N)$	Imputed estimator $\hat{\theta}_I$
Total	$y_i - n^{-1}\pi_i\theta_N$	$\hat{t}_I = \sum_{i \in S} w_i \tilde{y}_i$
Mean	$y_i - \theta_N$	$\hat{\bar{Y}}_I = \sum_{i \in S} w_i \tilde{y}_i / \sum_{i \in S} w_i$
Distribution function	$1(y_i \leq t) - \theta_N$	$\hat{F}_I(t) = \frac{\sum_{i \in S} w_i 1(\tilde{y}_i \leq t)}{\sum_{i \in S} w_i}$

- Imputed estimators:
  - obtained by applying complete data point estimation procedures to  $\tilde{y}_i$  instead of  $y_i$ ;
  - attractive from a data user point of view.
- How to obtain the imputed values  $y_i^*$ ?

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures**
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks



# The imputation model

- We assume that the  $y$ -variable obeys the following model:

$$y_i = m(\mathbf{x}_i; \beta) + \epsilon_i, \quad (1)$$

where

- $m(\mathbf{x}_i; \cdot)$  is an unknown function;
  - $\mathbf{x}_i$ : vector of **fully observed variables** associated with unit  $i$ ;
  - $\mathbb{E}_m(\epsilon_i) = 0$ ,  $\mathbb{E}_m(\epsilon_i \epsilon_j) = 0$  for  $i \neq j$  and  $\mathbb{V}_m(\epsilon_i) = \sigma^2 c_i$ , where  $c_i$  is a **known coefficient** attached to unit  $i$ .
- Model (1) is called **an imputation model or an outcome regression model**.
- $\mathbb{E}_m(y_i | \mathbf{x}_i) = m(\mathbf{x}_i; \beta)$ : **first moment** of the imputation model.
- $\mathbb{V}_m(y_i | \mathbf{x}_i) = \sigma^2 c_i$ : **second moment** of the imputation model.
- No distributional assumptions about the  $\epsilon_i$ 's



# The nonresponse model

- We make the Missing At Random (MAR) assumption (Rubin, 1976):

$$\Pr(r_i = 1 \mid \mathbf{x}_i, y_i) = \Pr(r_i = 1 \mid \mathbf{x}_i) \equiv p(\mathbf{x}_i; \alpha), \quad (2)$$

where  $p(\mathbf{x}_i; \cdot)$  is an unknown function.

- Model (2) is called a **nonresponse model**.
- Consequence of the MAR assumption:

$$f(y \mid \mathbf{x}_i, r_i = 1) = f(y \mid \mathbf{x}_i, r_i = 0).$$

- Under this assumption, we can generate the imputed values from  $f(y \mid \mathbf{x}_i, r_i = 1)$ , which can be estimated from the observed data.

# Semi-parametric imputation procedures

- **Deterministic:**

$$y_i^* = m(\mathbf{x}_i; \hat{\beta}_r), \quad i \in S_m,$$

where  $\hat{\beta}_r$  is a suitable estimator (e.g., MLE) of  $\beta$  based on the respondents.

- **Special cases:**

- $m(\mathbf{x}_i; \hat{\beta}_r) = \mathbf{x}_i^\top \hat{\beta}_r \longrightarrow$  Regression imputation

- $\mathbf{x}_i = \mathbf{c}_i = \mathbf{1}$  for all  $i \longrightarrow m(\mathbf{x}_i; \hat{\beta}_r) = \hat{\beta}_r \equiv \hat{Y}_r \longrightarrow$  Mean imputation

- **Random:**

$$y_i^* = m(\mathbf{x}_i; \hat{\beta}_r) + \hat{\sigma} \sqrt{c_i} e_i^*, \quad i \in S_m,$$

where  $e_i^*$  is randomly drawn with replacement from the empirical distribution of **standardized residuals** observed among the respondents,  $\hat{F}_e(t)$ .

- **Special case:**

- Mean imputation + added random residuals  $\longrightarrow$  Random hot-deck imputation

# Nonparametric imputation procedures

- Nearest-neighbour imputation :

$$y_i^* = y_j,$$

where  $j$  is the index of the nearest-neighbour of unit  $i$ , which satisfies

$$D(\mathbf{x}_j, \mathbf{x}_i) \leq D(\mathbf{x}_k, \mathbf{x}_i), \quad k \in S_r$$

and  $D(\cdot; \cdot)$  denotes a distance function.

- Predictive mean matching : Same as NNI except that the information contained  $\mathbf{x}_i$  is compressed into a single score  $\hat{y}_i = m(\mathbf{x}_i; \hat{\beta}_r)$ ; see Little, (1988).
- Other nonparametric procedures:
  - Kernel methods: Zhong and Chen (2014)
  - Smoothing splines and additive models: Hasler and Craiu (2016)
  - B-splines: Goga and Haziza (2017)

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators**
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks

# Properties of imputed estimators

[◀ Back](#)

- Decomposition of the total error for deterministic imputation:

$$\hat{\theta}_I - \theta_N = \underbrace{(\hat{\theta}_F - \theta_N)}_{\text{Sampling error}} + \underbrace{(\hat{\theta}_I - \hat{\theta}_F)}_{\text{Nonresponse error}}$$

- Sampling error:** only affected by the population type, the sampling design, the sample size, etc.
- Nonresponse error:** affected by the nonresponse mechanism, the response rate, the quality of the imputation model etc.
- Decomposition of the total error for random imputation:

$$\hat{\theta}_I - \theta_N = \underbrace{(\hat{\theta}_F - \theta_N)}_{\text{Sampling error}} + \underbrace{(\check{\theta}_I - \hat{\theta}_F)}_{\text{Nonresponse error}} + \underbrace{(\hat{\theta}_I - \check{\theta}_I)}_{\text{Imputation error}}$$

- Imputation error:** purely artificial  $\rightarrow$  parasitic error

# Properties of imputed estimators

- Different inferential frameworks may be used:
  - Which quantities among  $y_i$ ,  $\mathbf{x}_i$ ,  $l_i$ ,  $r_i$  are treated as fixed/random?
  - Usual framework: *mpq* (assuming MAR)
- Population means and totals
  - Deterministic methods such as regression imputation:  
 $\widehat{\bar{Y}}_I - \bar{Y} = O_p(n^{-1/2})$  as long as the first moment of the imputation model is correctly specified; e.g., Chauvet, Deville and Haziza (2011).
  - NNI: we have  $\widehat{\bar{Y}}_I - \bar{Y} = O_p(n^{1/2-1/q})$ , where  $q$  is the size of  $\mathbf{x}$ ; see Yang and Kim (2017a,b)  $\rightarrow$  Bias is not negligible for  $q \geq 2 \rightarrow$  curse of dimensionality
  - PMM:  $\widehat{\bar{Y}}_I - \bar{Y} = O_p(n^{-1/2})$ ; see Yang and Kim (2017a,b).
  - Random methods such as random regression imputation: same as deterministic except that  $\widehat{\bar{Y}}_I$  suffers from an additional variability due to the random selection of the residuals  $e_i^* \rightarrow$  Imputation variance

# Properties of imputed estimators

- Distribution function and quantiles
  - Deterministic methods such as regression imputation: distort the distribution of the variable being imputed → biased estimators of quantiles
  - NNI: tend to preserve the distribution function; see Yang and Kim (2017a,b).
  - Random methods such as random regression imputation: tend to preserve the distribution function;

$$\hat{F}_I(t) - F_N(t) = O_p(n^{-1/2});$$

see Chen, Rao and Sitter (2000), Chauvet, Deville and Haziza (2011) and Boistard, Chauvet and Haziza (2016).

# Eliminating the imputation error

- Random imputation for population totals/means: Is it possible to preserve the distribution of the variable being imputed without paying the price of increased variance?
- **Fractional imputation:** e.g., Kalton and Kish (1981), Fay (1996), Kim and Fuller (2004) and **Yang and Kim (2016)**.
- **Balanced imputation:** e.g., Kalton and Kish (1981), Chauvet, Deville and Haziza (2011) and Haziza, Nambeu and Chauvet (2014).
  - **Idea:** Select the residuals  $e_i^*$  at random with replacement so that

$$\underbrace{(\hat{\theta}_I - \check{\theta}_I)}_{\text{Imputation error}} = 0.$$

- **One Option:** Using the **Cube algorithm** (Deville and Tillé, 2004)



# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation**
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks

# Variance estimation

- In the multiple imputation literature, single imputation is often criticized. The most common criticism is:

*The most obvious limitation of single imputation is the underlying assumption that the imputed value is the true value. This limitation leads to underestimation of the variance which affects confidence intervals and statistical tests.*

- There is a wide literature on variance estimation procedures for singly imputed data in a survey sampling setting, developed in the last two decades.

# Variance estimation: Two-phase framework

- $\mathcal{P} \longrightarrow S \longrightarrow S_r$
- Total variance of  $\hat{\theta}_I$  (Särndal, 1992): ► Decomposition

$$V_{\text{tot}} = V_{\text{sam}} + V_{\text{nr}} + V_{\text{mix}};$$

- For random imputation, add the extra variance due to imputation,  $V_{\text{imp}}$ .
- Obtaining consistent estimators of each term requires **the first two moments of the imputation model** to be correctly specified.
- **Donor imputation**: Brick, Kalton and Kim (2004);
- **Nearest-neighbour imputation**: Beaumont and Bocci (2009);
- **Historical imputation**: Beaumont, Haziza and Bocci (2011);
- **Composite imputation**: Beaumont and Bissonnette (2011).
- **Generalized system SEVANI**: Beaumont and Bissonnette (2011).

# Variance estimation: The reverse framework

- $\mathcal{P} \longrightarrow (\mathcal{P}_r, \mathcal{P}_m) \longrightarrow (S_r, S_m)$
- Total variance of  $\hat{\theta}_I$ :

$$V_{tot} = V_1 + V_2,$$

where

$$V_1 = E_q E_m V_p(\hat{\theta}_I - \theta_N), \quad V_2 = E_q V_m E_p(\hat{\theta}_I - \theta_N).$$

- The contribution of  $V_2$  to the total variance,  $V_2/(V_1 + V_2)$  is of order  $O(n/N) \longrightarrow$  negligible when the sampling fraction,  $n/N$ , is negligible;
- For random imputation, add the extra variance due to imputation,  $V_{imp}$ .
- e.g., Fay (1991), Shao and Steel (1999), Haziza (2009) and Kim and Rao (2009).

## Resampling methods

Adjusted jackknife (Rao and Shao, 1992; Rao and Sitter, 1995; Yung and Rao, 2000; Chen and Shao, 2001)

- Calculated in the usual way except that whenever the  $i$ -th unit is deleted, the imputed values,  $y_i^*$  are adjusted to reflect the fact that the set of respondents is changed.
- Consistent estimator of  $V_1$  provided that the sampling fraction  $n/N$  is negligible

Bootstrap (Shao and Sitter, 1996)

- Performed in the usual way, except that the nonrespondents are re-imputed within each bootstrap samples
- Consistent estimator of  $V_1$  provided that the sampling fraction  $n/N$  is negligible.
- **Bootstrap for non-negligible sampling fractions:** more challenging (Mashreghi, Léger and Haziza, 2014; Chen, Haziza, Léger and Mashreghi, 2019).

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation**
- 7 Multiple imputation
- 8 Final remarks

# Multiply robust imputation procedures

- So far, the imputation procedures were based on a single imputation model → Resulting estimators vulnerable to model misspecification
- To provide additional protection, it may be of interest of fitting multiple imputation models and/or multiple nonresponse models
- Each model may be based on different functionals and/or different sets of predictors.
- Idea: Develop imputation procedures that make use of multiple models (Chen and Haziza, 2017a)
- A procedure is said to be multiply robust if the resulting imputed estimator is consistent if all the models but one are misspecified.
- Concept introduced by Han and Wang (2013); see also Chan and Yam (2014)
- Can be viewed as an extension of double robustness

# Two classes of models

- Two classes of models:
  - Nonresponse models:

$$\mathcal{C}_1 = \{p^j(\mathbf{x}_i; \alpha^j); j = 1, \dots, J\}$$

- Imputation models:

$$\mathcal{C}_2 = \{m^k(\mathbf{x}_i; \beta^k); k = 1, \dots, K\}$$

- We fit each of the  $J + K$  models to obtain

$$\hat{\alpha}^1, \dots, \hat{\alpha}^J \longrightarrow p^1(\mathbf{x}_i; \hat{\alpha}^1), \dots, p^J(\mathbf{x}_i; \hat{\alpha}^J)$$

and

$$\hat{\beta}^1, \dots, \hat{\beta}^K \longrightarrow m^1(\mathbf{x}_i; \hat{\beta}^1), \dots, m^K(\mathbf{x}_i; \hat{\beta}^K)$$

- Usual estimators: least squares estimators, maximum likelihood estimators etc.



# Multiply robust imputation

	$p^1(\mathbf{x}; \hat{\mathbf{a}}^1)$	$\cdot$	$p^J(\mathbf{x}; \hat{\mathbf{a}}^J)$	$m^1(\mathbf{x}; \hat{\boldsymbol{\beta}}^1)$	$\cdot$	$m^K(\mathbf{x}; \hat{\boldsymbol{\beta}}^K)$	$r$	$\tilde{y}$
1	$p^1(\mathbf{x}_1; \hat{\mathbf{a}}^1)$	$\cdot$	$p^J(\mathbf{x}_1; \hat{\mathbf{a}}^J)$	$m^1(\mathbf{x}_1; \hat{\boldsymbol{\beta}}^1)$	$\cdot$	$m^K(\mathbf{x}_1; \hat{\boldsymbol{\beta}}^K)$	1	$y_1$
2	$p^1(\mathbf{x}_2; \hat{\mathbf{a}}^1)$	$\cdot$	$p^J(\mathbf{x}_2; \hat{\mathbf{a}}^J)$	$m^1(\mathbf{x}_2; \hat{\boldsymbol{\beta}}^1)$	$\cdot$	$m^K(\mathbf{x}_2; \hat{\boldsymbol{\beta}}^K)$	0	$y_2^*$
$\cdot$	$\cdot$			$\cdot$			$\cdot$	$\cdot$
$\cdot$	$\cdot$			$\cdot$			$\cdot$	$\cdot$
$\cdot$	$\cdot$			$\cdot$			$\cdot$	$\cdot$
$n$	$p^1(\mathbf{x}_n; \hat{\mathbf{a}}^1)$	$\cdot$	$p^J(\mathbf{x}_n; \hat{\mathbf{a}}^J)$	$m^1(\mathbf{x}_n; \hat{\boldsymbol{\beta}}^1)$	$\cdot$	$m^K(\mathbf{x}_n; \hat{\boldsymbol{\beta}}^K)$	0	$y_n^*$

# Construction of imputed values: Two steps

- Two-steps for constructing the imputed values  $y_i^*$ :
  - (1) First step: compress the information included in the  $J$  nonresponse models and the  $K$  imputation models
  - (2) Second step: **Implementation step** through linear regression
- **Note 1:** There are different ways to perform Step 1.
- **Note 2:**
  - $J = 0$  and  $K = 1$ : customary imputation based on a single imputation model.
  - $J = 1$  and  $K = 1$ : doubly robust imputation based on a single imputation model and a single nonresponse model.

# Compressing through refitting

- For unit  $i$ , define

$$\hat{\mathbf{Q}}_{pi} = \left( p^1(\mathbf{x}_i; \hat{\alpha}^1), \dots, p^J(\mathbf{x}_i; \hat{\alpha}^J) \right), \quad \hat{\mathbf{Q}}_{mi} = \left( m^1(\mathbf{x}_i; \hat{\beta}^1), \dots, m^K(\mathbf{x}_i; \hat{\beta}^K) \right).$$

- We summarize the working models information by regressing

$$r_i \text{ on } \hat{\mathbf{Q}}_{pi} \quad \text{and} \quad y_i \text{ on } \hat{\mathbf{Q}}_{mi}.$$

This leads to the weighted least square regression coefficients

$$\hat{\eta}_p = \left( \sum_{i \in S} w_i \hat{\mathbf{Q}}_{pi} \hat{\mathbf{Q}}_{pi}^\top \right)^{-1} \sum_{i \in S} w_i \hat{\mathbf{Q}}_{pi} r_i$$

and

$$\hat{\eta}_m = \left( \sum_{i \in S} w_i r_i \hat{\mathbf{Q}}_{mi} \hat{\mathbf{Q}}_{mi}^\top \right)^{-1} \sum_{i \in S} w_i r_i \hat{\mathbf{Q}}_{mi} y_i.$$

# Compressing through refitting

- For unit  $i$ , define

$$\hat{p}_i = \hat{\mathbf{Q}}_{pi}^{\top} \hat{\boldsymbol{\eta}}_p \quad \text{and} \quad \hat{m}_i = \hat{\mathbf{Q}}_{mi}^{\top} \hat{\boldsymbol{\eta}}_m.$$

- The scores  $\hat{p}_i$  and  $\hat{m}_i$  compress respectively the information contained in the  $J$  nonresponse models and the  $K$  imputation models.
- $\hat{p}_i$ : consistent estimator of  $p(\mathbf{x}_i; \boldsymbol{\alpha})$  if one of the models in  $\mathcal{C}_1$  is correctly specified.
- $\hat{m}_i$ : consistent estimator of  $m(\mathbf{x}_i; \boldsymbol{\beta})$  if one of the models in  $\mathcal{C}_2$  is correctly specified.
- see Duan and Yin (2017) and Chen and Haziza (2017b).

# Implementation through linear regression

- We construct the imputed values  $y_i^*$  as follows:

$$y_i^* = \mathbf{h}_i^\top \hat{\boldsymbol{\tau}}, \quad i \in S_m,$$

where  $\mathbf{h}_i = (1, \hat{m}_i)^\top$  and

$$\hat{\boldsymbol{\tau}} = \left\{ \sum_{i \in S} w_i r_i (\hat{p}_i^{-1} - 1) \mathbf{h}_i \mathbf{h}_i^\top \right\}^{-1} \sum_{i \in S} w_i r_i (\hat{p}_i^{-1} - 1) \mathbf{h}_i y_i.$$

The resulting imputed estimator is given by

$$\hat{t}_{MR} = \sum_{i \in S} w_i r_i y_i + \sum_{i \in S} w_i (1 - r_i) \mathbf{h}_i^\top \hat{\boldsymbol{\tau}}.$$

# Theoretical properties

## Theorem

*If one of the imputation models is true, then the estimator  $\hat{t}_{MR}$  is consistent in the sense that  $\hat{t}_{MR}/t_y \xrightarrow{P} 1$  as  $n \rightarrow \infty$  and  $N \rightarrow \infty$ .*

## Theorem

*If one of the nonresponse models is true, then the estimator  $\hat{t}_{MR}$  is consistent in the sense that  $\hat{t}_{MR}/t_y \xrightarrow{P} 1$  as  $n \rightarrow \infty$  and  $N \rightarrow \infty$ .*

**Conclusion:**  $\hat{t}_{MR}$  is multiply robust

## Simulation study: the set-up

- We used the simulation setup of Kang and Schafer (2007), that was also used by Chan and Yam (2014) and Han (2014)
- We generated  $B = 1,000$  finite populations of size  $N = 10,000$  as follows:
  - For each unit, a vector  $\mathbf{x} = (x_1, x_2, x_3, x_4)^\top$  was generated from a standard multivariate normal distribution.
  - The survey variable  $y$  was generated according to

$$y = 210 + 27.4x_1 + 13.7(x_2 + x_3 + x_4) + \epsilon,$$

where the errors  $\epsilon$  were generated from a standard normal distribution.

- For each population unit, we generated a size variable  $\psi = 0.5\chi + 1$ , where  $\chi$  was drawn from a chi-square distribution with one degree of freedom.

# Simulation study: the set-up

- **Goal:** estimate the population mean

$$\bar{Y} = N^{-1} \sum_{i \in \mathcal{P}} y_i.$$

- From each finite population, we selected a sample, of size  $n = 800$ , according to **randomized systematic sampling with probability proportional to size**; That is,  $\pi_i = n (\psi_i / \sum_{i \in \mathcal{P}} \psi_i)$ .
- In each sample, the response indicators  $r_i$  were generated independently from a Bernoulli distribution with probability

$$p(\mathbf{x}_i) = \{1 + \exp(x_{1i} - 0.5x_{2i} + 0.25x_{3i} + 0.1x_{4i})\}^{-1}.$$

This led to an overall response rate approximately equal to 50%.



# Simulation study: the set-up

- As in Kang and Schafer (2007), we considered the following transformations of the  $x$ -variables:
  - $z_1 = \exp(x_1/2)$ ;
  - $z_2 = x_2 \{1 + \exp(x_1)\}^{-1} + 10$ ;
  - $z_3 = (x_1 x_3 / 25 + 0.6)^3$ ;
  - $z_4 = (x_2 + x_4 + 20)^2$ .

## 2 imputation models

$$m^1(\mathbf{x}; \beta^1) = \beta_0^1 + \beta_1^1 x_1 + \dots + \beta_5^1 x_4.$$

$$m^2(\mathbf{x}; \beta^1) = \beta_0^1 + \beta_1^1 z_1 + \dots + \beta_5^1 z_4$$

## 2 nonresponse models

$$p^1(\mathbf{x}; \alpha^1) = \left\{ 1 + \exp \left( \alpha_0^1 + \alpha_1^1 x_1 + \dots + \alpha_4^1 x_4 \right) \right\}^{-1}$$

$$p^2(\mathbf{x}; \alpha^1) = \left\{ 1 + \exp \left( \alpha_0^1 + \alpha_1^1 z_1 + \dots + \alpha_4^1 z_4 \right) \right\}^{-1}$$

# Simulation study: the set-up

We computed several estimators of  $\overline{Y}$  of the form:

$$\widehat{\overline{Y}}_I = \frac{1}{N} \left( \sum_{i \in S} w_i r_i y_i + \sum_{i \in S} w_i (1 - r_i) y_i^* \right)$$

1. Two imputed estimators through **nearest-neighbor imputation based on the set of matching variables  $x$  and  $z$** , respectively
2. Four doubly robust estimators (Haziza and Rao, 2006)
3. Nine multiply robust estimators

## Simulation study: results

Estimator	Percent relative bias	Root mean square error
$\hat{Y}_{NN}(00 10)$	-0.8	2.21
$\hat{Y}_{NN}(00 01)$	-3.8	8.21
$\hat{Y}_{DR}(10 10)$	-0.0	1.40
$\hat{Y}_{DR}(10 01)$	0.0	1.97
$\hat{Y}_{DR}(01 10)$	0.0	1.45
$\hat{Y}_{DR}(01 01)$	-2.5	16.88
$\hat{Y}_{MR}(10 10)$	-0.0	1.40
$\hat{Y}_{MR}(10 01)$	0.1	1.63
$\hat{Y}_{MR}(01 10)$	-0.0	1.40
$\hat{Y}_{MR}(01 01)$	-1.2	3.05
$\hat{Y}_{MR}(11 10)$	-0.0	1.40
$\hat{Y}_{MR}(11 01)$	0.0	1.64
$\hat{Y}_{MR}(10 11)$	-0.0	1.40
$\hat{Y}_{MR}(01 11)$	-0.0	1.40
$\hat{Y}_{MR}(11 11)$	-0.0	1.40

## Simulation study: More results

- How do multiply robust procedures perform when all the models are misspecified?

(MR1).  $p^2(\mathbf{z}; \alpha^2)$  and  $m^2(\mathbf{z}; \beta^2)$  :  $\hat{\bar{Y}}_{\text{MR1}}$ .

(MR2).  $p^2(\mathbf{z}; \alpha^2)$  and  $m^2(\mathbf{z}; \beta^2)$  : + one additional outcome regression model by including  $(z_1, z_2, z_3, z_4)$  and all their interactions as covariates:  $\hat{\bar{Y}}_{\text{MR2}}$ .

(MR3). One additional outcome regression model by using  $(\sqrt{z_1}, \sqrt{z_2}, \sqrt{z_3}, \sqrt{z_4})$  and all their interactions :  $\hat{\bar{Y}}_{\text{MR3}}$ .

(MR4). One additional outcome regression model by using  $(\log z_1, \log z_2, \log z_3, \log z_4)$  and all their interactions :  $\hat{\bar{Y}}_{\text{MR4}}$ .

(MR5). Three additional outcome regression models described in (MR2), (MR3) and (MR4):  $\hat{\bar{Y}}_{\text{MR5}}$ .

## Simulation study: more results

**Table:** Percent relative bias (RB), and root mean squared error (RMSE) with a response rate of 50%

Estimator	$\hat{Y}_{MR1}$	$\hat{Y}_{MR2}$	$\hat{Y}_{MR3}$	$\hat{Y}_{MR4}$	$\hat{Y}_{MR5}$
RB	-1.20	-0.49	-0.52	-0.55	-0.42
RMSE	3.05	2.03	1.91	1.90	1.92

**Results suggest:** Despite being inconsistent, multiply robust imputation procedures tend to have good numerical performance even if all the models are misspecified.

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation**
- 8 Final remarks

# Multiple imputation, Rubin (1978, 1987)

- $M \geq 2$  values are imputed for each missing value, resulting in the creation of  $M \geq 2$  imputed data files.
- The question is then: how to combine the  $M$  files to obtain a valid point estimate and a valid variance estimate?
- Although it is a popular method in many fields in statistics, single imputation is the norm in surveys. Why?
  - Survey statisticians and data analysts are used to work with a single imputed file;
  - There are some questions and debates about the validity of multiple imputation in a survey sampling setting.

# Rubin's Rule

- Point estimator of  $\theta_N$ :

$$\hat{\theta}_{I,M} = M^{-1} \sum_{m=1}^M \hat{\theta}_I^{(m)},$$

where  $\hat{\theta}_I^{(m)}$  the imputed estimator corresponding to the  $m$ th imputed data set.

- Variance estimator for  $\hat{\theta}_{I,M}$ :

$$T_M = \overline{W}_M + \left(1 + \frac{1}{M}\right) B_M,$$

where

$$\overline{W}_M = \frac{1}{M} \sum_{m=1}^M W^{(m)}, \quad B_M = \frac{1}{M-1} \sum_{m=1}^M \left( \hat{\theta}_I^{(m)} - \hat{\theta}_{I,M} \right)^2$$

and  $W^{(m)}$  is the variance estimator computed from the  $m$ th imputed data set **treating imputed values as if they were observed values.**



# Generating imputations

- Since multiple imputation was originally proposed under Bayesian considerations, Bayesian imputation procedures provide a natural option.
- **Idea:** generate imputed values from  $P(\mathbf{y}_{mis} \mid \mathbf{y}_{obs})$ , which can be expressed as

$$P(\mathbf{y}_{mis} \mid \mathbf{y}_{obs}) = \int P(\mathbf{y}_{mis} \mid \beta, \mathbf{y}_{obs}) P(\beta \mid \mathbf{y}_{obs}) d\beta,$$

where  $\beta$  is a parameter indexing the imputation model for the y-variable.

- Interpretation: For  $1 \leq m \leq M$ ,
  - First, draw a value  $\beta^{(m)}$  from  $P(\beta \mid \mathbf{y}_{obs})$ .
  - Then, sample  $\mathbf{y}_{mis}^{(m)}$  from  $P(\mathbf{y}_{mis} \mid \mathbf{y}_{obs}, \beta^{(m)})$ .

# Bayesian validity of multiple imputation

- If the imputation model holds, then multiple imputation yields valid inferences provided  $M$  is not too small.
- In practice, the imputer's model is generally different from the analyst's model, a situation known as **uncongeniality**.
- For instance, in business surveys, imputation may be done at the NAICS3 level, whereas the analyst is interested in producing estimates at the NAICS6 level → the analyst's model is more saturated than the imputer's model.
- Uncongeniality is especially problematic in this case; Xie and Meng (2017).
- Under uncongeniality, multiple imputation may lead to invalid inferences.
- What about the properties of multiple imputation from a frequentist perspective?

## Validity of multiple imputation: frequentist *pq*-approach

- In the frequentist *pq* approach, the properties of point and variance estimators are evaluated with respect to the joint distribution induced by the sampling design and the nonresponse mechanism. The population  $y$ -values are treated as fixed.
- Here, the imputation must be proper for multiple imputation to yield valid inferences.
- To be proper, the imputation must meet the following three conditions:

$$\text{C1: } \mathbb{E}_q(\hat{\theta}_{I,\infty}) = \hat{\theta}_F.$$

$$\text{C2: } \mathbb{E}_q(\overline{W}_\infty) = \mathbb{V}_p(\hat{\theta}_F).$$

$$\text{C3: } \mathbb{E}_q(B_\infty) \approx \mathbb{V}_q(\hat{\theta}_{I,\infty}).$$

- **Issue:** proper imputation may be difficult to achieve for complex sampling designs (Binder and Sun, 1996).

## Validity of multiple imputation: frequentist *pq*-approach

- To overcome the difficulty, Bjørnstad (2007) suggested a simple modification to the customary multiple imputation variance estimator:

$$T_M^* = \overline{W}_M + \left( c + \frac{1}{M} \right) B_M,$$

where  $c$  is such that

$$\mathbb{E}_{pqI} (T_M^*) = \mathbb{V}_{pqI} (\hat{\theta}_{I,M}).$$

- For instance, in the context of simple random sampling without replacement and random hot-deck imputation, he found  $c \approx \hat{p}_r^{-1}$ , where  $\hat{p}_r$  is the observed response rate.
- Drawback:** the value of  $c$  depends on the sampling design, the parameter of interest and the imputation procedure.
- Extensions to more complex cases are needed.

# Validity of multiple imputation: frequentist *mpq*-approach

- Recall that

$$V_{\text{tot}} = V_{\text{sam}} + V_{\text{nr}} + \mathbf{V}_{\text{mix}} + V_{\text{imp}}$$

and

$$T_M = \overline{W}_M + \left(1 + \frac{1}{M}\right) B_M.$$

- $\overline{W}_M$  is asymptotically unbiased for  $V_{\text{sam}}$ .
- $(1 + M^{-1})B_M$  is asymptotically unbiased for  $V_{\text{nr}} + V_{\text{imp}}$ .
- Issue: The estimator  $T_M$  does not track the component  $V_{\text{mix}}$ . Therefore,

$$\text{Bias}(T_M) = -V_{\text{mix}};$$

see Kott (1995) and Kim et al. (2006).

## Validity of multiple imputation: frequentist *mpq*-approach

- The term  $\mathbb{V}_{\text{mix}}$  is generally not equal to 0 and, in some cases, its contribution to the total variance may be large  $\rightarrow T_M$  may be considerably biased.
- Often (but not always), the term  $\mathbb{V}_{\text{mix}}$  is negative leading to conservative variance estimator.
- The condition  $\mathbb{V}_{\text{mix}} = 0$ , is called **the self-efficiency condition** (Meng, 1994; Meng and Romero, 2003).
- Whether or not the self-efficiency condition holds depends on the sampling design, on the parameter of interest and on the imputation procedure used to fill in the missing values.

## Validity of multiple imputation: survey features

- In the context of survey data, Reiter et al. (2006) showed that the point and variance estimators are generally biased when complex survey-design features are not accounted for in the imputation models. Survey-design features include unequal weighting, stratification and clustering.
- Reiter et al. (2006) and Schenker et al. (2006) suggested including design information (e.g., cluster membership indicators, stratum indicators and survey weights) as additional covariates in  $\mathbf{x}$ .
- Generally, incorporating design information into the imputation model may be challenging as the weighting process is typically complex.
- Kim and Yang (2017): Rubin's variance estimator may exhibit some bias even after incorporating design information in the imputation model → proposed an approach, whereby Rubin's rule can be safely used.

# Outline

- 1 Introduction
- 2 Estimation in ideal conditions
- 3 Imputation procedures
- 4 Properties of imputed estimators
- 5 Variance estimation
- 6 Multiply robust imputation
- 7 Multiple imputation
- 8 Final remarks**



## Other topics

- **Empirical likelihood confidence intervals under imputation:** Cai, Rao and Malgorzata (2019)
- **Imputation for multivariate parameters:** challenge is to preserve the relationships; e.g., Shao and Wang (2002), Andridge and Little (2010), Chauvet and Haziza (2012), Kim and Fuller (2012), Chaput et al. (2018)
- **NMAR:** challenging problem; e.g., Siddique et al. (2012), van Buuren (2012) and Sullivan and Andridge (2015).

Grazie mille!