

Statistica Sinica Preprint No: SS-2020-0053

| | |
|--|--|
| Title | Estimation for nonignorable missing response or covariate using semi-parametric quantile regression imputation and a parametric response probability model |
| Manuscript ID | SS-2020-0053 |
| URL | http://www.stat.sinica.edu.tw/statistica/ |
| DOI | 10.5705/ss.202020.0053 |
| Complete List of Authors | Emily Berg and Cindy Yu |
| Corresponding Author | Emily Berg |
| E-mail | emilyb@iastate.edu |
| Notice: Accepted version subject to English editing. | |

Estimation for nonignorable missing response or covariate using semi-parametric quantile regression imputation and a parametric response probability model

Emily Berg and Cindy Yu

Department of Statistics Iowa State University

Abstract: We address the problem of imputation when a response or covariate may subject to nonignorable (or, equivalently, missing not at random) nonresponse, meaning the response probability may depend on a variable that is not always observed. We discuss model identification and develop a novel estimator of the parameters of the response probability. We further utilize a propensity score adjustment to incorporate a subset for which both the response and the covariate are missing. We derive an approximation for the large sample variance and assess the finite sample properties of the variance estimator through simulations. The simulations also indicate that quantile regression offers a compromise between fully parametric and non-parametric alternatives. In an application to data from a 2011 survey of pet owners, quantile regression allows us to model complex relations between two types of veterinary expenditures, and we find evidence of nonignorable nonresponse.

Key words and phrases: B-spline, survey sampling, missing not at random

1. Introduction

A widely adopted remedy for missing data is to replace each missing value with one or more imputed values (Kim and Shao, 2014; Rubin, 1987). An imputation model defines (1) a relationship between a response (y) and a covariate (x), and (2) the nature of the dependence between the event of responding and (x, y) . A common simplifying assumption is that the data are missing at random (MAR), meaning that the probability of responding is independent of the missing variable after conditioning on fully observed variables. Under the MAR assumption, Kim (2011) and Wang and Chen (2009), respectively, develop fully parametric and non-parametric imputation procedures. Chen and Yu (2016) and Berg and Yu (2019) construct imputed values under the assumptions of a semiparametric quantile regression model, assuming MAR nonresponse. When the event of responding is not independent of missing values given observed values, the response mechanism is called missing not at random (MNAR) or nonignorable. (Hereafter, we use MNAR and nonignorable interchangeably.) We extend Chen and Yu (2016) to nonignorable nonresponse for a data structure in which neither the response nor the covariate is fully observed.

A condition for model identification in the presence of MNAR nonresponse is the existence of a nonresponse instrument, a variable that is

correlated with the response y but conditionally independent of the event of responding given y (Wang, Shao, and Kim, 2014). Tang, Little, and Rubin (2003) estimate a fully parametric model for y given x , using x as an instrument, without requiring a specific form for the response probability. Zhao and Shao (2015) extend the framework of Tang, Little, and Rubin (2003) to include an additional instrument. Other approaches, such as Wang, Shao, and Kim (2014) and Chang and Kott (2008), use an instrumental variable to estimate a parametrized propensity score model that depends on y but not on the instrument. Shao and Wang (2016) generalize the propensity score model of Wang, Shao, and Kim (2014) to include a non-parametric component. Riddles, Kim, and Im (2016) use likelihood-based methods to improve upon the efficiency of calibration estimation. Zhao and Ma (2019) use an instrumental variable but avoid estimating the response probability directly. Miao and Tchetgen Tchetgen (2016) develop a doubly robust estimator under the assumption that an instrumental variable (called a “shadow variable” in their work) exists. Fang, Zhao, and Shao (2018) use an instrumental variable assumption to estimate the coefficient associated with a missing covariate when the response probability depends on the covariate. However, it is well-known that identifying an instrumental variable in a given data set is nontrivial. Morikawa and Kim (2017) gener-

alize the instrumental variable condition of Wang, Shao, and Kim (2014) by deriving a necessary and sufficient condition for model identification under MNAR nonresponse. They develop an efficient propensity score estimator, assuming a univariate response variable is missing and a univariate covariate is fully observed. We extend the identification condition of Morikawa and Kim (2017) to accommodate missing covariates and construct a completed data set through imputation.

We propose to generate imputed values from a semi-parametric quantile regression model and then use estimates of the response probabilities to approximate required expectations for non-respondents. We augment the imputation procedure with a propensity score adjustment to incorporate a subset for which both the response and the covariate are missing. In our application, x and y represent two types of veterinary expenditures, neither of which is fully observed and either of which may influence the probability of responding. Semi-parametric quantile regression provides the needed flexibility to model nonlinear associations between the two types of veterinary expenditures. We define parametric and non-parametric alternatives for the purpose of comparison in the simulation study. As our data set has a univariate covariate, we focus on that case and briefly discuss an extension to multivariate covariates in Section 6.

We validate our proposed procedure through theory and simulation, and then apply the method to data from a survey of pet owners. In Section 2, we define the model assumptions, imputation, and estimation procedures. In Section 3, we define a variance estimator based on a linear approximation. In Section 4, we conduct simulation studies to compare alternative imputation models and assess the finite sample properties of the variance estimator. We apply the method to impute veterinary expenditures in Section 5. We summarize and discuss future work in Section 6.

2. Model Assumptions, Imputation and Estimation Procedures

Let x_i and y_i denote a continuous covariate and a continuous response variable, respectively, with a compact support on the box $[M_{1x}, M_{2x}] \times [M_{1y}, M_{2y}]$, where $i = 1, \dots, n$. Let δ_i denote a response indicator variable such that $\delta_i = 1$ if both x_i and y_i are observed, $\delta_i = 2$ if x_i is observed and y_i is missing, and $\delta_i = 3$ if y_i is observed and x_i is missing. We also use $\delta_{ki} = I[\delta_i = k]$ for $k = 1, 2, 3$. Table 1 shows the data structure.

Table 1: Structure of Missing Data

| Covariate x | Response y | Response Indicator δ |
|---------------|--------------|-----------------------------|
| ✓ | ✓ | 1 |
| ✓ | ? | 2 |
| ? | ✓ | 3 |

Assume that (x_i, y_i, δ_i) for $i = 1, \dots, n$ are *iid* realizations of the random variable (X, Y, Δ) with joint CDF $F(x, y, \delta)$. Further, assume X and Y are absolutely continuous and denote their corresponding conditional pdf's by $f(y|x, \delta)$ and $f(x|y, \delta)$, respectively. Assume Δ has parametric conditional pmf given by

$$P(\Delta = k | X = x, Y = y) = \frac{\exp(\phi_{k0} + \phi_{k1}x + \phi_{k2}y)}{\sum_{k=1}^3 \exp(\phi_{k0} + \phi_{k1}x + \phi_{k2}y)}, \quad (2.1)$$

for $k = 1, 2, 3$, where $(\phi_{10}, \phi_{11}, \phi_{12}) = (0, 0, 0)$.

To identify the parameters of (2.1), we require an additional assumption. By a direct extension of Theorem 3.1 of Morikawa and Kim (2017) to missing covariates, the additional assumption is that $F(x, y, \delta)$ is a joint CDF such that the condition

$$E[\exp(-\phi_{20} - \phi_{21}x_i - \phi_{22}Y) | x_i, \delta_i = 1] = E[\exp(-\phi'_{20} - \phi'_{21}x_i - \phi'_{22}Y) | x_i, \delta_i = 1] \quad (2.2)$$

almost everywhere implies $(\phi_{20}, \phi_{21}, \phi_{22}) = (\phi'_{20}, \phi'_{21}, \phi'_{22})$, and the condition

$$E[\exp(-\phi_{30} - \phi_{31}X - \phi_{32}y_i) \mid y_i, \delta_i = 1] = E[\exp(-\phi'_{30} - \phi'_{31}X - \phi'_{32}y_i) \mid y_i, \delta_i = 1] \quad (2.3)$$

almost everywhere implies $(\phi_{30}, \phi_{31}, \phi_{32}) = (\phi'_{30}, \phi'_{31}, \phi'_{32})$. If $\phi_{31} = \phi_{22} = 0$, then MAR holds and the model is automatically identified.

Sufficient conditions for (2.2) and (2.3) are that

$$h_y(\phi_{22}, x) = -\log(E[\exp\{-\phi_{22}Y\} \mid x, \delta = 1]) \quad (2.4)$$

is not in the column space of x , and

$$h_x(\phi_{31}, y) = -\log(E[\exp\{-\phi_{31}X\} \mid y, \delta = 1]) \quad (2.5)$$

is not in the column space of y . If $h_y(\phi_{22}, x)$ is in the column space of x , then ϕ_{21} is confounded with ϕ_{22} . Similarly, we require $h_x(\phi_{31}, y)$ to be not in the column space of y to prevent ϕ_{32} from being confounded with ϕ_{31} . Note that $-h_y(\phi_{22}, x)$ is the cumulant generating function of $f(y \mid x, \delta = 1)$, and likewise for $-h_x(\phi_{31}, y)$. An aspect of (2.4) and (2.5) that is of practical importance is that one can check these conditions using $\{(x_i, y_i) : \delta_i = 1\}$, as we illustrate in the data analysis of Section 5.

Let the parameter of interest be $\theta_0 = Eg(X, Y) = \sum_{\delta=1}^3 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) dF(x, y, \delta)$.

In the absence of any nonresponse, an estimator of $Eg(X, Y)$ is $\hat{\theta}_{full} =$

$n^{-1} \sum_{i=1}^n g(x_i, y_i)$. The estimator $\hat{\theta}_{full}$ is not directly applicable because of nonresponse. By Cheng (1994), a consistent estimator of θ_0 is

$$\tilde{\theta} = \frac{1}{n} \sum_{i=1}^n \{\delta_{1i}g(x_i, y_i) + \delta_{2i}E[g(x_i, Y)|x_i, \delta_i = 2] + \delta_{3i}E[g(X, y_i)|y_i, \delta_i = 3]\}. \quad (2.6)$$

We convert expectations given $\delta = 2$ or $\delta = 3$ in (2.6) to expectations given $\delta = 1$ using an “exponential tilting” relationship (Kim and Yu, 2011).

Under (2.1), it is straightforward to show that

$$f(y|x, \delta = 2) = \frac{f(y|x, \delta = 1)\exp(\phi_{22}y)}{E[\exp(\phi_{22}Y) | x, \delta = 1]}, \quad (2.7)$$

and

$$f(x|y, \delta = 3) = \frac{f(x|y, \delta = 1)\exp(\phi_{31}x)}{E[\exp(\phi_{31}X) | y, \delta = 1]}, \quad (2.8)$$

where ϕ_{22} and ϕ_{31} are the tilting parameters. The equality in (2.7) allows us to express the conditional expectation for the group with $\delta = 2$ in (2.6) as a function of different expectations given $\delta = 1$ by

$$E[g(x, Y)|x, \delta = 2] = \frac{E[g(x, Y)\exp(\phi_{22}Y)|x, \delta = 1]}{E[\exp(\phi_{22}Y)|x, \delta = 1]}. \quad (2.9)$$

Similarly, the third conditional expectation for the group with $\delta = 3$ in (2.6) converts to a ratio of two expectations given $\delta = 1$ as

$$E[g(X, y)|y, \delta = 3] = \frac{E[g(X, y)\exp(\phi_{31}X)|y, \delta = 1]}{E[\exp(\phi_{31}X)|y, \delta = 1]}, \quad (2.10)$$

2.1 Approximating Expectations with Estimated Quantiles 9

The expressions (2.9) and (2.10) show that we can estimate θ with (1) estimates of $f(y | x, \delta = 1)$ and $f(x | y, \delta = 1)$ and (2) estimates of ϕ_{22} and ϕ_{31} . In this paper, we focus on the use of semi-parametric quantile regression to estimate $f(y | x, \delta = 1)$ and $f(x | y, \delta = 1)$. We compare to non-parametric and fully parametric alternatives in the simulations. We first define our estimation method for known (ϕ_{22}, ϕ_{31}) in Section 2.1 and explain how to estimate unknown (ϕ_{22}, ϕ_{31}) in Section 2.2.

2.1 Approximating Expectations with Estimated Quantiles

We approximate $f(y|x, \delta = 1)$ and $f(x|y, \delta = 1)$ through their conditional quantile regression functions, denoted $q_\tau(x)$ and $q_\tau(y)$, respectively, for $\tau \in (0, 1)$. By definition, the quantile regression functions satisfy $\tau = P(Y \leq q_\tau(x)|x, \delta = 1)$ and $\tau = P(X \leq q_\tau(y)|y, \delta = 1)$. Assume $q_\tau(x)$ and $q_\tau(y)$ are one-to-one functions of x and y , respectively, for every τ . A well-known fact is that $q_\tau(x)$ and $q_\tau(y)$ satisfy $q_\tau(x) = \operatorname{argmin}_a \int \rho_\tau(y - a)f(y|x, \delta = 1)dy$ and $q_\tau(y) = \operatorname{argmin}_a \int \rho_\tau(x - a)f(x|y, \delta = 1)dx$, where $\rho_\tau(u)$ is the “check function” defined by, $\rho_\tau(u) = u(\tau - I[u < 0])$ (Koenker, 2005). We approximate $q_\tau(x)$ and $q_\tau(y)$ with a B-spline, allowing flexibility and computational efficiency. Let $\mathbf{B}(x)$ be a B-spline of degree $p_{y|x}$ and with $K_{n_1, y}$ interior knots, where n_1 is the sample size for $\delta = 1$. For any $\tau \in (0, 1)$,

2.1 Approximating Expectations with Estimated Quantiles 10

we estimate $q_\tau(x)$ by $\hat{q}_\tau(x) = \mathbf{B}(x)' \hat{\boldsymbol{\beta}}_{y|x}(\tau)$, where

$$\hat{\boldsymbol{\beta}}_{y|x}(\tau) = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \delta_{1i} \rho_\tau(y_i - \mathbf{B}(x_i)' \boldsymbol{\beta}) + \frac{\lambda_{n_1, y}}{2} \boldsymbol{\beta}' \mathbf{D}'_m \mathbf{D}_m \boldsymbol{\beta} \right\}, \quad (2.11)$$

\mathbf{D}_m is a difference matrix of order m , and $\lambda_{n_1, y} > 0$ is the smoothing parameter. See Chen and Yu (2016) and Berg and Yu (2019) for a precise definition of the B-spline and the difference matrix \mathbf{D}_m . In an analogous fashion, define the estimate of $q_\tau(y)$ by $\hat{q}_\tau(y) = \mathbf{B}(y)' \hat{\boldsymbol{\beta}}_{x|y}(\tau)$, where

$$\hat{\boldsymbol{\beta}}_{x|y}(\tau) = \underset{\boldsymbol{\beta}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \delta_{1i} \rho_\tau(x_i - \mathbf{B}(y_i)' \boldsymbol{\beta}) + \frac{\lambda_{n_1, x}}{2} \boldsymbol{\beta}' \mathbf{D}'_m \mathbf{D}_m \boldsymbol{\beta} \right\} \text{ for a given } \tau.$$

To approximate the full distributions of $f(y_i | x_i, \delta_i = 1)$ and $f(x_i | y_i, \delta_i = 1)$, we obtain estimates $\hat{\boldsymbol{\beta}}_{x|y}(\tau)$ and $\hat{\boldsymbol{\beta}}_{y|x}(\tau)$ for a grid of τ_j defined by $\tau_j = \tau_1 + (j - 1)/J$ for $j = 2, \dots, J$, where $\tau_1 \sim \operatorname{Unif}(0, 1/J)$. The resulting estimated quantiles, defined as $\mathbf{y}_i^* = \{y_{ij}^* = \hat{q}_{\tau_j}(x_i) : j = 1, \dots, J\}$, serve as imputed values for element i with $\delta_i = 2$. Likewise, $\mathbf{x}_i^* = \{x_{ij}^* = \hat{q}_{\tau_j}(y_i) : j = 1, \dots, J\}$ serve as imputed values for element i with $\delta_i = 3$.

The sequence of estimated quantiles permits us to approximate the expectations defining $\tilde{\theta}$. For any arbitrary function $m(x, y)$, a variable transformation implies

$$E[m(x, Y)|x, \delta = 1] = \int_0^1 m(x, F_{y|x, \delta=1}^{-1}(\tau)) \frac{f_{y|x, \delta=1}(F_{y|x}^{-1}(\tau) | x)}{f_{y|x, \delta=1}(F_{y|x}^{-1}(\tau) | x)} d\tau = \int_0^1 m(x, q_\tau(x)) d\tau.$$

We approximate $E[m(x, Y)|x, \delta = 1]$ by $\hat{E}[m(x, Y)|x, \delta = 1] = J^{-1} \sum_{j=1}^J m(x, \hat{q}_{\tau_j}(x))$.

2.2 Estimation of Response Probability 11

We approximate the numerator and denominator of (2.9) by replacing $m(x, Y)$ with $g(x, Y)\exp(\phi_{22}Y)$ and $\exp(\phi_{22}Y)$, respectively. Specifically, $\hat{E}[g(x, Y)\exp(\phi_{22}Y)|x, \delta = 1] = J^{-1} \sum_{j=1}^J g(x, \hat{q}_{\tau_j}(x))\exp(\phi_{22}\hat{q}_{\tau_j}(x))$, and $\hat{E}[\exp(\phi_{22}Y)|x, \delta = 1] = J^{-1} \sum_{j=1}^J \exp(\phi_{22}\hat{q}_{\tau_j}(x))$. Then, an approximation for (2.9) is

$$\hat{E}[g(x_i, Y)|x_i, \delta_i = 2] = \sum_{j=1}^J w_{2ij}(\boldsymbol{\phi}_2, \mathbf{y}_i^*)g(x_i, y_{ij}^*), \quad (2.12)$$

where $\boldsymbol{\phi}_2 = (\phi_{20}, \phi_{21}, \phi_{22})'$, and

$$w_{2ij}(\boldsymbol{\phi}_2, \mathbf{y}_i^*) = \frac{\exp(\phi_{22}y_{ij}^*)}{\sum_{j=1}^J \exp(\phi_{22}y_{ij}^*)}. \quad (2.13)$$

Analogously, we estimate the expectation in (2.10) as

$$\hat{E}[g(X, y_i)|y_i, \delta_i = 3] = \sum_{j=1}^J w_{3ij}(\boldsymbol{\phi}_3, \mathbf{x}_i^*)g(x_{ij}^*, y_i), \quad (2.14)$$

where $\boldsymbol{\phi}_3 = (\phi_{30}, \phi_{31}, \phi_{32})'$ and

$$w_{3ij}(\boldsymbol{\phi}_3, \mathbf{x}_i^*) = \frac{\exp(\phi_{31}x_{ij}^*)}{\sum_{j=1}^J \exp(\phi_{31}x_{ij}^*)}. \quad (2.15)$$

2.2 Estimation of Response Probability

The estimated expectations in (2.12) and (2.14) require estimators of ϕ_{22} and ϕ_{31} , the two tilting parameters. We estimate $\boldsymbol{\phi} = (\boldsymbol{\phi}'_2, \boldsymbol{\phi}'_3)'$ using conditional probabilities. Define for $k = 2, 3$,

$$\pi_{k1}(x_i, y_i, \boldsymbol{\phi}_k) := P(\delta_i = k | x_i, y_i, \boldsymbol{\phi}_k, \delta_i \in \{1, k\}) = \frac{\exp(\phi_{k0} + \phi_{k1}x_i + \phi_{k2}y_i)}{1 + \exp(\phi_{k0} + \phi_{k1}x_i + \phi_{k2}y_i)},$$

2.2 Estimation of Response Probability 12

and let $\pi_{1k\infty}(v) := P(\delta = 1|v, \delta \in \{1, k\})$ for $v = x$ if $k = 2$, and $v = y$ if $k = 3$. Based on a result of Morikawa and Kim (2017), we can show that

$$\pi_{12\infty}(x) = E[1 - \pi_{21}(x, Y, \phi_2) | x, \delta \in \{1, 2\}] = \frac{\exp(-\phi_{20} - \phi_{21}x + h_y(-\phi_{22}, x))}{1 + \exp(-\phi_{20} - \phi_{21}x + h_y(-\phi_{22}, x))}, \quad (2.16)$$

and

$$\pi_{13\infty}(y) = E[1 - \pi_{31}(X, y, \phi_3) | y, \delta \in \{1, 3\}] = \frac{\exp(-\phi_{30} + h_x(-\phi_{31}, y) - \phi_{32}y)}{1 + \exp(-\phi_{30} + h_x(-\phi_{31}, y) - \phi_{32}y)}, \quad (2.17)$$

where $h_y(\phi_{22}, x_i)$ and $h_x(\phi_{31}, y_i)$ are defined in (2.4) and (2.5), respectively.

Note that $\pi_{12\infty}(x)$ depends only on x and $\pi_{13\infty}(y)$ depends only on y . Thus equation (2.16) suggests an estimator of ϕ_2 defined as

$$\hat{\phi}_2 = \underset{\phi_2}{\operatorname{argmax}} \sum_{i:\delta_i=1} \log \left[\frac{\exp(-\phi_{20} - \phi_{21}x_i + \hat{h}_y(-\phi_{22}, \hat{\mathbf{q}}_{yi}))}{1 + \exp(-\phi_{20} - \phi_{21}x_i + \hat{h}_y(-\phi_{22}, \hat{\mathbf{q}}_{yi}))} \right] + \sum_{i:\delta_i=2} \log \left[1 - \frac{\exp(-\phi_{20} - \phi_{21}x_i + \hat{h}_y(-\phi_{22}, \hat{\mathbf{q}}_{yi}))}{1 + \exp(-\phi_{20} - \phi_{21}x_i + \hat{h}_y(-\phi_{22}, \hat{\mathbf{q}}_{yi}))} \right], \quad (2.18)$$

where $\hat{h}_y(\phi_{22}, \hat{\mathbf{q}}_{yi}) = -\log \left(J^{-1} \sum_{j=1}^J \exp\{-\phi_{22}y_{ij}^*\} \right)$. Likewise, we estimate

ϕ_3 as

$$\hat{\phi}_3 = \underset{\phi_3}{\operatorname{argmax}} \sum_{i:\delta_i=1} \log \left[\frac{\exp(-\phi_{30} + \hat{h}_x(-\phi_{31}, \hat{\mathbf{q}}_{xi}) - \phi_{32}y_i)}{1 + \exp(-\phi_{30} + \hat{h}_x(-\phi_{31}, \hat{\mathbf{q}}_{xi}) - \phi_{32}y_i)} \right] + \sum_{i:\delta_i=3} \log \left[1 - \frac{\exp(-\phi_{30} + \hat{h}_x(-\phi_{31}, \hat{\mathbf{q}}_{xi}) - \phi_{32}y_i)}{1 + \exp(-\phi_{30} + \hat{h}_x(-\phi_{31}, \hat{\mathbf{q}}_{xi}) - \phi_{32}y_i)} \right], \quad (2.19)$$

where $\hat{h}_x(\phi_{31}, \hat{\mathbf{q}}_{xi}) = -\log \left(J^{-1} \sum_{j=1}^J \exp\{-\phi_{31}x_{ij}^*\} \right)$. Note that \hat{h}_x and \hat{h}_y

are estimates of h_x and h_y using the imputed values y_{ij}^* and x_{ij}^* . In opera-

tion, we use the R function `optim` to find the maximum, where the initial value for ϕ_2 is from the logistic regression of $1 - \delta_{2i}$ on $(1, x_i, \mathbf{B}(x_i)' J^{-1} \sum_{j=1}^J \hat{\beta}_{y|x}(\tau_j))'$ for the set with $\delta_{3i} = 0$. We define the initial value for ϕ_3 from the logistic regression of $1 - \delta_{3i}$ on $(1, \mathbf{B}(y_i)' J^{-1} \sum_{j=1}^J \hat{\beta}_{x|y}(\tau_j), y_i)'$ for the set with $\delta_{2i} = 0$.

In summary, we define the basic steps of the estimation procedure:

1. Use $\{(x_i, y_i) : \delta_i = 1\}$ to estimate the quantile regression model and define imputed values y_{ij}^* and x_{ij}^* , as defined in Section 2.1.
2. Estimate ϕ_2 and ϕ_3 as defined in Section 2.2.
3. Define the imputed estimator $\hat{\theta}$ by

$$\hat{\theta} = n^{-1} \sum_{i=1}^n \left\{ \delta_{1i} g(x_i, y_i) + \delta_{2i} \sum_{j=1}^J w_{2ij} (\hat{\phi}_2, \mathbf{y}_i^*) g(x_i, y_{ij}^*) + \delta_{3i} \sum_{j=1}^J w_{3ij} (\hat{\phi}_3, \mathbf{x}_i^*) g(x_{ij}^*, y_i) \right\}. \quad (2.20)$$

This completes the description of our imputation and estimation procedures.

3. Large Sample Theories and Variance Estimation

As a pre-cursor to the statement of the large sample distributions of $\hat{\phi}_2$ and $\hat{\phi}_3$, we give the large sample distributions of the the estimates of the

quantile regression coefficients as Lemma 1. We state Lemma 1 without proof because lemma 1 is essentially an application of Yoshida (2013) to the set with $\delta_i = 1$. We use the linear approximation in lemma 1 in the subsequent derivation of the asymptotic properties of $(\hat{\phi}'_2, \hat{\phi}'_3)'$ and $\hat{\theta}$.

Lemma 1 uses the following property of Barrow and Smith (1978). The result is that the best L_∞ approximation to $q_\tau(x)$ (as a function of x), denoted $\mathbf{B}(x)' \boldsymbol{\beta}_{y|x}^*(\tau)$, satisfies $\sup_{x \in [M_{1x}, M_{2x}]} |q_\tau(x) + b_\tau^a(x) - \mathbf{B}(x)' \boldsymbol{\beta}_{y|x}^*(\tau)| = o(K_{n_1, y}^{-(p_{y|x}+1)})$, where $b_\tau^a(x)$ is the bias due to using a B-spline to approximate the true function $q_\tau(x)$, and is defined as in Yoshida (2013).

Lemma 1. Assume $q_\tau^{(p_{y|x}+1)}(x)$ is continuous, where $q_\tau^{(p_{y|x}+1)}(x)$ denotes the $p+1$ derivative of $q_\tau(x)$ with respect to x , $K_{n_1, y} = O(n_1^{1/(2p_{y|x}+3)})$, and $\lambda_{n_1, y} = O(n_1^{\nu_y})$ for $\nu_y < (p_{y|x} + m + 1)/(2p_{y|x} + 3)$. Then,

$$\sqrt{\frac{n_1}{K_{n_1, y}}} \left(\mathbf{B}(x)' \hat{\boldsymbol{\beta}}_{y|x}(\tau) - q_\tau(x) - b_\tau^a(x) - b_\tau^\lambda(x) \right) = W_{n_1} + o_p(1),$$

where

$$W_{n_1} = \sqrt{\frac{n_1}{K_{n_1, y}}} \mathbf{B}(x)' \mathbf{H}_{n_1, y|x}^{-1}(\tau) \frac{1}{n_1} \sum_{i: \delta_i=1} \mathbf{B}(x_i) \psi_\tau(e_{y|x, i}(\tau)),$$

$$\psi_\tau(u) = \tau - I[u < 0], \quad e_{y|x, i}(\tau) = y_i - q_\tau(x_i),$$

$$\mathbf{H}_{n_1, y|x}(\tau) = \boldsymbol{\Phi}_{y|x}(\tau) + n_1^{-1} \lambda_{n_1, y} \mathbf{D}'_m \mathbf{D}_m,$$

$$b_\tau^\lambda(x) = -\frac{\lambda_{n_1, y}}{n_1} \mathbf{B}(x)' \left(\boldsymbol{\Phi}_{y|x}(\tau) + \frac{\lambda_{n_1, y}}{n_1} \mathbf{D}'_m \mathbf{D}_m \right)^{-1} \mathbf{D}'_m \mathbf{D}_m \boldsymbol{\beta}_{y|x}^*(\tau), \quad (3.1)$$

3.1 Asymptotic Variance of $\hat{\phi}$ and $\hat{\theta}$ 15

and $\Phi_{y|x}(\tau) = \lim_{n_1 \rightarrow \infty} n_1^{-1} \sum_{i:\delta_i=1} f_{y|(x,\delta=1)}(x_i, q_\tau(x_i)) \mathbf{B}(x_i) \mathbf{B}(x_i)'$.

Lemma 1 holds for a given τ , but the order of approximation does not depend on τ . A result analogous to lemma 1 holds for $\hat{\beta}_{x|y}(\tau)$. We assume that the degree of $\mathbf{B}(y)$, denoted $p_{x|y}$, is such that $p_{x|y} \cdot p_{y|x}^{-1} = O(1)$. We also assume that the number of interior knots used to define $\mathbf{B}(y)$, denoted $K_{n_1,x}$, satisfies $K_{n_1,y} \cdot K_{n_1,x}^{-1} = O(1)$.

3.1 Asymptotic Variance of $\hat{\phi}$ and $\hat{\theta}$

We state the large sample distribution of $\hat{\phi}_2$ and $\hat{\theta}$ as Theorems 1 and 2, respectively. Section S1 of the supplement contains a result for $\hat{\phi}_3$ analogous to Theorem 1 as well as proofs.

Theorem 1. *In addition to the assumptions of Lemma 1, assume $\hat{\phi}_2 - \phi_2 = o_p(1)$, $J = O(n^{0.5+\delta})$ for some $\delta > 0$, and the conditions in the supplement hold. Then, $\hat{\phi}_2 - \phi_2 = O_p(n^{-0.5})$, $\hat{\phi}_2 - \phi_2 = n^{-1} \sum_{i=1}^n \mathbf{I}_{\phi_2}^{-1} \mathbf{U}_{\phi_2 i} + o_p(n^{-0.5})$, and $\sqrt{n} \mathbf{V}_{\phi_2}^{-1/2} (\hat{\phi}_2 - \phi_2) \xrightarrow{d} N(0, \mathbf{I}_3)$, where*

$$\mathbf{V}_{\phi_2} = \lim_{n \rightarrow \infty} n^{-1} \mathbf{I}_{\phi_2}^{-1} \left(\sum_{i=1}^n \mathbf{U}_{\phi_2 i} \mathbf{U}_{\phi_2 i}' \right) \mathbf{I}_{\phi_2}^{-1}, \quad (3.2)$$

$$\begin{aligned} \mathbf{I}_{\phi_2} &= \lim_{n \rightarrow \infty} \mathbf{I}_{n,\phi_2}(\mathbf{q}_y), \mathbf{I}_{n,\phi_2}(\mathbf{q}_y) = n^{-1} \sum_{i \in A_{12}} \pi_{12i}(\phi_2, \mathbf{q}_{yi}) (1 - \pi_{12i}(\phi_2, \mathbf{q}_{yi})) \mathbf{z}_{2i}(\phi_2, \mathbf{q}_{yi}) \mathbf{z}_{2i}(\phi_2, \mathbf{q}_{yi})', \\ \mathbf{U}_{\phi_2 i} &= (\delta_{1i} + \delta_{2i}) \mathbf{S}_{i\infty}(\phi_2) + \phi_{22} \delta_{1i} \int_{M_{1x}}^{M_{2x}} p_1^{-1} \pi_{12\infty}(x) (1 - \pi_{12\infty}(x)) \mathbf{z}_{2\infty}(x) \mathbf{B}(x)' \frac{\int_0^1 \exp(\phi_{22} q_\tau(x)) \ell_i(\tau) d\tau}{\int_0^1 \exp(\phi_{22} q_\tau(x)) d\tau} dF(x | \\ \delta_1 + \delta_2 &= 1), p_1 = \lim_{n \rightarrow \infty} n^{-1} n_{11}, \mathbf{S}_{i\infty}(\phi_2) = (\delta_{1i} - \pi_{12\infty}(x_i)) \mathbf{z}_{2i\infty}, \ell_i(\tau) = \end{aligned}$$

3.1 Asymptotic Variance of $\hat{\phi}$ and $\hat{\theta}$ 16

$$\begin{aligned} & \mathbf{H}_{n_1, y|x}^{-1}(\tau) \mathbf{B}(x_i) \psi_\tau(e_{y|x, i}(\tau)), \mathbf{z}_{2\infty}(x) = (-1, -x, -E_2(Y | x))', \mathbf{z}_{2i\infty} = \mathbf{z}_{2\infty}(x_i), \mathbf{z}_{2i}(\phi_2, \mathbf{q}_{yi}) = \\ & (-1, -x_i, -E_{2,J}(Y | x_i; \phi_2, \mathbf{q}_{yi})), \mathbf{q}_{yi} = \{q_{\tau_j}(x_i) : j = 1, \dots, J\}, E_{2,J}(Y | \\ & x_i; \phi_2, \mathbf{q}_{yi}) = \sum_{j=1}^J w_{2ij}(\phi_2, \mathbf{q}_{yi}) q_{\tau_j}(x_i), \\ & \pi_{12i}(\phi_2, \mathbf{q}_{yi}) = \left\{ 1 + \exp \left[\phi_{20} + \phi_{21} x_i + \log \left(J^{-1} \sum_{j=1}^J \exp\{\phi_{22} q_{\tau_j}(x_i)\} \right) \right] \right\}^{-1}, \\ & A_{12} = \{i : \delta_{1i} + \delta_{2i} = 1\}, \mathbf{q}_y = \{\mathbf{q}_{yi} : \delta_{1i} + \delta_{2i} = 1\}, \text{ and } E_2[Y | x] = E[Y | \\ & x, \delta = 2]. \end{aligned}$$

An estimator of the variance of $\hat{\phi}_2$ is

$$\hat{V}\{\hat{\phi}_2\} = n^{-2} \hat{\mathbf{I}}_{n, \phi_2}^{-1} \left(\sum_{i=1}^n \hat{\mathbf{U}}_{\phi_2 i} \hat{\mathbf{U}}'_{\phi_2 i} \right) \hat{\mathbf{I}}_{n, \phi_2}^{-1}, \quad (3.3)$$

where we substitute unknown parameters with their corresponding estimators to define $\hat{\mathbf{I}}_{n, \phi_2}$ and $\hat{\mathbf{U}}_{\phi_2 i}$, as defined explicitly in Section S2 of the supplement.

Theorem 2. *Continue to assume the conditions of Theorem 1. Also, assume $g(X, Y)$ has bounded $2 + c$ moments for $c > 0$ and has bounded second derivatives with respect to both x and y . Let $K_{n_1} = \max\{K_{n_1, y}, K_{n_1, x}\}$.*

Then, $\sqrt{n} V_g^{-0.5}(\hat{\theta} - E[g(X, Y)]) \xrightarrow{d} N(0, 1)$, where $V_g = \lim_{n \rightarrow \infty} (n-1)^{-1} \sum_{i=1}^n (r_i -$

3.1 Asymptotic Variance of $\hat{\phi}$ and $\hat{\theta}$ 17

$$\bar{r})^2, \bar{r} = n^{-1} \sum_{i=1}^n r_i,$$

$$r_i = g(x_i, y_i) - Eg(X, Y) + \delta_{2i}(E_2[g(x_i, Y) | x_i] - g(x_i, y_i)) + \delta_{3i}(E_3[g(X, y_i) | y_i] - g(x_i, y_i)) \quad (3.4)$$

$$+ (\delta_{1i} + \delta_{2i}) \{ \bar{C}_{2\infty} \} \mathbf{e}'_2 \mathbf{I}_{\phi_2}^{-1} \mathbf{U}_{\phi_2, i} + (\delta_{1i} + \delta_{3i}) \{ \bar{C}_{3\infty} \} \mathbf{e}'_3 \mathbf{I}_{\phi_3}^{-1} \mathbf{U}_{\phi_3, i} \\ + \delta_{1i} \left(\int_0^1 \int_{M_{1x}}^{M_{2x}} \mathbf{C}_y(x, \tau)' \boldsymbol{\ell}_i(\tau) dF(x | \delta = 2) d\tau + \int_0^1 \int_{M_{1y}}^{M_{2y}} \mathbf{C}_x(y, \tau)' \mathbf{m}_i(\tau) dF(y | \delta = 3) d\tau \right),$$

$$\bar{C}_{2\infty} = \lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n \delta_{2k} Cov_2(g(x_k, Y), Y | x_k), \bar{C}_{3\infty} = \lim_{n \rightarrow \infty} n^{-1} \sum_{k=1}^n \delta_{3k} Cov_3(g(X, y_k), X | y_k), \\ Cov_2(g(x, Y), Y | x) = Cov(g(x, Y), Y | X = x, \delta = 2), Cov_3(g(X, y), X | y) = Cov(g(X, y), X | Y = y, \delta = 3), \\ E_2[g(x, Y) | x] = E[g(x, Y) | \delta = 2, X = x], E_3[g(X, y) | y] = E[g(X, y) | \delta = 3, Y = y], \mathbf{e}_3 = (0, 0, 1)', \\ \mathbf{C}_y(x, \tau) = \tilde{c}_y(x, \tau) \mathbf{B}(x), \mathbf{C}_x(y, \tau) = \tilde{c}_x(y, \tau) \mathbf{B}(y),$$

$$\tilde{c}_y(x, \tau) = \frac{c_y(x, \tau)}{\int_0^1 \exp(\phi_{22} q_\tau(x))} - E_2[g(x, Y) | x] \frac{\phi_{22} \exp(\phi_{22} q_\tau(x))}{\int_0^1 \exp(\phi_{22} q_\tau(x))} \\ \tilde{c}_x(y, \tau) = \frac{c_x(y, \tau)}{\int_0^1 \exp(\phi_{31} q_\tau(y)) d\tau} - E_3[g(X, y) | y] \frac{\phi_{31} \exp(\phi_{31} q_\tau(y)) \phi_{31}}{\int_0^1 \exp(\phi_{31} q_\tau(y)) d\tau},$$

$c_y(x, \tau) = \exp(\phi_{22} q_\tau(x)) g'_y(x, q_\tau(x)) + g(x, q_\tau(x)) \exp(\phi_{22} q_\tau(x)) \phi_{22}$, $c_x(y, \tau) = \exp(\phi_{31} q_\tau(y)) g'_x(q_\tau(y), y) + g(q_\tau(y), y) \exp(\phi_{31} q_\tau(y)) \phi_{31}$, and $\mathbf{m}_i(\tau)$, \mathbf{I}_{ϕ_3} , and \mathbf{U}_{ϕ_3} are defined in the supplement for the linear approximation for $\hat{\phi}_3$.

A proof of Theorem 2 is presented in Section S1 of the supplement. An estimator of the variance of the imputed estimator is

$$\hat{V}\{\hat{\theta}\} = (n(n-1))^{-1} \sum_{i=1}^n (\hat{r}_i - \bar{r}_i)^2, \quad (3.5)$$

3.2 Propensity Score Adjusted Imputed Estimator 18

where \hat{r}_i is a plug-in estimator of r_i defined in Section S2.3 of the supplement, and $\bar{\hat{r}} = n^{-1} \sum_{i=1}^n \hat{r}_i$. In supplement Section S2.4, we define how to use a further Taylor linearization to estimate the variance of “composite” estimators of the form $\hat{\theta} = h(\hat{\theta}_1, \dots, \hat{\theta}_K)$ of a parameter $\theta = h(\theta_1, \dots, \theta_K)$, where each θ_k is of the form $Eg_k(X, Y)$, for some function $g_k(X, Y)$.

3.2 Propensity Score Adjusted Imputed Estimator

The data set may contain a fourth group for which both x_i and y_i are missing. Let $\delta_{4i} = 1$ if both x_i and y_i are missing. In this context, we interpret the probabilities (2.1) as conditional probabilities given that $\delta_{4i} = 0$. We apply the imputation procedure to $\{i : \delta_{4i} = 0\}$, as described in Section 2. We then apply a propensity score adjustment using a p -dimensional covariate \mathbf{v}_i known for all $i = 1, \dots, n$. Assume

$$P(\delta_{4i} = 0) = \exp(\phi_{40} + \phi'_{41} \mathbf{v}_i) [1 + \exp(\phi_{40} + \phi'_{41} \mathbf{v}_i)]^{-1} := p_{4i}(\phi_4). \quad (3.1)$$

Estimate the $(p + 1)$ -dimensional parameter $\phi'_4 = (\phi_{40}, \phi'_{41})'$ with $\hat{\phi}_4 = (\hat{\phi}_{40}, \hat{\phi}'_{41})'$ satisfying $S_4(\hat{\phi}_4) = \mathbf{0}$, where $S_4(\phi_4) = \sum_{i=1}^n (1, \mathbf{v}_i)' (1 - \delta_{4i} - p_{4i}(\phi_4))$. Then, let $\hat{p}_{4i} = p_{4i}(\hat{\phi}_4)$. The assumption (3.1) justifies the propensity score adjusted imputed estimator defined by

$$\hat{\theta}_{PSA-IMP} = \frac{1}{n} \left\{ \sum_{i=1}^n \delta_{1i} \frac{g(y_i, x_i)}{\hat{p}_{4i}} + \delta_{2i} \frac{\sum_{j=1}^J w_{2ij}(\hat{\phi}_2, \mathbf{y}_i^*) g(x_i, y_{ij}^*)}{\hat{p}_{4i}} + \delta_{3i} \frac{\sum_{j=1}^J w_{3ij}(\hat{\phi}_3, \mathbf{x}_i^*) g(x_{ij}^*, y_i)}{\hat{p}_{4i}} \right\}.$$

The propensity weights \hat{p}_{4i}^{-1} extrapolate the set $\{i : \delta_{1i} + \delta_{2i} + \delta_{3i} = 1\}$ onto the full sample $\{i = 1, \dots, n\}$. In supplement Section S3, we define an estimator of the variance of $\hat{\theta}_{PSA-Imp}$ as a straightforward extension of (3.5), and we verify through simulation that $\hat{\theta}_{PSA-Imp}$ and the corresponding variance estimator are approximately unbiased.

4. Simulation Study

We assess the finite-sample properties of the proposed estimator. We first compare the estimator of Section 2 to competitive alternatives. We then assess the properties of the variance estimator proposed in Section 3.

4.1 Comparison of Alternative Imputation Estimators

We consider two distributions for $F(y, x, \delta)$. For both, the parameter of interest is $\boldsymbol{\theta} = (EY, EX, V(Y), V(X), C(X, Y))'$, where $V(Y)$ (or $V(X)$) and $C(X, Y)$, respectively, denote the variance of Y (or X) and the correlation between X and Y . We compute the estimators for a Monte Carlo (MC) sample size of 500 and define $\boldsymbol{\theta}$ based on a separate simulation of size 500,000.

We compare the estimator proposed in Section 2 (abbreviated “Imp”) to three alternatives. To assess the impact of accounting for MNAR nonre-

4.1 Comparison of Alternative Imputation Estimators 20

sponse, we consider an ignorable (Ign) estimator that is essentially that of Chen and Yu (2016) and is obtained by setting $\phi = \mathbf{0}$ so that $w_{2ij}(\phi_2, \mathbf{q}_{yi}) = w_{3ij}(\phi_3, \mathbf{q}_{xi}) = J^{-1}$. We define parametric (Par) and non-parametric (NP) alternatives that involve implementing the 3 steps of Section 2.3, including estimation of ϕ , but generating the imputed values differently. For Par, we assume that $y_i = \beta_{0,y} + \beta_{1,y}x_i + \beta_{2,y}x_i^2 + \beta_{3,y}x_i^3 + e_{i,y}$, where $e_{i,y} \stackrel{iid}{\sim} N(0, \sigma_{e,y}^2)$, and likewise, $x_i = \beta_{0,x} + \beta_{1,x}y_i + \beta_{2,x}y_i^2 + \beta_{3,x}y_i^3 + e_{i,x}$, where $e_{i,x} \stackrel{iid}{\sim} N(0, \sigma_{e,x}^2)$. The imputed values for Par are $v_{ij}^* = \hat{v}_i + e_{vij}^*$, where for $\nu = x, y$, $e_{vij}^* \stackrel{iid}{\sim} N(0, \hat{\sigma}_{e,\nu}^2)$, \hat{v}_i is the predicted mean using the ordinary least squares coefficients $(\hat{\beta}_{0,\nu}, \hat{\beta}_{1,\nu}, \hat{\beta}_{2,\nu}, \hat{\beta}_{3,\nu})$, and $\hat{\sigma}_{e,\nu}^2 = (n - 4)^{-1} \sum_{i=1}^n (\nu_i - \hat{v}_i)^2$. For NP, we generate imputed values independently and with replacement from the set of observed values such that $P\{y_{ij}^* = y_k\} = K(x_k - x_i) [\sum_{\ell=1}^n \delta_{1\ell} K(x_\ell - x_i)]^{-1}$, $P\{x_{ij}^* = x_k\} = K(y_k - y_i) [\sum_{\ell=1}^n \delta_{1\ell} K(y_\ell - y_i)]^{-1}$, where $K(\cdot)$ is a Gaussian kernel with bandwidth defined by applying the R function `bw.ucv` to the sets $\{x_i : \delta_{1i} = 1\}$ and $\{y_i : \delta_{1i} = 1\}$ individually. Due to the adjustment for MNAR nonresponse, through estimation of ϕ , the Par and NP estimators proposed above are themselves innovations upon Kim (2011) and Wang and Chen (2009), respectively.

We define the *FlippedExp* simulation model by

$$y_i = h(x_i) + 1.25(1 + x_i)(\epsilon_i - 0.2), \quad \text{and} \quad \epsilon_i \stackrel{iid}{\sim} \text{Beta}(1, 4), \quad (4.2)$$

4.1 Comparison of Alternative Imputation Estimators 21

where $h(x_i) = \{2\exp(-2) - \exp(-2(x_i - 1))\}I[x_i < 2] + \{2\exp(2) - \exp(-2(x_i - 5))\}I[2 < x_i < 4] + \exp(-2(x_i - 3))I[4 < x_i < 6]$, $x_i \stackrel{iid}{\sim} Unif(0, 6)$ for $i = 1, \dots, n$, and $(\phi_{20}, \phi_{21}, \phi_{22}, \phi_{30}, \phi_{31}, \phi_{32}) = (-1, 0.033, 0.12, -0.800, 0.1, 0.033)$. We consider $n = 100, 1000$, and 5000 . The penalties $(\lambda_{n_1, y}, \lambda_{n_1, x})$ are $(0.2, 2)$, $(1, 10)$, and $(3, 30)$ for $n = 100, 1000$, and 5000 , respectively. They are based on a rule of $(\lambda_{n_1, y}, \lambda_{n_1, x}) \approx (0.1, 0.01)n^{6/9}$, determined from an exploratory analysis of simulated data using generalized cross-validation (Chen and Yu, 2016) and the relation between $\lambda_{n_1, y}$ and n in Lemma 1. We define $J \approx n^{0.5}$, giving $J = 10, 30$, and 70 for $n = 100, 1000$, and 5000 , respectively. The knots are the $k/(K + 1)$ quantiles of $\{x_i : \delta_{1i} + \delta_{2i} = 1, i = 1, \dots, n\}$ and $\{y_i : \delta_{1i} + \delta_{3i} = 1, i = 1, \dots, n\}$, where $k = 1, \dots, K$, and $K = 20, 30$, and 35 for $n = 100, 1000$, and 5000 , respectively. The values of K are based loosely on the rule of thumb, $K = \min\{n/4, 35\}$ (Ruppert, Wand, and Carroll, 2003).

Tables 2 and 3 contain the MC biases and RMSE's of the estimators of θ and ϕ , respectively, with the smallest absolute value among competitors in bold. For $n = 100$, variation from estimating additional parameters causes the RMSE of Imp to exceed those of Par and Ign, except for EX and $Cor(X, Y)$. For $n = 1000$, the Imp procedure is efficient. As n increases to 5000 , the efficiency NP improves. The Imp estimator of ϕ typically has

4.1 Comparison of Alternative Imputation Estimators 22

smallest absolute bias and RMSE.

To construct a model that better satisfies the assumptions of the Par estimator, we define the *Exp* configuration by (4.2) with $h(x_i) = \exp(2x_i)$, where $x_i \stackrel{iid}{\sim} Unif(-1, 1)$, and $(\phi_{20}, \phi_{21}, \phi_{22}, \phi_{30}, \phi_{31}, \phi_{32}) = (-0.9, 0.15, 0.2, -0.8, 0.15, 0.1)$. A rule of $\lambda_{n_1, y} = \lambda_{n_1, x} \approx n^{6/9}$ gives penalties of 20 and 100 for $n = 100$ and 1000, respectively. We define knots and τ_j the same as for FlippedExp.

4.1 Comparison of Alternative Imputation Estimators 23

Table 2: MC bias and RMSE of alternative estimators of θ for *FlippedExp*.

| | True | Bias | | | | RMSE | | | |
|------------------------------|--------|---------------|---------------|--------------|---------------|--------------|--------------|--------------|--------------|
| | | Ign | Imp | Par | NP | Ign | Imp | Par | NP |
| <u>$n = 100$</u> | | | | | | | | | |
| <i>EY</i> | 4.412 | -0.034 | 0.022 | 0.041 | -0.125 | 0.637 | 0.644 | 0.751 | 0.751 |
| <i>EX</i> | 3.000 | -0.838 | -0.001 | -0.010 | -0.014 | 0.862 | 0.180 | 0.186 | 0.196 |
| <i>V(Y)</i> | 42.865 | -0.182 | 0.619 | 3.469 | -1.610 | 5.545 | 6.235 | 14.460 | 5.575 |
| <i>V(X)</i> | 3.000 | 0.954 | -0.018 | 0.011 | -0.040 | 1.010 | 0.315 | 0.308 | 0.312 |
| <i>C(X, Y)</i> | 0.939 | -0.341 | -0.000 | -0.015 | -0.039 | 0.352 | 0.014 | 0.030 | 0.059 |
| <u>$n = 1000$</u> | | | | | | | | | |
| <i>EY</i> | 4.412 | -0.039 | 0.004 | -0.044 | 0.011 | 0.209 | 0.208 | 0.229 | 0.211 |
| <i>EX</i> | 3.000 | -0.839 | -0.000 | -0.003 | -0.002 | 0.841 | 0.053 | 0.056 | 0.053 |
| <i>V(Y)</i> | 42.865 | -0.564 | 0.003 | -0.114 | -0.038 | 1.472 | 1.430 | 1.809 | 1.441 |
| <i>V(X)</i> | 3.000 | 1.001 | -0.002 | 0.015 | -0.005 | 1.006 | 0.089 | 0.095 | 0.090 |
| <i>C(X, Y)</i> | 0.939 | -0.341 | 0.000 | -0.003 | -0.005 | 0.342 | 0.004 | 0.005 | 0.007 |
| <u>$n = 5000$</u> | | | | | | | | | |
| <i>EY</i> | 4.412 | -0.035 | 0.007 | -0.037 | 0.007 | 0.101 | 0.097 | 0.110 | 0.096 |
| <i>EX</i> | 3.000 | -0.838 | 0.001 | 0.001 | 0.001 | 0.838 | 0.025 | 0.025 | 0.025 |
| <i>V(Y)</i> | 42.865 | -0.574 | 0.011 | -0.190 | -0.010 | 0.833 | 0.654 | 0.833 | 0.637 |
| <i>V(X)</i> | 3.000 | 1.002 | -0.004 | 0.012 | -0.004 | 1.003 | 0.039 | 0.042 | 0.039 |
| <i>C(X, Y)</i> | 0.939 | -0.341 | 0.000 | -0.002 | -0.001 | 0.341 | 0.002 | 0.003 | 0.002 |

4.1 Comparison of Alternative Imputation Estimators 24

Table 3: MC bias and RMSE of alternative estimators of ϕ for *FlippedExp*.

| | n | True | <u>Bias</u> | | | <u>RMSE</u> | | |
|-------------|------|---------|----------------|----------------|----------------|---------------|--------|---------------|
| | | | Imp | Par | NP | Imp | Par | NP |
| ϕ_{20} | 100 | -1.0000 | 0.1008 | -0.2918 | -0.3151 | 1.0756 | 1.3114 | 1.0886 |
| ϕ_{21} | 100 | 0.0333 | -0.0853 | 0.1236 | 0.1046 | 0.5482 | 0.6408 | 0.5548 |
| ϕ_{22} | 100 | 0.1200 | 0.0277 | -0.0323 | -0.0176 | 0.1475 | 0.1712 | 0.1663 |
| ϕ_{30} | 100 | -0.8000 | -0.0923 | -0.0709 | -0.0852 | 1.3860 | 1.5112 | 1.7724 |
| ϕ_{31} | 100 | 0.1000 | -0.0023 | -0.0231 | -0.0477 | 0.7272 | 0.7588 | 0.8893 |
| ϕ_{32} | 100 | 0.0333 | 0.0061 | 0.0103 | 0.0223 | 0.1928 | 0.2037 | 0.2286 |
| ϕ_{20} | 1000 | -1.0000 | 0.0021 | -0.1877 | -0.0042 | 0.3077 | 0.4324 | 0.3001 |
| ϕ_{21} | 1000 | 0.0333 | -0.0049 | 0.0995 | -0.0059 | 0.1531 | 0.2242 | 0.1510 |
| ϕ_{22} | 1000 | 0.1200 | 0.0025 | -0.0248 | 0.0042 | 0.0400 | 0.0598 | 0.0402 |
| ϕ_{30} | 1000 | -0.8000 | 0.0075 | 0.0664 | 0.0549 | 0.3613 | 0.4875 | 0.3721 |
| ϕ_{31} | 1000 | 0.1000 | -0.0045 | -0.0379 | -0.0296 | 0.1824 | 0.2535 | 0.1892 |
| ϕ_{32} | 1000 | 0.0333 | 0.0014 | 0.0097 | 0.0078 | 0.0479 | 0.0669 | 0.0498 |
| ϕ_{20} | 5000 | -1.0000 | 0.0014 | -0.1442 | -0.0016 | 0.1411 | 0.2273 | 0.1402 |
| ϕ_{21} | 5000 | 0.0333 | -0.0012 | 0.0798 | 0.0000 | 0.0691 | 0.1198 | 0.0688 |
| ϕ_{22} | 5000 | 0.1200 | 0.0006 | -0.0199 | 0.0006 | 0.0180 | 0.0311 | 0.0179 |
| ϕ_{30} | 5000 | -0.8000 | 0.0057 | 0.0191 | 0.0221 | 0.1630 | 0.2210 | 0.1654 |
| ϕ_{31} | 5000 | 0.1000 | -0.0027 | -0.0100 | -0.0113 | 0.0829 | 0.1139 | 0.0843 |
| ϕ_{32} | 5000 | 0.0333 | 0.0008 | 0.0025 | 0.0030 | 0.0214 | 0.0295 | 0.0218 |

4.1 Comparison of Alternative Imputation Estimators 25

Results for Exp in Table 4 favor Par because the assumed cubic approximates the Exp function well. An exception is for $Var(X)$, where Imp has smaller RMSE than Par for $n = 100$ and $n = 1000$. Imp and Par are superior to NP in Table 4 due to the small sample size. Results for $\hat{\phi}$ and $n = 5000$ (omitted for brevity) lead to similar conclusions.

Table 4: Comparison of imputation procedures for Exp .

| | <u>Bias</u> | | | | | <u>RMSE</u> | | | |
|------------------------------|-------------|--------|---------------|---------------|---------------|--------------|--------------|--------------|-------|
| | True | Ign | Imp | Par | NP | Ign | Imp | Par | NP |
| <u>$n = 100$</u> | | | | | | | | | |
| EY | 1.813 | -0.012 | 0.000 | 0.007 | -0.025 | 0.203 | 0.202 | 0.199 | 0.223 |
| EX | 0.000 | -0.008 | 0.002 | 0.009 | -0.015 | 0.053 | 0.059 | 0.063 | 0.069 |
| $V(Y)$ | 3.613 | -0.211 | -0.161 | -0.069 | -0.197 | 0.685 | 0.675 | 0.648 | 0.846 |
| $V(X)$ | 0.333 | -0.084 | 0.001 | 0.024 | -0.005 | 0.089 | 0.032 | 0.095 | 0.033 |
| $C(X, Y)$ | 0.888 | -0.129 | 0.004 | 0.005 | -0.075 | 0.139 | 0.015 | 0.018 | 0.107 |
| <u>$n = 1000$</u> | | | | | | | | | |
| EY | 1.813 | -0.015 | -0.007 | -0.008 | -0.006 | 0.063 | 0.062 | 0.061 | 0.062 |
| EX | 0.000 | -0.009 | -0.001 | -0.001 | -0.003 | 0.018 | 0.019 | 0.019 | 0.019 |
| $V(Y)$ | 3.613 | -0.099 | -0.069 | -0.070 | -0.046 | 0.216 | 0.208 | 0.200 | 0.206 |
| $V(X)$ | 0.333 | -0.084 | -0.001 | 0.002 | -0.002 | 0.084 | 0.009 | 0.010 | 0.010 |
| $C(X, Y)$ | 0.888 | -0.124 | 0.001 | 0.001 | -0.010 | 0.125 | 0.006 | 0.005 | 0.014 |

4.2 Variance Estimator for Imputed Estimator 26

4.2 Variance Estimator for Imputed Estimator

Table 5: Properties of variance estimator for Imp for *FlippedExp*.

| | $n = 100$ | | | $n = 1000$ | | | $n = 5000$ | | |
|-------------|---|---------|------|---|--------|------|---|---------|------|
| | $V_{MC}(\hat{\theta})$ $\times 10^3$ | RB% | CR% | $V_{MC}(\hat{\theta})$ $\times 10^3$ | RB% | CR% | $V_{MC}(\hat{\theta})$ $\times 10^3$ | RB% | CR% |
| EY | 517.943 | -9.061 | 93.4 | 43.164 | 2.387 | 94.6 | 8.854 | 0.771 | 95.6 |
| EX | 37.112 | -10.713 | 93.8 | 3.102 | -0.216 | 94.0 | 0.633 | -2.467 | 94.8 |
| $V(Y)$ | 36090.740 | -11.852 | 93.0 | 1962.765 | -1.616 | 94.4 | 435.427 | -10.830 | 93.6 |
| $V(X)$ | 84.684 | 33.543 | 96.6 | 7.271 | 10.397 | 95.2 | 1.623 | -3.466 | 94.2 |
| $C(X, Y)$ | 0.185 | 29.068 | 96.0 | 0.014 | 0.177 | 92.8 | 0.003 | 12.921 | 95.6 |
| ϕ_{20} | 1121.574 | 5.837 | 97.4 | 100.770 | -3.255 | 94.2 | 19.616 | -2.589 | 95.2 |
| ϕ_{21} | 298.448 | 1.864 | 97.0 | 25.131 | -5.082 | 93.8 | 4.974 | -6.440 | 94.6 |
| ϕ_{22} | 21.847 | -0.004 | 96.0 | 1.759 | -5.536 | 93.6 | 0.343 | -6.006 | 93.4 |
| ϕ_{30} | 1707.902 | -5.553 | 97.4 | 131.168 | -3.995 | 95.0 | 25.236 | -2.336 | 94.2 |
| ϕ_{31} | 467.293 | -8.326 | 96.2 | 33.185 | -3.753 | 95.0 | 6.163 | 1.473 | 95.0 |
| ϕ_{32} | 32.245 | -5.849 | 97.0 | 2.351 | -5.450 | 95.0 | 0.411 | 5.636 | 95.4 |

Table 5 contains the MC variances ($V_{MC}(\hat{\theta})$) of the Imp estimators, the percent relative biases (RB%) of the variance estimator ($100(E_{MC}[\hat{V}] - V_{MC}(\hat{\theta}))/V_{MC}(\hat{\theta})$), where $E_{MC}[\hat{V}]$ denotes the MC mean of the variance estimator (3.5), and the percent of normal theory confidence intervals that

contain the true parameter values (CR%). For $n = 100$, the absolute RB% can exceed 15% and CR% can exceed 97%. For $n \in \{100, 5000\}$, the absolute RB% is below 15% and the CR% is within 2% of 95%.

5. Data Analysis

We analyze data from the 2011 Pet Demographic Survey (PDS), a national survey that collects information about pet ownership. The Iowa State Center for Survey Statistics and Methodology (CSSM) received the data as an agreement to plan for the 2017 survey. Variables of interest on the PDS include the number and type of pets owned, body types of those pets, and expenditures on veterinary services. We consider X^* , the sum of the most recent vet visit expenditures for a dog and cat combined, as a covariate for Y^* , the average vet visit expenditures in 2011 for dogs and cats. Table 6 has the number of observations for X^* and Y^* with four missing data patterns. We apply the propensity-score adjusted imputed estimator to estimate the veterinary expenditures for dogs and cats.

Table 6: Number of records in each group for pet data.

| Group | Count | Group | Count |
|-----------------------------|-------|----------------------------|-------|
| 1: X^* and Y^* observed | 3338 | 3: Only Y^* observed | 262 |
| 2: Only X^* observed | 2461 | 4: X^* and Y^* missing | 1169 |

The nature of the relationship between X^* and Y^* as well as extreme values preclude us from finding a quantile regression model that fits sample data well in the original scale. Further, the 75 zeros for X^* and 64 zeros for Y^* make a log transformation problematic. After exploring several transformations, including square root, cube root, and fifth root, we find that the cube root transformation allows us to construct a quantile regression model that appears adequate.

We apply the quantile regression procedure to first construct imputed values for $X = (X^*)^{1/3}$ and $Y = (Y^*)^{1/3}$ for groups 2 and 3. The generalized cross-validation criterion of Chen and Yu (2016) suggests $\lambda_{n_1,y} = 100$. The rule used for the *Exp* configuration of $\lambda_{n_1,y} \approx n_{123}^{6/9}$, where n_{123} is the number of observations in groups 1, 2, and 3, suggests $\lambda_{n_1,y} \approx 330$. At first, we tried the approximate mid-point of $\lambda_{n_1,y} \approx 200$ and obtained negative estimated quantiles for y_i for τ_1 and small values of x_i . Increasing the penalty to $\lambda_{n_1,y} = 300$ successfully avoided negatives. We present results for $\lambda_{n_1,y} = 300$. We use a fixed sequence of $\tau_j = j/(J + 1)$ for $j = 1, \dots, J$ with $J = 80 \approx n_{123}^{0.5}$. The fixed sequence avoids extreme quantiles and ensures that the data analysis is reproducible. (Chen and Yu (2016) compare results for fixed and random τ_j .) We define knots at the $k/(K + 1)$ quantiles of $\{x_i : \delta_{1i} + \delta_{2i} = 1 : i = 1, \dots, n\}$ and $\{y_i : \delta_{1i} + \delta_{3i} = 1 : i = 1, \dots, n\}$ for

$k = 1, \dots, K$, where $K = 35$.

We assess the model identification conditions (2.4) and (2.5) using the estimated functions $\hat{h}_y(\hat{\phi}_{22}, x)$ and $\hat{h}_x(\hat{\phi}_{31}, y)$ plotted in Figure 1. To construct the left plot in Figure 1, we first define an estimate of $h_y(\hat{\phi}_{22}, x_i)$ in equation (2.4) as the negative logarithm of the lowess regression of $\exp(-\hat{\phi}_{22}y_i)$ on x_i for the $\{i : \delta_i = 1\}$, where $\hat{\phi}_{22}$ is the estimated exponential tilting parameter in (2.7) obtained using the method described in Section 2.3. The right plot is constructed analogously, interchanging the roles of x_i and y_i and replacing $\hat{\phi}_{22}$ with $\hat{\phi}_{31}$. The nonlinearities seen in Figure 1 support the model identification conditions (2.4) and (2.5).

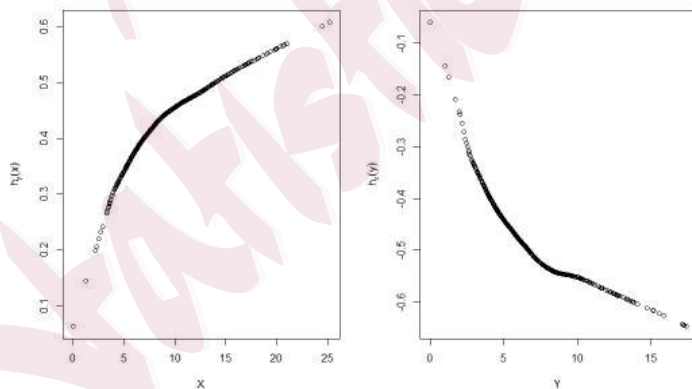


Figure 1: Estimated $\hat{h}_y(\hat{\phi}_{22}, x)$ (left) and $\hat{h}_x(\hat{\phi}_{31}, y)$ (right).

Table 7 gives estimates and corresponding standard errors for the propensity score model. The covariates, given in the column headings, are

selected with step-wise selection, starting with a model that contains all fully observed covariates and using the BIC criterion. The gender variable is 1 for females and 0 for males. The other covariates (defined in Section S3.3 of the supplement) are defined by ordered categories and are treated as numeric. The response variable is the indicator that unit i is not in group 4. Therefore, a positive coefficient is associated with a higher probability of providing a response. As such, we estimate that women with higher income and education who live alone or with one other person are more likely to provide a response to at least one of the questions about veterinary expenses.

Table 7: Estimated $\hat{\phi}_4$ and SE for propensity score model.

| | Intercept | Age | Gender | Income | Education | Household Size |
|------|-----------|---------|---------|---------|-----------|----------------|
| Est. | 0.16252 | 0.10355 | 0.38652 | 0.38395 | 0.21250 | -0.31212 |
| SE | 0.20897 | 0.02687 | 0.09208 | 0.03056 | 0.03671 | 0.05129 |

Table 8 contains estimates of ϕ_2 and ϕ_3 (obtained using (2.18) and (2.19)) along with associated standard errors (defined in (3.3)). The estimator of ϕ_{21} differs significantly from zero at the 5% level, but after accounting for x_i, y_i is no longer significantly associated with the response indicator, δ_{2i} . Interestingly, the estimate of ϕ_{31} is more than double the

standard error. The component of the model that accounts for nonignorable nonresponse is important for δ_{3i} .

Table 8: Estimates and standard errors for $\phi = (\phi'_2, \phi'_3)'$ for the pet data.

| | ϕ_{2j} Est. | ϕ_{2j} SE | ϕ_{3j} Est. | ϕ_{3j} SE |
|---------|------------------|----------------|------------------|----------------|
| $j = 0$ | 0.5136 | 0.2782 | -1.0677 | 0.4484 |
| $j = 1$ | -0.0561 | 0.0271 | -0.2590 | 0.1037 |
| $j = 2$ | -0.0810 | 0.0903 | 0.0609 | 0.0587 |

Table 9: Complete case and Imp-PSA estimators of selected parameters, along with standard errors for the Imp-PSA estimator.

| | EY | EX | $Var(Y)$ | $Var(X)$ | $Cor(X, Y)$ | EY^3 | EX^3 |
|------------------|-------|-------|----------|----------|-------------|---------|---------|
| Complete Case | 5.210 | 7.359 | 3.729 | 7.073 | 0.420 | 208.336 | 575.560 |
| SE Complete Case | 0.032 | 0.035 | 0.161 | 0.225 | 0.033 | 5.664 | 11.409 |
| Imp-PSA | 5.052 | 7.274 | 3.269 | 6.979 | 0.442 | 185.164 | 566.653 |
| SE Imp-PSA | 0.077 | 0.035 | 0.184 | 0.218 | 0.016 | 7.297 | 10.912 |

Table 9 compares the propensity-score adjusted imputed estimator (Imp-PSA) to the complete case estimator, which naively ignores missing values. The parameters EY^3 and EX^3 represent the mean expenditures in the original scale and are defined by $g(x, y) = y^3$ and $g(x, y) = x^3$, respectively.

We also estimate means and the correlation in the cube root scale. The comparison of complete case and imputed estimators suggests that ignoring the missing data would overstate the expenditures and understate the correlation between X and Y . As a result of the nonignorable nonresponse, the complete-case standard errors are also invalid. Imputation requires estimating additional parameters and can therefore leads to an increase in SE relative to the complete-case SE. The sample size for the complete-case estimator of the correlation is smaller than the sample size used to estimate the other parameters because the complete-case estimator of the correlation only uses pairs where both x_i and y_i are simultaneously observed.

6. Discussion

The theory, simulations, and data analysis demonstrate that the proposed semiparametric quantile regression imputation procedure is a viable method of constructing imputed values when the probability of responding may depend on the value of a missing response or covariate. We prove that the imputed estimator is asymptotically normal and verify through simulation that an estimate of the large sample covariance matrix has reasonable finite-sample properties. The simulations also show that failure to account for nonignorable nonresponse can lead to a severe bias. The squared bias of the

ignorable predictor can account for over 90% of MSE. In contrast, the ratio of the squared bias to MSE for the proposed (Imp) estimator is consistently below 1%. In our simulations, quantile regression is more robust than fully parametric imputation and more efficient than non-parametric imputation at small sample sizes. We do not have theoretical support for the superiority of semi-parametric quantile regression relative to non-parametric regression and therefore do not expect these results to hold broadly. A further advantage of quantile regression over the non-parametric estimator of Wang and Chen (2009) is that quantile regression permits a linearization-based variance estimator. In the application, the proposed procedure allows us to use one type of veterinary expenditure to impute the other, while allowing for nonignorable nonresponse and modeling complex patterns in the data. Further, we develop a propensity score adjustment to incorporate a set for which neither veterinary expenditure is observed.

In this paper, we use a fully parametric model for the response probability. As demonstrated in Robins and Ritov (1997), identification for nonignorable nonresponse is elusive without any restrictions. Nonetheless, relaxing the parametric assumptions of the response probability model, along the lines of Shao and Wang (2016), is a possible avenue for future work.

REFERENCES

In principle, our approach of modeling the conditional distribution of the covariate given a response extends to multivariate covariates. One must ensure that the quantile regression model adequately describes each full univariate conditional and that identification conditions are satisfied. We define an identification condition for multivariate covariates in Section S4 of the supplement. An alternative approach for missing covariates is to use Bayes rule to deduce $f(x | y)$ from a specification of $f(y | x)$ and $f(x)$ (Yang and Kim, 2017). Our preliminary studies suggest that an extension of Yang and Kim (2017) to nonignorable nonresponse and quantile regression is a promising direction for future work.

Supplementary Materials have details omitted for brevity.

Acknowledgements We acknowledge the grant NSF MMS-1733572.

References

- BARROW, D.L. AND SMITH, P.W. (1978). Asymptotic properties of best $L_2[0, 1]$ approximation by spline with variable knots. *Quart. Appl. Math.* **36**, 293–304.
- BERG, E., AND YU, C. (2019). Semiparametric quantile regression imputation for a complex survey with application to the Conservation Effects Assessment Project, *Survey Methodology* **45**, 249–270.

REFERENCES

- CHEN, S. AND YU, C.L. (2016). Parameter estimation through semiparametric quantile regression imputation. *Electron. J. Statist.*, **10**, 3621–3647.
- CHENG, P. E. (1994). Nonparametric estimation of mean functionals with data missing at random. *J. Amer. Statist. Assoc.*, **89**, 81–87.
- FANG, F., ZHAO, J., AND SHAO, J. (2018). Imputation-based adjusted score equations in generalized linear models with nonignorable missing covariate values. *Statistica Sinica*, **28**, 1677–1701.
- KIM, J. K. (2011). Parametric fractional imputation for missing data analysis. *Biometrika*, **98**, 119–132.
- KIM, J. K. AND SHAO, J. (2014). *Statistical Methods for Handling Incomplete Data*. Chapman & Hall, Raton, FL.
- KIM, J.K AND YU, C.L. (2011). *A Semiparametric Estimation of Mean Functionals With Nonignorable Missing Data*, *JASA*, **106**, 157–165.
- KOENKER, R. (2005). *Quantile Regression*. Cambridge University Press, New York.
- CHANG, T. AND KOTT, P. S. (2008). *Using calibration weighting to adjust for nonresponse under a plausible model*. *Biometrika*, **95**, 555–571.
- MIAO, W. AND TCHETGEN TCHETGEN, E. J. (2016). *On varieties of doubly robust estimators under missingness not at random with a shadow variable*. *Biometrika*, **103**, 475–482.
- MORIKAWA, K. AND KIM, J.K. (2017). Semiparametric adaptive estimation with nonignorable

REFERENCES

- nonresponse data. *arXiv preprint arXiv:1612.09207* .
- RIDDLES, M.K., KIM, J.K., AND IM, J. (2016). A Propensity-score-adjustment method for nonignorable nonresponse. *Journal of Survey Statistics and Methodology*, **4**, 215–245.
- ROBINS, J. M., AND RITOV, Y. A. (1997). Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semi-parametric models. *Statistics in medicine*, **16**, 285–319.
- RUBIN, D.B. (1987). *Multiple Imputation for Nonresponse in Surveys*. New York: Wiley.
- RUPPERT, D., WAND, M. P., & CARROLL, R. J. (2003). *Semiparametric regression (No. 12)*. Cambridge university press.
- SHAO, J., AND WANG, L. (2016). *Semiparametric inverse propensity weighting for nonignorable missing data*. *Biometrika*, **103**, 175–187.
- TANG, G., LITTLE, R. J., AND RAGHUNATHAN, T.E. (2003). Analysis of multivariate missing data with nonignorable nonresponse. *Biometrika*, **90**, 747–764.
- WANG, D. AND CHEN, S.X. (2009). Empirical likelihood for estimating equations with missing values. *Ann. Statist.* **37**, 490–517.
- WANG, S., SHAO, J., AND KIM, J. K. (2014). An instrumental variable approach for identification and estimation with nonignorable nonresponse. *Statistica Sinica*, 1097-1116.
- YANG, S. AND KIM, J. K. (2017). A semiparametric inference to regression analysis with missing covariates in survey data. *Statistica Sinica*, 261-285.

REFERENCES

- YOSHIDA, T. (2013) Asymptotics for Penalized Spline Estimators in Quantile Regression. *Communications in Statistics - Theory and Methods*.
- ZHAO, J., AND MA, Y. (2019). A Versatile Estimation Procedure without Estimating the Nonignorable Missingness Mechanism. arXiv preprint arXiv:1907.03682.
- ZHAO, J., AND SHAO, J. (2015). Semiparametric pseudo-likelihoods in generalized linear models with nonignorable missing data. *Journal of the American Statistical Association*, **110**, 1577–1590.