

# EMPIRICAL LIKELIHOOD ESTIMATION FOR SAMPLES WITH NONIGNORABLE NONRESPONSE

Fang Fang, Quan Hong, and Jun Shao

*GE Consumer Finance, Eli Lilly and Company, and  
University of Wisconsin-Madison and East China Normal University*

*Abstract:* Nonresponse is very common in survey sampling. Nonignorable nonresponse, a response mechanism in which the response probability of a survey variable  $Y$  depends directly on the value of  $Y$  regardless of whether  $Y$  is observed or not, is the most difficult type of nonresponse to handle. The population mean estimators ignoring the nonrespondents typically have heavy biases. This paper studies an empirical likelihood-based estimation method, with samples under nonignorable nonresponse, when an observed auxiliary categorical variable  $Z$  is available. The likelihood is semiparametric: we assume a parametric model on the response mechanism and the conditional probability of  $Z$  given  $Y$ , and a nonparametric model on the distribution of  $Y$ . When the number of  $Z$  categories is not small, a pseudo empirical likelihood method is applied to reduce the computational intensity. Asymptotic distributions of the proposed population mean estimators are derived. For variance estimation, we consider a bootstrap procedure and its consistency is established. Some simulation results are provided to assess the finite sample performance of the proposed estimators.

*Key words and phrases:* Empirical likelihood, Pseudo likelihood, Nonignorable nonresponse, Sample survey, Semiparametric likelihood, Stratified samples.

## 1. Introduction

Nonresponse is a common phenomenon in sample surveys. Let  $Y$  be a variable of interest having nonrespondents and  $Z$  be a covariate with no nonresponse. If the propensity  $P(\delta = 1|Y, Z)$ , where  $\delta$  is the response indicator for  $Y$ , depends not only on  $Z$  and observed  $Y$ , but also on unobserved  $Y$ , then the nonresponse mechanism is nonignorable. Nonignorable nonresponse creates a great challenge in the estimation of the mean of  $Y$  based on incomplete survey data. Ignoring the dependence of nonresponse probability on unobserved  $Y$  typically leads to heavy bias.

Greenlees, Reece, and Zieschang (1982) studied maximum likelihood estimators for survey data with nonignorable nonresponse, based on a parametric model on the propensity  $P(\delta = 1|Y, Z)$  and a parametric (normal) model on  $L(Y|Z)$ , the distribution of  $Y$  conditional on  $Z$ . However, parametric models (especially normal models) on  $L(Y|Z)$  for survey data are often not valid. In fact, Greenlees, Reece, and Zieschang (1982) admitted that the normality assumption on  $L(Y|Z)$  was not valid for the data in their example, even though their method was better than the method of ignoring the fact that nonresponse was nonignorable.

On the other hand, it is impossible to develop a pure nonparametric method that produces a consistent estimator of the mean of  $Y$  in the presence of nonignorable nonresponse. Thus, some semiparametric methods assuming a parametric model on one of  $P(\delta = 1|Y, Z)$  and  $L(Y|Z)$  have been proposed in the literature. Tang, Little, and Rachunathan (2003) developed a likelihood method by assuming a parametric model on  $L(Y|Z)$ ; they assumed that  $P(\delta = 1|Y, Z) = P(\delta = 1|Y)$  but otherwise is nonparametric. Qin, Leung, and Shao (2002) proposed an empirical likelihood method by assuming a parametric model on  $P(\delta = 1|Y, Z)$  and a nonparametric model on  $L(Y|Z)$ ; the resulting estimator of the mean of  $Y$  is similar to the estimator obtained by weighting each respondent by the inverse of an estimated propensity  $P(\delta = 1|Y, Z)$  (Robins, Rotnitzky, and Zhao (1994)). For survey data, finding a suitable parametric model for  $P(\delta = 1|Y, Z)$  is much easier than finding an appropriate parametric model for  $L(Y|Z)$ . However, the estimation of  $P(\delta = 1|Y, Z)$  is still difficult under a parametric assumption on  $P(\delta = 1|Y, Z)$  because of the presence of unobserved  $Y$  values.

In many survey problems the covariate  $Z$  is categorical, e.g., age group, sex, race, education level, type of industry etc., while the main variable  $Y$  is continuous. If there is an appropriate parametric model on the conditional distribution  $L(Z|Y)$  given  $Y$  (e.g., the logistic model), then we can improve the approach in Qin, Leung, and Shao (2002). The purpose of this paper is to study an empirical likelihood approach under parametric models on  $P(\delta = 1|Y, Z)$  and  $L(Z|Y)$  with a discrete  $Z$ , and under a nonparametric model on the distribution of  $Y$ . Our approach works for a stratified sampling design with a superpopulation within each stratum, which is commonly used in practice. Furthermore, we study a pseudo empirical likelihood to reduce the amount of computation when the num-

ber of  $Z$  categories is not small. Although losing some efficiency, the estimators based on the pseudo empirical likelihood are still consistent and asymptotically normal. Note that the same technique has been applied to the case of ignorable nonresponse (Fang, Hong, and Shao (2009)).

This paper is organized as follows. Section 2 presents details on the sampling design and model, and gives results for estimation without imputation. In addition to the derivation of empirical likelihood estimators, their consistency and asymptotic normality are established. Section 3 discusses the pseudo empirical likelihood estimators. Section 4 considers variance estimation by bootstrapping. In Section 5, we consider two imputation methods related to the results in Sections 2 and 3. Section 6 examines by simulation the finite sample performance of the proposed estimators, under some response patterns and models. The Appendix contains proofs or sketched proofs.

## 2. Empirical Likelihood Approach

We consider the following sampling design commonly used in such business surveys as the Current Employment Survey conducted by the U.S. Bureau of Labor Statistics (Wolter, Shao, and Huff (1998)), the Transportation Annual Survey conducted by the U.S. Census Bureau (Census Bureau (1987)), and the Financial Farm Survey conducted by Statistics Canada (Rancourt (1999)). The finite population  $\mathcal{P}$  is stratified into  $H$  (a fixed positive integer) strata and samples are taken independently across the strata. Within each stratum, a large number of units are either independently sampled with replacement according to a probability sampling plan, or selected as a simple random sample without replacement with a negligible sampling fraction. According to the sampling plan, survey weights  $\{\omega_i\}$  are constructed so that for any set of values  $\{x_i\}$ ,

$$E_{\mathcal{S}} \left( \sum_{i \in \mathcal{S}} \omega_i x_i \right) = \sum_{i \in \mathcal{P}} x_i,$$

where  $\mathcal{S}$  is the sample and  $E_{\mathcal{S}}$  is the expectation with respect to sampling.

Let  $Y$  be the variable of interest in the survey and  $Z$  be a categorical covariate taking values in  $\{z_1, \dots, z_s\}$ . We assume that values of  $(Y, Z)$  are iid from a superpopulation within each stratum, and are independent across strata. To present the main idea, we first consider the special case of one stratum so that

the subscript for stratum is omitted.

Under the superpopulation model (within each stratum),  $Y$  has an unknown nonparametric distribution  $F$ , and we assume a parametric probability function

$$P(Z = z|Y = y) = f(y, z, \beta), \quad (1.1)$$

where  $\beta$  is an unknown parameter vector and  $f$  is a known function. For each sampled unit, the  $Z$  value is always observed, but the  $Y$  value may be a nonrespondent. We assume that the probability that an individual responds on  $Y$  can depend on both  $Y$  and  $Z$  according to

$$\phi(Y, Z, \gamma) = P(\delta = 1|Y, Z), \quad (1.2)$$

where  $\delta$  is the response indicator for  $Y$ ,  $\phi$  is a known function, and  $\gamma$  is an unknown parameter vector.

Without loss of generality, we assume that the first  $r$  sampled units are respondents and the rest of  $n - r$  sampled units are nonrespondents. Thus, the observed data set is

$$\{(Y_i, Z_i), i = 1, \dots, r\} \cup \{Z_i, i = r + 1, \dots, n\}.$$

Let  $p_i = dF(Y_i)$  be the point mass that  $F$  places on  $Y_i$ . For observed  $Y_i$ , the likelihood is

$$\phi(Y_i, Z_i, \gamma)f(Y_i, Z_i, \beta)p_i.$$

For a nonrespondent  $Y_i$ , the likelihood is

$$\int [1 - \phi(y, Z_i, \gamma)]f(y, Z_i, \beta)dF(y).$$

Together with the survey weights (see, e.g., Chen and Qin (1993)), we obtain the following log-likelihood for the entire sample

$$\begin{aligned} & \sum_{i=1}^r w_i \log(\phi(Y_i, Z_i, \gamma)f(Y_i, Z_i, \beta)p_i) \\ & + \sum_{i=r+1}^n w_i \log \left( \int [1 - \phi(y, Z_i, \gamma)]f(y, Z_i, \beta)dF(y) \right), \end{aligned}$$

where  $w_i = \omega_i/N$  and  $N$  is the finite population size. The use of  $w_i$ , instead of  $\omega_i$ , does not change the maximization of the log-likelihood over the parameters. Since  $Z$  takes values  $z_1, \dots, z_s$ , this log-likelihood can be written as

$$\sum_{i=1}^r w_i \log(\phi(Y_i, Z_i, \gamma)f(Y_i, Z_i, \beta)p_i) + \sum_{j=1}^s a_j \log(\pi_j), \quad (1.3)$$

where  $a_j = \sum_{i=r+1}^n w_i I_{\{Z_i=z_j\}}$ ,  $I_A$  is the indicator function of the event  $A$ , and

$$\pi_j = P(\delta = 0, Z = z_j) = \int [1 - \phi(y, z_j, \gamma)]f(y, z_j, \beta)dF(y).$$

Note that  $\pi_j$  is a function of  $\gamma$ ,  $\beta$ , and  $F$ . Maximizing (1.3) over  $\gamma$ ,  $\beta$ , and  $F$  is equivalent to maximizing (1.3) over  $\gamma$ ,  $\beta$ ,  $p_i$ 's, and  $\pi_j$ 's subject to

$$p_i \geq 0, \quad \sum_{i=1}^r p_i = 1, \quad \pi_j = \sum_{i=1}^r p_i [1 - \phi(Y_i, z_j, \gamma)]f(Y_i, z_j, \beta), \quad j = 1, \dots, s. \quad (1.4)$$

By introducing Lagrange multipliers, we can derive that

$$p_i = \frac{w_i}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]}, \quad i = 1, \dots, r, \quad (1.5)$$

where  $\hat{N}_r = \sum_{k=1}^r w_k$  and  $\lambda_j$ 's are Lagrange multipliers satisfying

$$\sum_{i=1}^r \frac{w_i [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} = 0, \quad j = 1, \dots, s. \quad (1.6)$$

Treating  $p_i$  in (1.5) as a function of  $\beta$ ,  $\gamma$ ,  $\pi = (\pi_1, \dots, \pi_s)$ , and  $\lambda = (\lambda_1, \dots, \lambda_s)$ , and substituting  $p_i$  into (1.3), the profile log-likelihood with Lagrange multipliers is

$$\begin{aligned} l(\beta, \gamma, \pi, \lambda) &= \sum_{i=1}^r w_i \log(\phi(Y_i, Z_i, \gamma)f(Y_i, Z_i, \beta)) + \sum_{j=1}^s a_j \log(\pi_j) \\ &+ \sum_{i=1}^r w_i \log \left( \frac{w_i}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} \right), \end{aligned}$$

Differentiating  $l(\beta, \gamma, \pi, \lambda)$  with respect to  $\pi$ ,  $\lambda$ ,  $\beta$ , and  $\gamma$ , and setting the partial

derivatives to 0, we have

$$\frac{a_j}{\pi_j} + \sum_{i=1}^r \frac{w_i \lambda_j}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} = 0, \quad j = 1, \dots, s, \quad (1.7)$$

$$\sum_{i=1}^r \frac{w_i [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} = 0, \quad j = 1, \dots, s, \quad (1.8)$$

$$\sum_{i=1}^r \left\{ \frac{w_i \partial \log f(Y_i, Z_i, \beta)}{\partial \beta} - \frac{w_i \sum_{j=1}^s \lambda_j (1 - \phi(Y_i, z_j, \gamma)) \partial f(Y_i, z_j, \beta) / \partial \beta}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} \right\} = 0, \quad (1.9)$$

$$\sum_{i=1}^r \left\{ \frac{w_i \partial \log \phi(Y_i, Z_i, \gamma)}{\partial \gamma} + \frac{w_i \sum_{j=1}^s \lambda_j \partial \phi(Y_i, z_j, \gamma) / \partial \gamma f(Y_i, z_j, \beta)}{\hat{N}_r + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma))f(Y_i, z_j, \beta) - \pi_j]} \right\} = 0. \quad (1.10)$$

From (1.5), (1.7), and the fact that  $\sum_{i=1}^r p_i = 1$ , we have

$$\lambda_j = -a_j / \pi_j, \quad j = 1, \dots, s. \quad (1.11)$$

Let  $(\hat{\beta}, \hat{\gamma}, \hat{\pi}, \hat{\lambda})$  be a solution to equations (1.8)-(1.11). The maximum empirical likelihood estimator (MELE) of  $(\beta, \gamma)$  is  $(\hat{\beta}, \hat{\gamma})$ , and the MELE of  $F$  is the empirical distribution  $\hat{F}$  putting mass  $\hat{p}_i$  at  $Y_i$ ,  $i = 1, \dots, r$ , where  $\hat{p}_i$  is given by (1.5) with  $(\beta, \gamma, \pi, \lambda)$  replaced by  $(\hat{\beta}, \hat{\gamma}, \hat{\pi}, \hat{\lambda})$ . If the parameter of interest is the finite population mean  $\bar{Y} = \sum_{i \in \mathcal{P}} Y_i / N$ , its MELE is

$$\hat{Y} = \sum_{i=1}^r \hat{p}_i Y_i. \quad (1.12)$$

If the parameter of interest is the cell mean  $\bar{Y}_j$ , the finite population mean of  $Y$  given  $Z = z_j$ , the MELE is

$$\hat{Y}_j = \sum_{i=1}^r \hat{p}_i f(Y_i, z_j, \hat{\beta}) Y_i / \sum_{i=1}^r \hat{p}_i f(Y_i, z_j, \hat{\beta}). \quad (1.13)$$

Let  $\theta = (\beta, \gamma, \pi)$ ,  $\hat{\theta} = (\hat{\beta}, \hat{\gamma}, \hat{\pi})$ , and  $\hat{\nu} = \hat{\lambda} / \hat{N}_r + J / (1 - \sum_{j=1}^s \hat{\pi}_j)$ , where  $J$  is the  $s$ -vector of ones. The following result shows that  $(\hat{\theta}, \hat{\nu})$  converges to  $(\theta_0, 0)$ , where  $\theta_0 = (\beta_0, \gamma_0, \pi_0)$  is the true value of  $(\beta, \gamma, \pi)$ . Also,  $\hat{Y}$  and  $\hat{Y}_j$  are consistent for  $\bar{Y}$  and  $\bar{Y}_j$ , respectively. Furthermore,  $(\hat{\theta}, \hat{\nu})$ ,  $\hat{Y}$ , and  $\hat{Y}_j$  are

asymptotically normal. The proof is given in the Appendix.

**THEOREM 1.** *Assume the following.*

(i) *The sample from the finite population is selected with replacement according to a probability sampling plan or selected as a simple random sample without replacement. The values of  $(Y, Z)$  in the population is iid from a superpopulation according to (1.1)-(1.2) with a nonparametric  $Y$ -marginal  $F$ .*

(ii) *As  $n \rightarrow \infty$ ,  $N \rightarrow \infty$ ,  $n/N \rightarrow 0$ ,  $\max_{i \leq N} w_i = O(1/n)$ , and  $n \sum_{i=1}^N w_i/N \rightarrow d$  for some constant  $d$ .*

(iii)  *$f(y, z, \beta)$  and  $\phi(y, z, \gamma)$  are twice continuously differentiable in  $\beta$  and  $\gamma$  for any  $y$  and  $z$ , and  $\|\frac{\partial \log f(y, z, \beta)}{\partial \beta}\|^2$ ,  $\|\frac{\partial \log \phi(y, z, \gamma)}{\partial \gamma}\|^2$ ,  $\|\frac{\partial^2 f(y, z, \beta)}{\partial \beta \partial \beta^\tau}\|^2$ ,  $\|\frac{\partial^2 \phi(y, z, \gamma)}{\partial \gamma \partial \gamma^\tau}\|^2$ ,  $\|\frac{\partial f(y, z_j, \beta)}{\partial \beta}\|^3$ ,  $\|\frac{\partial \phi(y, z_j, \gamma)}{\partial \gamma}\|^3$ ,  $\|[\frac{\partial f(y, z_j, \beta)}{\partial \beta}][\frac{\partial f(y, z_k, \beta)}{\partial \beta}]^\tau\|^2$ ,  $\|[\frac{\partial \phi(y, z_j, \gamma)}{\partial \gamma}][\frac{\partial \phi(y, z_k, \gamma)}{\partial \gamma}]^\tau\|^2$ , and  $\|[\frac{\partial f(y, z_j, \beta)}{\partial \beta}][\frac{\partial \phi(y, z_k, \gamma)}{\partial \gamma}]^\tau\|^2$  are bounded by some integrable functions in a neighborhood of  $\beta_0$  and  $\gamma_0$ ,  $j, k = 1, \dots, s$ .*

(iv) *For any nonzero vector  $c \in \mathcal{R}^{p+q}$ , the value of  $c^\tau \begin{pmatrix} \partial \log f(y, z_j, \beta_0)/\partial \beta \\ \partial \log \phi(y, z_j, \gamma_0)/\partial \gamma \end{pmatrix}$  depends on  $j$ , where  $p$  and  $q$  are the dimensions of  $\beta$  and  $\gamma$ .*

(v)  *$\theta_0$  is a unique root of  $E[g(Y, \theta)|\delta = 1] = 0$  and  $E[g(Y, \theta_0)g(Y, \theta_0)^\tau|\delta = 1]$  is positive definite, where  $g = (g_1, \dots, g_s)^\tau$  and*

$$g_j(y, \theta) = \frac{(1 - \sum_{k=1}^s \pi_k) [(1 - \phi(y, z_j, \gamma))f(y, z_j, \beta) - \pi_j]}{\sum_{k=1}^s \phi(y, z_k, \gamma)f(y, z_k, \beta)}. \quad (1.14)$$

(vi)  *$\phi(y, z, \gamma)$  has a positive lower bound.*

*Then, there exists a sequence  $\{\hat{\theta}, \hat{\nu}, n = 1, 2, \dots\}$  such that as  $n \rightarrow \infty$ ,*

$$P(\hat{\theta} \text{ is a solution to (1.8)-(1.10)}) \rightarrow 1, \quad (1.15)$$

$$\sqrt{n} \begin{pmatrix} \hat{\nu} - 0 \\ \hat{\theta} - \theta_0 \end{pmatrix} \rightarrow_d N(0, \Sigma), \quad (1.16)$$

*where the probability  $P$  and  $\rightarrow_d$  (convergence in distribution) are with respect to the sampling and the superpopulation, and  $\Sigma$  is a positive definite matrix. Furthermore, if functions  $\|y \frac{\partial f(y, z_j, \beta)}{\partial \beta}\|^2$  and  $\|y \frac{\partial \phi(y, z_j, \gamma)}{\partial \gamma}\|^2$  are bounded by some*

integrable functions in a neighborhood of  $\beta_0$  and  $\gamma_0$  for each  $j$ , then

$$\sqrt{n}(\hat{Y} - \bar{Y}) \rightarrow_d N(0, \sigma^2) \quad \text{and} \quad \sqrt{n}(\hat{Y}_j - \bar{Y}_j) \rightarrow_d N(0, \sigma_j^2), \quad j = 1, \dots, s, \quad (1.17)$$

where  $\sigma^2$  and  $\sigma_j^2$  are some constants.

In condition (v),  $E[g(Y, \theta)|\delta = 1] = 0$  has a unique root  $\theta_0$  is equivalent to  $\pi_j = \int [1 - \phi(y, z_j, \gamma)] f(y, z_j, \beta) dF(y)$  is uniquely defined by  $(\beta, \gamma)$ , i.e., a condition of identifiability of the  $\pi$  by  $(\beta, \gamma)$ .

We now consider the stratified sample described in the beginning of this section. If  $(\beta, \gamma)$  in conditions (1.1) and (1.2) has different values in different strata, then we can solve (1.7)-(1.10) within each stratum to obtain an estimator of  $(\beta, \gamma)$  for each stratum. If  $(\beta, \gamma)$  is common for all strata, then constraint (1.4) is within each stratum, the sums in (1.7)-(1.8) are over each stratum, and the sums in (1.9)-(1.10) are over all strata. In any case, the marginal distribution of  $Y$  for stratum  $h$  is the empirical distribution putting mass  $\hat{p}_i$  at  $Y_i$  with  $i$  in stratum  $h$ ; the estimator of  $\bar{Y}$  is the weighted average of estimators given by (1.12) over all strata with the weights  $W_h = N_h/N$ , where  $N_h$  is the population size for stratum  $h$  and  $N = \sum_h N_h$ ; the estimator of  $\bar{Y}_j$  is the ratio of the averages of the numerators and denominators in (1.13) with the weights  $W_h$ . Theorem 1 still holds if all conditions are given within each stratum and  $n_h/n$  converges to a positive constant, where  $n_h$  is the sample size in stratum  $h$  and  $n = \sum_h n_h$ .

### 3. Pseudo Empirical Likelihood

When  $s$  (the number of  $Z$  categories) is not small, numerical solutions to (1.8)-(1.11) may be computationally intensive. Hence, we apply the idea of pseudo likelihood (Gong and Samaniego (1981)). That is, we substitute each  $\pi_j$  in (1.8)-(1.11) by a consistent estimator  $\tilde{\pi}_j$ . Note that consistent estimators of  $\pi_j$ 's are easy to construct. For example, we may estimate  $\pi_j$  by

$$\tilde{\pi}_j = \frac{\sum_{i=1}^n w_i I_{\{\delta_i=0, Z_i=z_j\}}}{\sum_{i=1}^n w_i}. \quad (1.18)$$

Let  $\tilde{\pi} = (\tilde{\pi}_j, j = 1, \dots, s)$ ,  $\tilde{\lambda}_j = -a_j/\tilde{\pi}_j$ , and  $\tilde{\lambda} = (\tilde{\lambda}_1, \dots, \tilde{\lambda}_s)$ . Maximizing the pseudo empirical likelihood  $l(\beta, \gamma, \tilde{\pi}, \tilde{\lambda})$  over  $(\beta, \gamma)$  results in the maximum pseudo empirical likelihood estimator (MPELE)  $(\tilde{\beta}, \tilde{\gamma})$ . Note that the MPELE



is different from MELE since  $\tilde{\pi}$  is not  $\hat{\pi}$ . However, we can directly establish the consistency and asymptotic normality of the MPELE.

Let  $\tilde{p}_i$  be the estimator of  $p_i$  obtained by using (1.5) with  $\beta$ ,  $\gamma$ ,  $\pi_j$ , and  $\lambda_j$  replaced by  $\tilde{\beta}$ ,  $\tilde{\gamma}$ ,  $\tilde{\pi}_j$ , and  $\tilde{\lambda}_j$ , respectively. Because the MPELE is used,  $\sum_{i=1}^r \tilde{p}_i \neq 1$ , although  $\sum_{i=1}^r \tilde{p}_i \rightarrow_p 1$ . The MPELE of  $\bar{Y}$  is

$$\tilde{\bar{Y}} = \sum_{i=1}^r \tilde{p}_i Y_i / \sum_{i=1}^r \tilde{p}_i, \quad (1.19)$$

and the MPELE of  $\bar{Y}_j$  is

$$\tilde{\bar{Y}}_j = \sum_{i=1}^r \tilde{p}_i f(Y_i, z_j, \tilde{\beta}) Y_i / \sum_{i=1}^r \tilde{p}_i f(Y_i, z_j, \tilde{\beta}). \quad (1.20)$$

Estimators under stratified sampling can be obtained as described in the end of Section 2, with the sums in (1.18) within each stratum.

The following result shows that the MPELE is consistent and asymptotically normal.

**THEOREM 2.** *Assume the conditions in Theorem 1. There exists a sequence  $\{\tilde{\beta}, \tilde{\gamma}, n = 1, 2, \dots\}$  such that, as  $n \rightarrow \infty$ ,*

$$P \left( \frac{\partial l(\tilde{\beta}, \tilde{\gamma}, \tilde{\pi}, \tilde{\lambda})}{\partial(\beta, \gamma)} = 0 \right) \rightarrow 1 \text{ and } \sqrt{n} \begin{pmatrix} \tilde{\beta} - \beta_0 \\ \tilde{\gamma} - \gamma_0 \end{pmatrix} \rightarrow_d N(0, \Sigma_p), \quad (1.21)$$

where  $\Sigma_p$  is a positive definite matrix. Furthermore,

$$\sqrt{n}(\tilde{\bar{Y}} - \bar{Y}) \rightarrow_d N(0, \sigma_p^2) \quad \text{and} \quad \sqrt{n}(\tilde{\bar{Y}}_j - \bar{Y}_j) \rightarrow_d N(0, \sigma_{pj}^2), \quad j = 1, \dots, s, \quad (1.22)$$

where  $\sigma_p^2$  and  $\sigma_{pj}^2$  are some constants.

#### 4. Variance Estimation by Bootstrapping

It is a common practice in sample surveys to report a variance estimate for each estimate of the parameter of interest. We focus on the most commonly used estimators, the mean estimators  $\hat{\bar{Y}}$ ,  $\tilde{\bar{Y}}$  in (1.12) and (1.19), and the cell mean estimators  $\hat{\bar{Y}}_j$ ,  $\tilde{\bar{Y}}_j$  in (1.13) and (1.20). Because the formulation of these estimators is complicated, it is difficult to derive an analytic form of their asymptotic variances,  $\sigma^2$ ,  $\sigma_j^2$  in (1.17), and  $\sigma_p^2$ ,  $\sigma_{pj}^2$  in (1.22). Thus, we apply the bootstrap

method that consists of the following steps. In the following,  $\hat{\eta}$  denotes any of  $\hat{\beta}$ ,  $\hat{\gamma}$ ,  $\hat{\pi}$ ,  $\hat{\nu}$ ,  $\hat{Y}$ ,  $\hat{Y}_j$ ,  $\hat{\beta}$ ,  $\tilde{\gamma}$ ,  $\tilde{\pi}$ ,  $\tilde{Y}$ , and  $\tilde{Y}_j$ .

1. Within stratum  $h$ , draw a simple random sample of size  $n_h$  with replacement from the set of sampled units (respondents or nonrespondents). Carry out this procedure independently across strata. For each unit in the bootstrap sample, the bootstrap data are the  $Z$  and  $Y$  values (if the  $Y$  is missing, the bootstrap datum is treated as missing) and their survey weights.
2. Compute  $\hat{\eta}^*$ , which is the same as  $\hat{\eta}$  but with the original data replaced by the bootstrap data generated in Step 1.
3. Repeat the previous steps independently  $B$  times and obtain  $\hat{\eta}^{*1}, \dots, \hat{\eta}^{*B}$ . Estimate the variance of  $\hat{\eta}$  by the sample variance of  $\hat{\eta}^{*1}, \dots, \hat{\eta}^{*B}$ .

The following result establishes the asymptotic validity of the bootstrap.

**THEOREM 3.** *Assume the conditions in Theorem 1.*

(i) *Let (1.8\*)-(1.11\*) be the bootstrap analog of (1.8)-(1.11). Then there exists a sequence  $\{\hat{\theta}^*, \hat{\nu}^*, n = 1, 2, \dots\}$  such that, as  $n \rightarrow \infty$ ,*

$$P_*(\hat{\theta}^* \text{ is a solution to (1.8*)-(1.10*)}) \rightarrow_p 1, \quad (1.23)$$

$$\sqrt{n} \begin{pmatrix} \hat{\nu}^* - \hat{\nu} \\ \hat{\theta}^* - \hat{\theta} \end{pmatrix} \rightarrow_{d^*} N(0, \Sigma), \quad (1.24)$$

where  $\Sigma$  is given in (1.16),  $P_*$  denotes the bootstrap probability conditional on the data, and  $\vartheta_n^* \rightarrow_{d^*} \vartheta$  means  $P_*(\vartheta_n^* \in B) - P(\vartheta \in B) \rightarrow_p 0$  for any Borel set  $B$ . Furthermore,

$$\sqrt{n}(\hat{Y}^* - \hat{Y}) \rightarrow_{d^*} N(0, \sigma^2) \quad \text{and} \quad \sqrt{n}(\hat{Y}_j^* - \hat{Y}_j) \rightarrow_{d^*} N(0, \sigma_j^2), \quad (1.25)$$

where  $\sigma^2$  and  $\sigma_j^2$  are defined in (1.17).

(ii) *Let  $\tilde{\pi}^* = (\tilde{\pi}_1^*, \dots, \tilde{\pi}_s^*)$ , with  $\tilde{\pi}_j^*$  being the bootstrap analog of  $\tilde{\pi}_j$  in (1.18). Then there exists a sequence  $\{\tilde{\beta}^*, \tilde{\gamma}^*, n = 1, 2, \dots\}$  such that, as  $n \rightarrow \infty$ ,*

$$P_* \left( \frac{\partial l^*(\tilde{\beta}^*, \tilde{\gamma}^*, \tilde{\pi}^*, \tilde{\lambda}^*)}{\partial(\beta, \gamma)} = 0 \right) \rightarrow_p 1 \quad \text{and} \quad \sqrt{n} \begin{pmatrix} \tilde{\beta}^* - \tilde{\beta} \\ \tilde{\gamma}^* - \tilde{\gamma} \end{pmatrix} \rightarrow_{d^*} N(0, \Sigma_p), \quad (1.26)$$

where  $\Sigma_p$  is given in (1.21). Further,

$$\sqrt{n}(\tilde{Y}^* - \tilde{Y}) \rightarrow_{d^*} N(0, \sigma_p^2) \quad \text{and} \quad \sqrt{n}(\tilde{Y}_j^* - \tilde{Y}_j) \rightarrow_{d^*} N(0, \sigma_{pj}^2), \quad (1.27)$$

where  $\sigma_p^2$  and  $\sigma_{pj}^2$  are defined in (1.22).

## 5. Imputation

Imputation is often carried out for practical reasons (Kalton and Kasprzyk (1986)). After imputation, estimates of parameters are computed by treating imputed values as observed data and using the standard formulas for the case of no nonresponse. In this section we consider imputation for the estimation of the population mean  $\bar{Y}$  and the population cell mean  $\bar{Y}_j$ . Let  $\hat{Y}_i = Y_i$  if  $Y_i$  is a respondent and  $\hat{Y}_i$  be an imputed value if  $Y_i$  is a nonrespondent. After imputation, the population mean  $\bar{Y}$  and cell mean  $\bar{Y}_j$  are estimated by

$$\hat{Y}_I = \sum_{i=1}^n w_i \hat{Y}_i, \quad (1.28)$$

$$\hat{Y}_{jI} = \sum_{i=1}^n w_i \hat{Y}_i I_{\{Z_i=z_j\}} / \sum_{i=1}^n w_i I_{\{Z_i=z_j\}}, \quad (1.29)$$

respectively. Under stratified sampling, (1.28)-(1.29) should be modified as described at the end of Section 2.

The naive mean imputation method imputes each nonrespondent with  $Z = z_j$  by the cell sample mean  $\sum_{i=1}^r w_i Y_i I_{\{Z_i=z_j\}} / \sum_{i=1}^r w_i I_{\{Z_i=z_j\}}$ . The naive random imputation method imputes each nonrespondent with  $Z = z_j$  by a random sample with replacement from respondents with  $Z = z_j$ , where each  $Y_i$  with  $Z_i = z_j$  has probability  $w_i I_{\{Z_i=z_j\}} / \sum_{i=1}^r w_i I_{\{Z_i=z_j\}}$  to be selected,  $i = 1, \dots, r$ . The population mean estimators based on the naive imputation methods are inconsistent since they do not consider the difference between the respondents and the nonrespondents.

Using the MELE estimators developed in Section 2, we consider the following two imputation procedures.

1. Empirical Likelihood Mean Imputation. For each nonrespondent with  $Z = z_j$ , the imputed  $Y$  value is

$$\frac{\sum_{i=1}^r \hat{p}_i [1 - \phi(Y_i, z_j, \hat{\gamma})] f(Y_i, z_j, \hat{\beta}) Y_i}{\sum_{i=1}^r \hat{p}_i [1 - \phi(Y_i, z_j, \hat{\gamma})] f(Y_i, z_j, \hat{\beta})}.$$

2. Empirical Likelihood Random Imputation. Each nonrespondent with  $Z = z_j$  is imputed by a random sample with replacement from all respondents, where the probability of each  $Y_i$  to be selected is

$$\frac{\hat{p}_i[1 - \phi(Y_i, z_j, \hat{\gamma})]f(Y_i, z_j, \hat{\beta})}{\sum_{i=1}^r \hat{p}_i[1 - \phi(Y_i, z_j, \hat{\gamma})]f(Y_i, z_j, \hat{\beta})}.$$

For stratified sampling, imputation should be carried out within each stratum.

Similarly, using the MPELE estimators developed in Section 3, we can develop Pseudo Empirical Likelihood Mean Imputation and Random Imputation. They are similar to the Empirical Likelihood Mean Imputation and Random Imputation that we described above. We just need to replace  $\hat{\beta}$ ,  $\hat{\gamma}$ , and  $\hat{p}_i$  by  $\tilde{\beta}$ ,  $\tilde{\gamma}$ , and  $\tilde{p}_i$ , respectively.

The following result shows that the estimators of  $\bar{Y}$  and  $\bar{Y}_j$  based on these four imputation procedures are consistent and asymptotically normal.

**THEOREM 4:** *Under the conditions of Theorem 1, for empirical likelihood mean imputation, empirical likelihood random imputation, pseudo empirical likelihood mean imputation, or pseudo empirical likelihood random imputation,*

$$\sqrt{n}(\hat{Y}_I - \bar{Y}) \rightarrow_d N(0, \sigma_I^2), \quad \text{and} \quad \sqrt{n}(\hat{Y}_{jI} - \bar{Y}_j) \rightarrow_d N(0, \sigma_{jI}^2), \quad j = 1, \dots, s,$$

where  $\sigma_I^2$  and  $\sigma_{jI}^2$  are some constants.

The asymptotic variances  $\sigma_I^2$  and  $\sigma_{jI}^2$  do not have simple analytic forms. Variance estimation can be carried out using the bootstrap procedure described in Section 4. It should be emphasized that, to address the variability caused by imputation, nonrespondents in each bootstrap data set must be imputed using the bootstrap data and the same imputation method as that used to impute the original data set, as suggested by Shao and Sitter (1996).

## 6. Simulation Results

In this section, we report on simulation of the finite-sample properties of the MELE, MPELE, the empirical likelihood imputation, and the pseudo empirical likelihood imputation. We created a finite population similar to the Current Establishment Survey conducted by the U.S. Bureau of Labor Statistics. We chose four different industries as four strata with sizes  $N_1 = 3370$ ,  $N_2 = 2910$ ,  $N_3 =$

5430, and  $N_4 = 4110$ . The variable  $Y$  is the total pay for each establishment and values of  $Y$  in stratum  $h$  were generated from a superpopulation  $F_h$ . The form of  $F_h$  was chosen to be the gamma distribution and  $F_1 = \Gamma(43, 0.20)$ ,  $F_2 = \Gamma(42, 0.19)$ ,  $F_3 = \Gamma(38, 0.20)$ , and  $F_4 = \Gamma(50, 0.17)$ , where  $\Gamma(a, b)$  denotes the gamma distribution with shape parameter  $a$  and scale parameter  $b$ . The parameters in  $F_h$ 's were chosen to match the mean and variance of a data set from the Current Establishment Survey.

The covariate  $Z \in \{1, 2, 3, 4, 5\}$  was generated by the logistic model

$$P(Z = j|Y = y) = \frac{\exp\{\beta_j + \beta_5 y\}}{1 + \sum_{k=1}^4 \exp\{\beta_k + \beta_5 y\}}, \quad j = 1, 2, 3, 4,$$

$$P(Z = 5|Y = y) = \frac{1}{1 + \sum_{k=1}^4 \exp\{\beta_k + \beta_5 y\}},$$

where  $\beta_k$ ,  $k = 1, 2, 3, 4, 5$ , are unknown parameters whose values in the simulation are 0.25, 0.5, 0.75, 1, and  $-0.1$ , respectively.

The sampling plan was stratified simple random sampling. In each stratum, the sampling fraction was 0.05. For each sampled unit, the  $Y$  respondent was generated according to the response probability function

$$P(\delta = 1|Y = y, Z = j) = \frac{\exp\{-10 - j + \gamma y\}}{1 + \exp\{-10 - j + \gamma y\}}$$

with a parameter  $\gamma = 1.8$  or 2, or

$$P(\delta = 1|Y = y, Z = j) = \frac{\exp\{10 + j + \gamma y\}}{1 + \exp\{10 + j + \gamma y\}}$$

with  $\gamma = -1.4$ . The following table lists the response rate for each  $Z$  and the mean response rate  $E[P(\delta = 1|Z)]$ .

$\gamma$	1.8	2	-1.4
$P(\delta = 1 Z = 1)$	0.888	0.951	0.457
$P(\delta = 1 Z = 2)$	0.803	0.910	0.621
$P(\delta = 1 Z = 3)$	0.697	0.842	0.751
$P(\delta = 1 Z = 4)$	0.560	0.749	0.856
$P(\delta = 1 Z = 5)$	0.469	0.675	0.908
$E[P(\delta = 1 Z)]$	0.651	0.804	0.756

For each of the three  $\gamma$ , Table 1-3 respectively reports the relative bias (RB) and variance (VAR) of the MELE estimators in (1.12) and (1.13), the MPELE estimators in (1.19) and (1.20), the naive estimators that simply ignore nonrespondents, and the imputation estimators in (1.28) and (1.29) based on empirical, pseudo empirical, or naive mean imputation and random imputation. We also report their bootstrap variance estimators (Vboot) based on the bootstrap replication size  $B = 200$ , the coverage probabilities (CP) and the lengths (LEN) of the bootstrap confidence intervals of the form

$$\text{point estimate} \pm 1.96\sqrt{\text{Vboot}}$$

that approximately have nominal coverage probability 95%.

Table 4 reports the mean and the variance (VAR) of the parameter estimates. Table 5 reports the ratios of the mean squared errors. Each MPELE is compared with its counterpart; that is,  $\tilde{Y}$  in (1.19) is compared with  $\hat{Y}$  in (1.12),  $\tilde{Y}_j$  in (1.20) is compared with  $\hat{Y}_j$  in (1.13), and  $\tilde{Y}_l$  in (1.28) (or  $\hat{Y}_{j_l}$  in (1.29)) with pseudo empirical likelihood mean (or random) imputation is compared with  $\hat{Y}_l$  in (1.28) (or  $\hat{Y}_{j_l}$  in (1.29)) with empirical likelihood mean (or random) imputation described in Section 5.

The computation was done using MATLAB in a UNIX at the Department of Statistics, University of Wisconsin-Madison. For each  $\gamma$  and a single simulation, it took about 12 seconds to compute the MELE, MPELE, and imputed estimates for  $\tilde{Y}$  and  $\tilde{Y}_l$ ,  $l = 1, \dots, 5$ . Because of the bootstrap, however, each simulation with a given  $\gamma$  took about 40 minutes. For each  $\gamma$ , we ran the simulation 250 times.

The simulation results can be summarized as follows.

1. In all cases, the proposed population mean and population cell mean estimators based on empirical likelihood or pseudo empirical likelihood (with imputation or not) performed well in terms of the relative bias (less than 1%) and variance, while the naive methods had heavy relative biases up to 10.31%.
2. The bootstrap variance estimate for our proposed estimators worked well in most cases in terms of its bias and the coverage probability of the bootstrap confidence interval. For the naive estimators, the coverage probability of

the confidence interval was very low.

3. Although the MPELE estimators required less computational intensities, they were less efficient in terms of larger MSE compared with the MELE estimators. Most of the MSE ratios were greater than 1 (Table 5). For the estimators without imputation, the ratios were all greater than 5, and some of them were even greater than 20. The lengths of confidence intervals of the MPELE estimators were all greater than those of the MELE estimators, especially for the estimators without imputation.
4. Although the variances of the  $\beta$  and  $\gamma$  parameter estimates were a little bit large, the estimation of the population mean and population cell means, which is our major interest, was still good.

### Acknowledgments

The authors wish to thank the referees for their comments and suggestions. Jun Shao's work was partially supported by the NSF Grants DMS-0404535 and SES-0705033.

### References

- Census Bureau (1987). Noncertainty sample specification. BSR-87 Action Memo D.06, the U.S. Census Bureau.
- Chen, J. and Qin, J. (1993). Empirical likelihood estimation for finite population and the effective usage of auxiliary information. *Biometrika* **80**, 107-116.
- Fang, F., Hong, Q., and Shao, J. (2009). A pseudo empirical likelihood approach for stratified samples with nonresponse. *Annals of Statistics* **37**, 371-393.
- Gong, G. and Samaniego, F. (1981). Pseudo maximum likelihood estimation: theory and application. *Annals of Statistics* **9**, 861-869.
- Greenlees, J.S., Reece, W.S., and Zieschang, K.Y. (1982). Imputation of missing values when the probability of response depends on the variable being imputed. *Journal of the American Statistical Association* **77**, 251-261.
- Kalton, G. and Kasprzyk, D. (1986). The treatment of missing data. *Survey Methodology* **12**, 1-16.

- Qin, J., Leung, D., and Shao, J. (2002). Estimation with survey data under non-ignorable nonresponse or informative sampling. *Journal of the American Statistical Association* **97**, 193-200.
- Rancourt, E. (1999). Estimation with nearest-neighbor imputation at Statistics Canada. *Proceedings of the Survey Research Methods Section*, American Statistical Association, 446-451.
- Robins, J.M., Rotnitzky, A., and Zhao, L.P. (1994). Estimation of regression coefficients when some regressors are not always observed. *Journal of American Statistical Association* **89**, 846-86.
- Shao, J. and Sitter, R.R. (1996). Bootstrap for imputed survey data. *Journal of the American Statistical Association* **91**, 1278-1288.
- Tang, G., Little, R.J., and Raghunathan, T.E. (2003). Analysis of multivariate missing data with nonignorable nonresponse. *Biometrika* **90**, 747-764.
- Wolter, K., Shao, J., and Huff, L. (1998). Variance estimation for the Current Employment Statistics Program. *Proceedings of the Section on Survey Research Methods*. American Statistical Association, 775-780.

## Appendix

The proofs in this appendix are for the special case of one stratum. The proof for the case of  $H > 1$  is similar.

LEMMA 1: Let  $\psi(x, \theta)$  be a function satisfying  $E(\psi(x, \theta)) = 0$ . Assume that  $E[\psi(x, \theta_0)\psi^\tau(x, \theta_0)]$  is positive definite,  $\partial\psi(x, \theta)/\partial\theta$  is continuous in a neighborhood of  $\theta_0$ ,  $\|\partial\psi(x, \theta)/\partial\theta\|$  and  $\|\psi(x, \theta)\|^3$  are bounded by some integrable functions in the neighborhood. Under the conditions (i)-(ii) of Theorem 1, with probability 1, there exists a  $\nu$  such that  $\sum_{i=1}^n \frac{w_i \psi(x_i, \theta)}{1 + \nu^\tau \psi(x_i, \theta)} = 0$ . Furthermore, let  $l(\theta, \nu) = -\sum_{i=1}^n w_i \log\{1 + \nu^\tau \psi(x_i, \theta)\}$ , then in an  $O_p(n^{-1/3})$  neighborhood of  $\theta_0$ ,  $\nu = \nu(\theta)$  is a function of  $\theta$ , and  $l(\theta, \nu(\theta))$  attains its maximum value at some interior point of the ball  $\|\theta - \theta_0\| \leq n^{-1/3}$ .

PROOF. Consider the problem of maximizing  $\sum_{i=1}^n w_i \log p_i$  under the constraints  $p_i \geq 0$ ,  $\sum_{i=1}^n p_i = 1$ , and  $\sum_{i=1}^n p_i \psi(x_i, \theta) = 0$ . Since  $E(\psi(x, \theta)) = 0$ , it follows from the arguments of Owen (1990) that, as  $n \rightarrow \infty$ , 0 is contained in the convex hull of  $\{\psi(x_i, \theta), i = 1, \dots, n\}$  with probability 1. For a given  $\theta$ , when 0 is inside of the convex hull, a unique maximum exists, which can be found via Lagrange multipliers as follows. Let

$$H = \sum_{i=1}^n w_i \log p_i + \lambda(1 - \sum_{i=1}^n p_i) - \sum_{k=1}^n w_k \nu^\tau \sum_{i=1}^n p_i \psi(x_i, \theta)$$



where  $\lambda$  and  $\nu$  are Lagrange multipliers. Taking derivatives with respect to  $p_i$ , we have

$$\frac{\partial H}{\partial p_i} = \frac{w_i}{p_i} - \lambda - \sum_{k=1}^n w_k \nu^\tau \psi(x_i, \theta) = 0.$$

Then

$$\sum_{i=1}^n p_i \frac{\partial H}{\partial p_i} = \sum_{i=1}^n w_i - \lambda = 0,$$

which leads to

$$p_i = \frac{w_i}{1 + \nu^\tau \psi(x_i, \theta)} \bigg/ \sum_{i=1}^n w_i$$

with  $\nu$  satisfying

$$\sum_{i=1}^n \frac{w_i \psi(x_i, \theta)}{1 + \nu^\tau \psi(x_i, \theta)} = 0.$$

This proves the first conclusion for Lemma 1.

Note that it is necessary that  $0 \leq p_i \leq 1$ , which implies that  $\nu$  and  $\theta$  must satisfy  $1 + \nu^\tau \psi(x_i, \theta) \geq w_i / \sum_{i=1}^n w_i$  for each  $i$ . For fixed  $\theta$ , let  $D_\theta = \{\nu : 1 + \nu^\tau \psi(x_i, \theta) \geq w_i / \sum_{i=1}^n w_i\}$ ;  $D_\theta$  is convex and closed, and it is bounded when 0 is inside the convex hull of the  $\psi(x_i, \theta)$ 's. Notice that

$$\frac{\partial}{\partial \nu} \left\{ \sum_{i=1}^n w_i \frac{\psi(x_i, \theta)}{1 + \nu^\tau \psi(x_i, \theta)} \right\} = - \sum_{i=1}^n w_i \frac{\psi(x_i, \theta) \psi^\tau(x_i, \theta)}{[1 + \nu^\tau \psi(x_i, \theta)]^2}$$

is negative definite. By the inverse function theorem  $\nu = \nu(\theta)$  is a differentiable function. Let  $\psi_i = \psi(x_i, \theta)$ . Since

$$0 = \sum_{i=1}^n w_i \frac{\psi_i}{1 + \nu^\tau \psi_i} = \sum_{i=1}^n w_i \left( \psi_i - \frac{\psi_i \psi_i^\tau \nu}{1 + \nu^\tau \psi_i} \right),$$

we have

$$\left\| \sum_{i=1}^n w_i \psi_i \right\| = \|\nu\| \sum_{i=1}^n w_i \frac{\psi_i \psi_i^\tau}{1 + \nu^\tau \psi_i} \geq \frac{\|\nu\|}{1 + \|\nu\| \psi^*} \sum_{i=1}^n w_i \psi_i \psi_i^\tau,$$

where  $\psi^* = \max_{1 \leq i \leq n} \|\psi_i\| = o(n^{1/3})$ , a.s., by lemma 3 of Owen (1990) and the condition  $E\|\psi(x_i, \theta)\|^3 < \infty$ . When  $\|\theta - \theta_0\| = n^{-1/3}$ ,  $\left\| \sum_{i=1}^n w_i \psi_i \right\| = O_p(n^{-1/3})$  and  $\sum_{i=1}^n w_i \psi_i \psi_i^\tau = O_p(1)$ . Then

$$\frac{\|\nu\|}{1 + \|\nu\| o(n^{1/3})} = O_p(n^{-1/3}),$$

and  $\nu = \nu(\theta) = O_p(n^{-1/3})$ . Furthermore, similar to the proof of Owen (1990), we have

$$\nu(\theta) = \left[ \sum_{i=1}^n w_i \psi_i \psi_i^\tau \right]^{-1} \left[ \sum_{i=1}^n w_i \psi_i \right] + o_p(n^{-1/3}), \quad \|\theta - \theta_0\| = n^{-1/3}.$$

It follows from the arguments of Qin and Lawless (1994) that

$$l(\theta, \nu(\theta)) < l(\theta_0, \nu(\theta_0)) \text{ in probability,}$$

if  $\|\theta - \theta_0\| = n^{-1/3}$ . Then  $l(\theta, \nu(\theta))$  attains its local maximum value at some interior point of the ball  $\|\theta - \theta_0\| \leq n^{-1/3}$ .

PROOF OF THEOREM 1. Let  $\nu_j = \lambda_j / \hat{N}_r + (1 - \sum_{j=1}^s \pi_j)^{-1}$ ,  $\nu = (\nu_1, \dots, \nu_s)^\tau$  and  $g(Y_i, \theta)$  be defined by (1.14). Then

$$\begin{aligned} 1 + \nu^\tau g(Y_i, \theta) &= 1 + \frac{\sum_{j=1}^s \left[ \frac{(1 - \sum_{j=1}^s \pi_j) \lambda_j}{\sum_{k=1}^r w_k} + 1 \right] [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &= 1 + \frac{\sum_{j=1}^s [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &\quad + \frac{\sum_{j=1}^s (1 - \sum_{j=1}^s \pi_j) \lambda_j [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{k=1}^r w_k \sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &= \frac{\sum_{j=1}^s [f(Y_i, z_j, \beta) - \pi_j]}{\sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &\quad + \frac{(1 - \sum_{j=1}^s \pi_j) \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{k=1}^r w_k \sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &= \frac{1 - \sum_{j=1}^s \pi_j}{\sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &\quad + \frac{(1 - \sum_{j=1}^s \pi_j) \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{k=1}^r w_k \sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \\ &= \frac{(1 - \sum_{j=1}^s \pi_j) \{ \sum_{k=1}^r w_k + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j] \}}{\sum_{k=1}^r w_k \sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta)} \end{aligned}$$

Then the function  $l(\beta, \gamma, \pi, \lambda)$  can be written as

$$\begin{aligned} &l(\beta, \gamma, \pi, \lambda) \\ &= \sum_{i=1}^r w_i \log(\phi_i f_i) + \sum_{j=1}^s a_j \log(\pi_j) + \sum_{i=1}^r w_i \log \frac{w_i (1 - \sum_{j=1}^s \pi_j)}{(1 + \nu^\tau g(Y_i, \theta)) \sum_{k=1}^r w_k \sum_{j=1}^s \phi_{ij} f_{ij}} \\ &= - \sum_{i=1}^r w_i \log \{1 + \nu^\tau g(Y_i, \theta)\} + \sum_{i=1}^r w_i \log(\phi_i f_i) - \sum_{i=1}^r w_i \log \left( \sum_{j=1}^s \phi_{ij} f_{ij} \right) \\ &\quad + \sum_{j=1}^s a_j \log(\pi_j) + \sum_{i=1}^r w_i \log(1 - \sum_{j=1}^s \pi_j) + \sum_{i=1}^r w_i \log \frac{w_i}{\sum_{k=1}^r w_k}, \end{aligned}$$

where  $\phi_i = \phi(Y_i, Z_i, \gamma)$ ,  $f_i = f(Y_i, Z_i, \beta)$ ,  $\phi_{ij} = \phi(Y_i, z_j, \gamma)$ , and  $f_{ij} = f(Y_i, z_j, \beta)$ . Therefore  $l(\beta, \gamma, \pi, \lambda)$  is equal to

$$l(\theta, \nu) = l_1(\theta, \nu) + l_2(\theta) + l_3(\theta)$$

plus a term that does not depend on the parameters, where

$$\begin{aligned} l_1(\theta, \nu) &= -\sum_{i=1}^r w_i \log \{1 + \nu^\tau g(Y_i, \theta)\}, \\ l_2(\theta) &= \sum_{i=1}^r w_i \log \left( \phi(Y_i, Z_i, \gamma) f(Y_i, Z_i, \beta) \right) \\ &\quad - \sum_{i=1}^r w_i \log \left( \sum_{j=1}^s \phi(Y_i, z_j, \gamma) f(Y_i, z_j, \beta) \right), \\ l_3(\theta) &= \sum_{j=1}^s a_j \log \pi_j + \sum_{i=1}^r w_i \log \left( 1 - \sum_{j=1}^s \pi_j \right). \end{aligned}$$

Notice that

$$\sum_{i=1}^r \frac{w_i g_j(Y_i, \theta)}{1 + \nu^\tau g(Y_i, \theta)} = \sum_{i=1}^r \frac{w_i [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]}{\sum_{k=1}^r w_k + \sum_{j=1}^s \lambda_j [(1 - \phi(Y_i, z_j, \gamma)) f(Y_i, z_j, \beta) - \pi_j]} \sum_{k=1}^r w_k.$$

Then constraint (1.6) becomes

$$\sum_{i=1}^r \frac{w_i g(Y_i, \theta)}{1 + \nu^\tau g(Y_i, \theta)} = 0. \quad (1.30)$$

Since  $E[g(Y_i, \theta) | \delta = 1] = 0$ , it follows from (1.30) and Lemma 1 that in an  $O_p(n^{-1/3})$  neighborhood, we can determine uniquely a differentiable implicit function

$$\nu = \nu(\theta) = O_p(n^{-1/3}) \text{ if } \|\theta - \theta_0\| \leq O_p(n^{-1/3}),$$

and

$$l_1(\theta, \nu(\theta)) < l_1(\theta_0, \nu(\theta_0)) \text{ in probability,} \quad (1.31)$$

if  $\theta$  is in the set  $B_n = \{\theta : \|\theta - \theta_0\| = n^{-\frac{1}{3}}\}$ .

For  $l_2(\beta, \gamma) = l_2(\theta)$ , denote  $E_c$  as the conditional expectation of  $(Y, Z)$  given  $\delta = 1$ , which is

$$E_c = \sum_{z \in \{z_1, \dots, z_s\}} \int \frac{\phi(y, z, \gamma_0) f(y, z, \beta_0)}{P(\delta = 1)} dF(y).$$

Then

$$\begin{aligned} \frac{\partial l_2(\beta_0, \gamma_0)}{\partial \beta} &= \sum_{i=1}^r w_i \frac{\partial f(Y_i, Z_i, \beta_0) / \partial \beta}{f(Y_i, Z_i, \beta_0)} - \sum_{i=1}^r w_i \frac{\sum_{j=1}^s \phi(Y_i, z_j, \gamma_0) \frac{\partial f(Y_i, z_j, \beta_0)}{\partial \beta}}{\sum_{j=1}^s \phi(Y_i, z_j, \gamma_0) f(Y_i, z_j, \beta_0)} \\ &\rightarrow_p P(\delta = 1) E_c \frac{\partial f(y, z, \beta_0) / \partial \beta}{f(y, z, \beta_0)} - P(\delta = 1) E_c \frac{\sum_{j=1}^s \phi(y, z_j, \gamma_0) \frac{\partial f(y, z_j, \beta_0)}{\partial \beta}}{\sum_{j=1}^s \phi(y, z_j, \gamma_0) f(y, z_j, \beta_0)} \\ &= 0 \end{aligned}$$

By similar calculation, we can show that  $\frac{\partial l_2(\beta_0, \gamma_0)}{\partial \gamma} \rightarrow_p 0$ ,  $\frac{\partial^2 l_2(\beta^*, \gamma^*)}{\partial (\beta, \gamma)^2} \rightarrow_p -U$ , where  $(\beta^*, \gamma^*)$  is

between  $(\beta, \gamma)$  and  $(\beta_0, \gamma_0)$ , and  $U$  is defined as

$$\begin{aligned} U &= \begin{pmatrix} U_{11} & U_{12} \\ U_{12}^\tau & U_{22} \end{pmatrix}, \\ U_{11} &= \int \left( \sum_j \frac{\phi_j}{f_j} \left( \frac{\partial f_j}{\partial \beta} \right)^2 - \frac{(\sum_j \phi_j \frac{\partial f_j}{\partial \beta})^2}{\sum_j \phi_j f_j} \right) dF(y), \\ U_{12} &= \int \left( \sum_j \frac{\partial f_j}{\partial \beta} \frac{\partial \phi_j}{\partial \gamma^\tau} - \frac{[\sum_j \phi_j \frac{\partial f_j}{\partial \beta}][\sum_j \frac{\partial \phi_j}{\partial \gamma^\tau} f_j]^\tau}{\sum_j \phi_j f_j} \right) dF(y), \\ U_{22} &= \int \left( \sum_j \frac{f_j}{\phi_j} \left( \frac{\partial \phi_j}{\partial \gamma} \right)^2 - \frac{(\sum_j \frac{\partial \phi_j}{\partial \gamma} f_j)^2}{\sum_j \phi_j f_j} \right) dF(y), \end{aligned} \quad (1.32)$$

where  $\phi_j = \phi(y, z_j, \gamma_0)$  and  $f_j = f(y, z_j, \beta_0)$ . For any nonzero vector  $c = (c_1, c_2)$ , by Cauchy's inequality,

$$c^\tau U c = \int \left\{ \sum \left( \sqrt{\frac{\phi}{f}} \cdot c_1^\tau \frac{\partial f}{\partial \beta} + \sqrt{\frac{f}{\phi}} \cdot c_2^\tau \frac{\partial \phi}{\partial \gamma} \right)^2 - \frac{\left( \sum \phi \cdot c_1^\tau \frac{\partial f}{\partial \beta} + \sum f \cdot c_2^\tau \frac{\partial \phi}{\partial \gamma} \right)^2}{\sum \phi f} \right\} dF(y) \geq 0$$

If the equation holds, then  $c_1^\tau \frac{\partial \log f(y, z_j, \beta_0)}{\partial \beta} + c_2^\tau \frac{\partial \log \phi(y, z_j, \gamma_0)}{\partial \gamma} = c(y)$  for  $j = 1, \dots, s$ , a.s., which contradicts condition (iv). Hence the equation does not hold and  $U$  is positive definite. By central limit theorem and delta method, we can show that  $\sqrt{n} \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)}$  is asymptotical normal.

When  $\theta \in B_n$ , we have  $(\beta, \gamma) = (\beta_0, \gamma_0) + n^{-\frac{1}{3}} u^\tau$ ,  $\|u\| \leq 1$ ,

$$\begin{aligned} l_2(\theta) - l_2(\theta_0) &= l_2(\beta, \gamma) - l_2(\beta_0, \gamma_0) \\ &= n^{-\frac{1}{3}} u^\tau \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)} + \frac{1}{2} n^{-\frac{2}{3}} u^\tau \frac{\partial^2 l_2(\beta^*, \gamma^*)}{\partial(\beta, \gamma)^2} u \\ &= n^{-\frac{2}{3}} \left( n^{\frac{1}{3}} u^\tau \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)} - \frac{1}{2} u^\tau U u + o_p(1) \right) \end{aligned}$$

Denote  $\lambda_{min}$  is the smallest eigenvalue of  $U$ . Since  $U$  is positive definite,  $\lambda_{min} > 0$ . Then

$$\begin{aligned} P(\|n^{\frac{1}{3}} u^\tau \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)}\| \leq \frac{\lambda_{min}}{4} \|u\|) &= P(\|u^\tau \sqrt{n} \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)}\| \leq \frac{\lambda_{min}}{4} \|u\| n^{\frac{1}{6}}) \\ &\geq P(\|\sqrt{n} \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)}\| \leq \frac{\lambda_{min}}{4} n^{\frac{1}{6}}) \\ &\rightarrow 1 \end{aligned}$$

where the last convergence holds since  $\sqrt{n} \frac{\partial l_2(\beta_0, \gamma_0)}{\partial(\beta, \gamma)}$  is asymptotical normal and  $\frac{\lambda_{min}}{4} n^{\frac{1}{6}} \rightarrow \infty$ . Since  $\frac{1}{2} u^\tau U u - \frac{\lambda_{min}}{4} \|u\| \geq \frac{\lambda_{min}}{4} \|u\| \geq 0$  and the last equation holds if and only if  $\|u\|=0$ , we have

$$l_2(\theta) \leq l_2(\theta_0) \text{ in probability if } \theta \in B_n, \quad (1.33)$$

and the equation holds if and only if  $\|u\|=0$ .

For  $l_3(\pi) = l_3(\theta)$ ,

$$\begin{aligned}
\frac{\partial l_3(\pi_0)}{\partial \pi_j} &= \frac{a_j}{\pi_{j0}} - \sum_{i=1}^r w_i \frac{1}{1 - \sum_{j=1}^s \pi_{j0}} \\
&= \sum_{i=r+1}^n w_i \frac{I\{Z_i = z_j\}}{\pi_{j0}} - \sum_{i=1}^r w_i \frac{1}{1 - \sum_{j=1}^s \pi_{j0}} \\
&\rightarrow_p E\left((1-\delta) \frac{I\{Z = z_j\}}{\pi_{j0}}\right) - E\left(\delta \frac{1}{1 - \sum_{j=1}^s \pi_{j0}}\right) \\
&= \frac{P(\delta = 0, Z = z_j)}{\pi_{j0}} - \frac{P(\delta = 1)}{1 - \sum_{j=1}^s P(\delta = 0, Z = z_j)} \\
&= 1 - 1 \\
&= 0
\end{aligned}$$

By similar calculation we can show that

$$\frac{\partial^2 l_3(\pi^*)}{\partial \pi^2} \rightarrow_p -\text{diag} \left\{ \frac{1}{\pi_{10}}, \dots, \frac{1}{\pi_{s0}} \right\} - \frac{1}{1 - \sum_{j=1}^s \pi_{j0}} J J^T$$

where  $\pi^*$  is between  $\pi$  and  $\pi_0$  and  $J$  is a column vector of 1 with length  $s$ . By central limit theorem and delta method, we can show that  $\sqrt{n} \frac{\partial l_3(\pi_0)}{\partial \pi}$  is asymptotically normal. If we denote  $\|\pi - \pi_0\| = n^{-\frac{1}{2}} v$ , then by similar arguments for  $l_2(\theta)$ , we can show that

$$l_3(\theta) \leq l_3(\theta_0) \text{ in probability if } \theta \in B_n, \quad (1.34)$$

and the equation holds if and only if  $\|v\|=0$ .

Therefore, by (1.31), (1.33) and (1.34), we show that, in the set  $B_n$ ,

$$l(\theta, \nu(\theta)) < l(\theta_0, \nu(\theta_0)) \text{ in probability.}$$

Because  $l(\theta, \nu(\theta))$  is continuous and differentiable, it must attain local maximum at some point  $\hat{\theta}$  inside the ball with surface  $B_n$  and  $\hat{\theta}$  and  $\hat{\nu} = \nu(\hat{\theta})$  satisfy

$$Q_{1n}(\hat{\theta}, \hat{\nu}) = 0, \quad Q_{2n}(\hat{\theta}, \hat{\nu}) = 0, \quad (1.35)$$

where

$$\begin{aligned}
Q_{1n}(\theta, \nu) &= \sum_{i=1}^r w_i \frac{g(Y_i, \theta)}{1 + \nu^T g(Y_i, \theta)}, \\
Q_{2n}(\theta, \nu) &= \sum_{i=1}^r w_i \frac{(\partial g(Y_i, \theta) / \partial \theta)^T}{1 + \nu^T g(Y_i, \theta)} \nu - \frac{\partial l_2(\theta)}{\partial \theta} - \frac{\partial l_3(\theta)}{\partial \theta}.
\end{aligned}$$

Notice that (1.35) is equivalent to that  $(\hat{\theta}, \hat{\nu})$  is the solution to (1.8)-(1.11). This proves (1.15).

The consistency of  $(\hat{\theta}, \hat{\nu})$  follows from the fact that  $B_n$  shrinks to  $\theta_0$  as  $n \rightarrow \infty$ .

Expanding  $Q_{1n}(\hat{\theta}, \hat{\nu})$ ,  $Q_{2n}(\hat{\theta}, \hat{\nu})$  at  $(\theta_0, 0)$ , we have

$$\begin{aligned} 0 &= Q_{1n}(\hat{\theta}, \hat{\nu}) = Q_{1n}(\theta_0, 0) + \frac{\partial Q_{1n}(\theta_0, 0)}{\partial \theta}(\hat{\theta} - \theta_0) + \frac{\partial Q_{1n}(\theta_0, 0)}{\partial \nu^\tau}(\hat{\nu} - 0) + o_p(\Delta_n), \\ 0 &= Q_{2n}(\hat{\theta}, \hat{\nu}) = Q_{2n}(\theta_0, 0) + \frac{\partial Q_{2n}(\theta_0, 0)}{\partial \theta}(\hat{\theta} - \theta_0) + \frac{\partial Q_{2n}(\theta_0, 0)}{\partial \nu^\tau}(\hat{\nu} - 0) + o_p(\Delta_n), \end{aligned}$$

where  $\Delta_n = \|\hat{\theta} - \theta_0\| + \|\hat{\nu}\|$ . Then

$$\begin{pmatrix} \hat{\nu} \\ \hat{\theta} - \theta_0 \end{pmatrix} = S_n^{-1} \begin{pmatrix} -Q_{1n}(\theta_0, 0) + o_p(\Delta_n) \\ -Q_{2n}(\theta_0, 0) + o_p(\Delta_n) \end{pmatrix}, \quad (1.36)$$

where

$$\begin{aligned} S_n &= \begin{pmatrix} \frac{\partial Q_{1n}}{\partial \nu^\tau} & \frac{\partial Q_{1n}}{\partial \theta} \\ \frac{\partial Q_{2n}}{\partial \nu^\tau} & \frac{\partial Q_{2n}}{\partial \theta} \end{pmatrix} \rightarrow_p S = \begin{pmatrix} S_{11} & S_{12} \\ S_{12}^\tau & S_{22} \end{pmatrix}, \\ S_{11} &= -(1 - \sum_{j=1}^s \pi_{j0})E(gg^\tau | \delta = 1), \\ S_{12} &= (1 - \sum_{j=1}^s \pi_{j0})E\left(\frac{\partial g}{\partial \theta} \middle| \delta = 1\right)^\tau, \\ S_{22} &= \text{diag}\{U, V\}, \\ V &= \text{diag}\left\{\frac{1}{\pi_{10}}, \dots, \frac{1}{\pi_{s0}}\right\} + \frac{1}{1 - \sum_j \pi_{j0}} J J^\tau, \end{aligned} \quad (1.37)$$

$\pi_{j0}$  is the true value of  $\pi_j$ ,  $U$  is defined in (1.32), and  $J$  is a column vector of 1 with length  $s$ . By central limit theorem,

$$\sqrt{n} \begin{pmatrix} Q_{1n}(\theta_0, 0) \\ Q_{2n}(\theta_0, 0) \end{pmatrix} \rightarrow_d N(0, T) = N\left(0, \begin{pmatrix} -T_{11} & 0 \\ 0 & T_{22} \end{pmatrix}\right), \quad (1.38)$$

where

$$T_{11} = dS_{11}, \quad T_{22} = \text{diag}\{dU, dV\},$$

and  $d$  is defined in condition (ii). Then by (1.36), (1.37) and (1.38), we conclude that (1.16) holds with  $\Sigma = S^{-1}TS$ .

Let  $k(\theta, \nu) = \sum_{i=1}^r p_i Y_i$ . Then by Taylor expansion,

$$\hat{Y} = k(\hat{\theta}, \hat{\nu}) = k(\theta_0, 0) + \frac{\partial k(\theta^*, \nu^*)}{\partial \nu}(\hat{\nu} - 0) + \frac{\partial k(\theta^*, \nu^*)}{\partial \theta}(\hat{\theta} - \theta_0), \quad (1.39)$$

where  $(\theta^*, \nu^*)$  is between  $(\hat{\theta}, \hat{\nu})$  and  $(\theta_0, 0)$ . By the convergence of  $(\hat{\theta}, \hat{\nu})$ , we can show that  $\left(\frac{\partial k(\theta^*, \nu^*)}{\partial \nu}, \frac{\partial k(\theta^*, \nu^*)}{\partial \theta}\right)$  is consistent for a constant vector  $c$ . Then by (1.36), (1.37), (1.38) and

(1.39), we have

$$\begin{aligned}\sqrt{n}(\hat{Y} - \bar{Y}) &= \sqrt{n}(\hat{Y} - EY) + o_p(1) \\ &= \sqrt{n} \left( k(\theta_0, 0) - c^\tau S^{-1} \begin{pmatrix} Q_{1n}(\theta_0, 0) \\ Q_{2n}(\theta_0, 0) \end{pmatrix} - EY \right) + o_p(1)\end{aligned}\quad (1.40)$$

$$= \sqrt{n}t \left( \sum_{i=1}^n w_i \phi(x_i, \theta_0) \right) + o_p(1), \quad (1.41)$$

where  $x_i = (\delta_i, Y_i, Z_i)$ ,

$$\begin{aligned}\phi(x_i, \theta_0) &= \left\{ \frac{\delta_i Y_i (1 - \sum_j \pi_{j0})}{\sum_j \phi_{ij} f_{ij}}, \delta_i, \delta_i g(Y_i, \theta_0), \delta_i \left( \frac{\partial \log f_i}{\partial \beta} - \frac{\sum_j \phi_{ij} \frac{\partial f_{ij}}{\partial \beta}}{\sum_j \phi_{ij} f_{ij}} \right), \right. \\ &\left. \delta_i \left( \frac{\partial \log \phi_i}{\partial \gamma} - \frac{\sum_j \frac{\partial \phi_{ij}}{\partial \gamma} f_{ij}}{\sum_j \phi_{ij} f_{ij}} \right), \frac{(1 - \delta_i) I\{Z_i = z_j\}}{\pi_{j0}} - \frac{\delta_i}{1 - \sum_j \pi_{j0}} \Big|_{j=1, \dots, s} \right\},\end{aligned}$$

$f_i = f(Y_i, Z_i, \beta_0)$ ,  $\phi_i = \phi(Y_i, Z_i, \gamma_0)$ ,  $f_{ij} = f(Y_i, z_j, \beta_0)$ ,  $\phi_{ij} = \phi(Y_i, z_j, \gamma_0)$ , and function  $t$  is defined as

$$t(\zeta, \eta, \kappa, \xi, \varsigma, \varrho_1, \dots, \varrho_s) = \frac{\zeta}{\eta} - c^\tau S^{-1} \{\kappa, -\xi, -\varsigma, -\varrho_1, \dots, -\varrho_s\}^\tau - EY$$

with  $\kappa$  being  $s$ -dimensional,  $\xi$  being  $p$ -dimensional,  $\varsigma$  being  $q$ -dimensional, and  $\zeta, \eta, \varrho_1, \dots, \varrho_s$  being real numbers. Denote  $\bar{\phi} = \sum_{i=1}^n w_i \phi(x_i, \theta_0)$  and  $E\phi = E(\bar{\phi})$ . By central limit theorem and  $\delta$ -method,

$$\sqrt{n} [t(\bar{\phi}) - t(E\phi)] \rightarrow_d N(0, \sigma^2), \quad (1.42)$$

where

$$\sigma^2 = [t'(E\phi)] [dE\phi\phi^\tau - E\phi E\phi^\tau] [t'(E\phi)]^\tau.$$

Then by (1.41) and (1.42), we have

$$\sqrt{n} (\hat{Y} - \bar{Y} - t(E\phi)) \rightarrow_d N(0, \sigma^2). \quad (1.43)$$

On the other hand,  $\hat{Y} - \bar{Y} \rightarrow_p 0$  by (1.40) and the fact that  $k(\theta_0, 0) \rightarrow_p EY$ . Then it follows (1.43) that  $t(E\phi) = 0$  and  $\sqrt{n}(\hat{Y} - \bar{Y}) \rightarrow_d N(0, \sigma^2)$ . The proof of  $\sqrt{n}(\hat{Y}_j - \bar{Y}_j) \rightarrow_d N(0, \sigma_j^2)$  is similar. This shows (1.17) and completes the proof of Theorem 1.

**PROOF OF THEOREM 2.** Notice that  $l(\beta, \gamma, \tilde{\pi}, \tilde{\lambda}) = l_2(\theta)$ . By the proof of Theorem 1, (1.21) holds with  $\Sigma_p = dU^{-1}$ . (1.22) can be shown similarly to the proof of (1.17).

**PROOF OF THEOREM 3.** The proof is similar to the proof of Theorems 1 and 2, but we replace the functions and the parameters with their bootstrap analog. First of all, in Lemma 1, if we

denote  $\{x_1^*, \dots, x_n^*\}$  as a bootstrap sample, since

$$E(g(x^*, \theta_0)) = E\left(\frac{1}{n} \sum_{i=1}^n g(x_i, \theta_0)\right) = E(g(x, \theta_0)) = 0,$$

we know that when  $\|\theta - \theta_0\| = O_p(n^{-1/3})$ , as  $n \rightarrow \infty$ , 0 is contained in the convex hull of  $\{g(x_i^*, \theta), i = 1, \dots, n\}$  with probability 1. The bootstrap analog of Lemma 1 follows. Then, similar to the proof of Theorem 1, we can show (1.23) and  $(\hat{\theta}^*, \hat{\nu}^*)$  satisfies

$$Q_{1n}^*(\hat{\theta}^*, \hat{\nu}^*) = 0, \quad Q_{2n}^*(\hat{\theta}^*, \hat{\nu}^*) = 0,$$

where  $Q_{1n}^*$  and  $Q_{2n}^*$  are the bootstrap analog of  $Q_{1n}$  and  $Q_{2n}$ . Then

$$\begin{pmatrix} \hat{\nu}^* - \hat{\nu} \\ \hat{\theta}^* - \hat{\theta} \end{pmatrix} = S_n^{*-1} \begin{pmatrix} -Q_{1n}^*(\hat{\theta}, \hat{\nu}) + o_p(\Delta_n^*) \\ -Q_{2n}^*(\hat{\theta}, \hat{\nu}) + o_p(\Delta_n^*) \end{pmatrix}, \quad (1.44)$$

where  $\Delta_n^* = \|\hat{\theta}^* - \hat{\theta}\| + \|\hat{\nu}^* - \hat{\nu}\|$  and

$$S_n^* = \begin{pmatrix} \frac{\partial Q_{1n}^*(\hat{\theta}, \hat{\nu})}{\partial \nu^\tau} & \frac{\partial Q_{1n}^*(\hat{\theta}, \hat{\nu})}{\partial \theta} \\ \frac{\partial Q_{2n}^*(\hat{\theta}, \hat{\nu})}{\partial \nu^\tau} & \frac{\partial Q_{2n}^*(\hat{\theta}, \hat{\nu})}{\partial \theta} \end{pmatrix}. \quad (1.45)$$

By Lemma 1 of Fang, Hong, and Shao (2008),  $S_n^* \rightarrow_p S$ , where  $S$  is given in (1.37), and

$$\sqrt{n} \begin{pmatrix} Q_{1n}^*(\hat{\theta}, \hat{\nu}) - Q_{1n}(\hat{\theta}, \hat{\nu}) \\ Q_{2n}^*(\hat{\theta}, \hat{\nu}) - Q_{2n}(\hat{\theta}, \hat{\nu}) \end{pmatrix} \rightarrow_{d^*} N(0, T), \quad (1.46)$$

where  $T$  is given in (1.38). Notice that  $Q_{1n}(\hat{\theta}, \hat{\nu}) = 0$  and  $Q_{2n}(\hat{\theta}, \hat{\nu}) = 0$ . Then by (1.44), (1.45) and (1.46), we show (1.24). The proofs of (1.25), (1.26) and (1.27) are similar.

**PROOF OF THEOREM 4.** The proofs for the mean imputation estimators are similar to that of Theorem 1. Conditional on the sample, the mean of the random imputation estimators are equal to the mean imputation estimators. Then the results for the random imputation estimators follow from those for the mean imputation estimators and Lemma 1 of Schenker and Welsh (1988).



Table 1: For  $\gamma = 1.8$ : Relative Bias (RB) in % and Variance (VAR) of the Estimators, Bootstrap Variance Estimates (Vboot), Coverage Probability (CP) in %, and Length (LEN) of 95% Confidence Interval

Method		Naive					MELE					MPELE				
		RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN
Without																
Imputation	Y	5.92	.0026	.0025	0	.19	.23	.0042	.0041	91.5	.24	.16	.0446	.0522	96.5	.84
	Y <sub>1</sub>	1.95	.0175	.0160	74.0	.49	.26	.0043	.0044	93.5	.25	.19	.0362	.0414	96.5	.75
	Y <sub>2</sub>	3.45	.0112	.0127	30.0	.44	.26	.0043	.0044	92.8	.25	.17	.0362	.0414	96.9	.75
	Y <sub>3</sub>	5.75	.0106	.0107	0	.40	.24	.0043	.0044	91.8	.25	.18	.0362	.0414	96.5	.75
	Y <sub>4</sub>	8.63	.0105	.0097	0	.38	.23	.0043	.0044	90.8	.25	.20	.0362	.0414	96.9	.75
	Y <sub>5</sub>	10.59	.0151	.0141	0	.46	.16	.0146	.0150	94.1	.47	.41	.1208	.1333	94.2	1.37
Mean																
Imputation	Y	6.74	.0026	.0025	0	.19	.15	.0031	.0034	91.2	.22	.15	.0055	.0059	95.8	.29
	Y <sub>1</sub>	2.30	.0173	.0157	66.4	.48	.10	.0151	.0165	96.2	.50	.22	.0173	.0162	95.8	.50
	Y <sub>2</sub>	3.62	.0137	.0126	26.8	.43	.10	.0138	.0127	94.4	.43	.18	.0127	.0128	96.2	.44
	Y <sub>3</sub>	5.75	.0105	.0106	0	.40	.31	.0101	.0102	93.2	.39	.23	.0122	.0120	91.9	.42
	Y <sub>4</sub>	8.45	.0103	.0097	0	.38	.13	.0070	.0079	91.5	.34	.04	.0141	.0143	95.8	.45
	Y <sub>5</sub>	10.31	.0149	.0145	0	.46	.09	.0131	.0145	93.6	.47	.14	.0183	.0236	96.2	.58
Random																
Imputation	Y	6.72	.0032	.0030	0	.21	.13	.0033	.0037	91.0	.23	.12	.0058	.0062	95.4	.30
	Y <sub>1</sub>	2.24	.0194	.0169	70.8	.50	.08	.0154	.0170	95.7	.50	.19	.0176	.0166	95.8	.50
	Y <sub>2</sub>	3.64	.0124	.0142	30.8	.46	.10	.0147	.0134	94.0	.45	.20	.0138	.0135	96.9	.45
	Y <sub>3</sub>	5.70	.0124	.0125	1.2	.43	.29	.0109	.0111	91.5	.41	.16	.0135	.0130	94.6	.44
	Y <sub>4</sub>	8.47	.0117	.0118	0	.42	.12	.0081	.0092	91.5	.37	.00	.0155	.0156	96.2	.48
	Y <sub>5</sub>	10.26	.0185	.0175	0	.51	.04	.0163	.0168	94.0	.50	.10	.0213	.0258	96.5	.61

Table 2: For  $\gamma = 2$ : Relative Bias (RB) in % and Variance (VAR) of the Estimators, Bootstrap Variance Estimates (Vboot), Coverage Probability (CP) in %, and Length (LEN) of 95% Confidence Interval

Method		Naive					MELE					MPELE				
		RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN
Without																
Imputation	Y	3.61	.0021	.0020	0	.17	.18	.0026	.0026	94.8	.20	.13	.0243	.0323	94.6	.67
	Y <sub>1</sub>	.81	.0158	.0162	92.4	.49	.17	.0032	.0032	90.4	.22	.03	.0196	.0232	96.8	.57
	Y <sub>2</sub>	1.64	.0120	.0120	80.4	.42	.18	.0032	.0032	94.8	.22	.01	.0196	.0232	93.5	.57
	Y <sub>3</sub>	3.19	.0105	.0096	23.6	.38	.17	.0032	.0032	93.6	.22	.02	.0196	.0232	94.2	.57
	Y <sub>4</sub>	5.09	.0079	.0078	.4	.34	.20	.0032	.0032	94.8	.22	.04	.0196	.0232	94.6	.57
	Y <sub>5</sub>	6.58	.0108	.0103	0	.39	.15	.0118	.0116	95.6	.42	.26	.0752	.0915	94.2	1.14
Mean																
Imputation	Y	3.96	.0022	.0021	0	.17	.14	.0026	.0024	94.4	.19	.06	.0035	.0035	96.2	.23
	Y <sub>1</sub>	1.11	.0153	.0161	88.8	.49	.12	.0171	.0168	94.8	.50	.18	.0169	.0165	94.6	.50
	Y <sub>2</sub>	1.79	.0118	.0119	76.8	.42	.17	.0124	.0126	95.2	.43	.05	.0132	.0126	93.5	.44
	Y <sub>3</sub>	3.21	.0105	.0095	22.0	.38	.13	.0097	.0098	94.4	.38	.27	.0108	.0104	95.0	.40
	Y <sub>4</sub>	4.97	.0080	.0078	.4	.34	.15	.0079	.0076	94.8	.34	.03	.0103	.0107	96.2	.40
	Y <sub>5</sub>	6.61	.0106	.0104	0	.39	.12	.0119	.0116	95.6	.42	.15	.0153	.0162	96.5	.49
Random																
Imputation	Y	3.94	.0026	.0023	0	.18	.12	.0027	.0025	94.4	.19	.08	.0036	.0037	95.8	.23
	Y <sub>1</sub>	1.15	.0159	.0167	87.6	.50	.12	.0176	.0169	95.2	.50	.18	.0172	.0167	94.6	.50
	Y <sub>2</sub>	1.78	.0126	.0128	79.6	.44	.18	.0129	.0129	95.2	.44	.08	.0136	.0129	94.6	.44
	Y <sub>3</sub>	3.20	.0116	.0106	26.8	.40	.12	.0105	.0103	94.8	.39	.26	.0110	.0108	94.2	.41
	Y <sub>4</sub>	4.92	.0096	.0091	.4	.37	.15	.0089	.0083	94.8	.35	.07	.0109	.0112	96.2	.41
	Y <sub>5</sub>	6.58	.0133	.0123	0	.43	.06	.0135	.0127	96.0	.44	.15	.0159	.0172	95.0	.50

Table 3: For  $\gamma = -1.4$ : Relative Bias (RB) in % and Variance (VAR) of the Estimators, Bootstrap Variance Estimates (Vboot), Coverage Probability (CP) in %, and Length (LEN) of 95% Confidence Interval

Method		Naive					MELE					MPELE				
		RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN	RB	VAR	Vboot	CP	LEN
Without																
Imputation	Y	-3.96	.0018	.0019	0	.17	.20	.0031	.0031	94.0	.21	.57	.1084	.1195	95.4	1.18
	Y <sub>1</sub>	-9.89	.0201	.0204	0	.55	.20	.0043	.0039	90.0	.24	.58	.1224	.1369	96.2	1.28
	Y <sub>2</sub>	-7.11	.0135	.0131	0	.44	.25	.0043	.0039	92.5	.24	.63	.1224	.1369	95.4	1.28
	Y <sub>3</sub>	-4.84	.0085	.0090	2.0	.37	.19	.0043	.0039	91.5	.24	.58	.1224	.1369	95.4	1.28
	Y <sub>4</sub>	-2.86	.0065	.0068	19.2	.32	.22	.0043	.0039	94.5	.24	.60	.1224	.1369	95.8	1.28
	Y <sub>5</sub>	-1.72	.0085	.0086	66.8	.36	.15	.0112	.0106	94.0	.40	.30	.0823	.0776	96.9	.96
Mean																
Imputation	Y	-4.57	.0018	.0019	0	.17	.25	.0027	.0028	92.7	.20	.15	.0054	.0059	94.2	.29
	Y <sub>1</sub>	-9.46	.0211	.0208	0.4	.55	.26	.0124	.0128	95.4	.44	.08	.0279	.0269	97.3	.60
	Y <sub>2</sub>	-6.81	.0132	.0130	0	.44	.36	.0112	.0113	94.2	.41	.03	.0180	.0190	97.3	.52
	Y <sub>3</sub>	-4.80	.0086	.0089	2.4	.36	.14	.0096	.0097	93.1	.38	.20	.0130	.0134	96.5	.44
	Y <sub>4</sub>	-2.97	.0067	.0068	17.2	.32	.26	.0070	.0078	94.6	.34	.16	.0098	.0100	95.4	.39
	Y <sub>5</sub>	-1.93	.0085	.0086	60.4	.36	.27	.0102	.0099	94.2	.38	.35	.0102	.0099	93.5	.39
Random																
Imputation	Y	-4.58	.0022	.0022	0	.18	.24	.0031	.0031	93.8	.21	.17	.0059	.0063	95.0	.29
	Y <sub>1</sub>	-9.47	.0025	.0024	0.8	.60	.25	.0167	.0186	96.5	.53	.01	.0337	.0327	97.3	.67
	Y <sub>2</sub>	-6.83	.0152	.0155	0	.48	.30	.0132	.0142	93.5	.46	.08	.0218	.0222	96.9	.56
	Y <sub>3</sub>	-4.80	.0105	.0103	3.6	.39	.18	.0118	.0112	93.5	.41	.18	.0147	.0149	95.8	.47
	Y <sub>4</sub>	-2.95	.0076	.0075	20.0	.33	.23	.0086	.0084	95.0	.35	.15	.0106	.0106	94.2	.41
	Y <sub>5</sub>	-1.94	.0091	.0092	60.8	.37	.25	.0104	.0103	95.0	.39	.34	.0109	.0104	92.7	.40

Table 4: The Mean(Mean) and the Variance (VAR) of the Parameter Estimates. The true values are  $\beta_1 = 0.25$ ,  $\beta_2 = 0.5$ ,  $\beta_3 = 0.75$ ,  $\beta_4 = 1$ , and  $\beta_5 = -0.1$ .

	MELE						MPELE					
	$\gamma = 1.8$		$\gamma = 2$		$\gamma = -1.4$		$\gamma = 1.8$		$\gamma = 2$		$\gamma = -1.4$	
	Mean	VAR	Mean	VAR	Mean	VAR	Mean	VAR	Mean	VAR	Mean	VAR
$\beta_1$	.2554	.4074	.2425	.2987	.2986	.3863	0.3935	2.1449	.1479	1.5564	.2592	.5873
$\beta_2$	.5173	.3873	.5216	.3025	.5644	.3883	0.6505	2.0736	.4004	1.4972	.4986	.6183
$\beta_3$	.7601	.3990	.7667	.2984	.8052	.3792	0.9012	1.9614	.6580	1.4382	.7515	.6906
$\beta_4$	1.0171	.3878	1.0136	.2945	1.0621	.3713	1.1373	1.8363	.9075	1.3843	.9882	.7577
$\beta_5$	-.1006	.0057	-.1015	.0041	-.1075	.0054	-.1164	.0214	-.0896	.0167	-.0994	.0127
$\gamma$	1.7952	.0004	1.9960	.0005	-1.4038	.0002	1.8236	.0597	2.0350	.1310	-1.3912	.0163

Table 5: The Ratio of MSE:  $mse(MPELE)/mse(MELE)$ .

	Without Imputation						Mean Imputation						Random Imputation					
	Y	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>	Y <sub>5</sub>	Y	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>	Y <sub>5</sub>	Y	Y <sub>1</sub>	Y <sub>2</sub>	Y <sub>3</sub>	Y <sub>4</sub>	Y <sub>5</sub>
$\gamma = 1.8$	9.85	7.38	7.41	7.51	7.54	7.78	1.23	0.97	0.84	1.11	1.39	1.27	1.16	0.95	0.84	1.06	1.28	1.15
$\gamma = 2$	8.01	5.38	5.34	5.35	5.24	7.39	1.01	0.85	1.05	1.21	1.16	1.17	1.00	0.85	1.04	1.15	1.09	1.28
$\gamma = -1.4$	23.91	24.90	24.13	25.16	24.74	4.34	1.46	1.69	1.04	1.13	1.33	1.00	1.52	1.47	1.03	1.04	1.40	1.04