

Pseudo-likelihood and bootstrapped pseudo-likelihood
inference in logistic regression model with
misclassified responses

A. CHATTERJEE, T. BANDYOPADHYAY AND S. ADHYA

Indian Statistical Institute, Delhi Centre
7, SJSS Marg, New Delhi-110 016, India

Pseudo-likelihood and bootstrapped pseudo-likelihood inference in logistic regression model with misclassified responses

A. Chatterjee*

Theoretical Statistics & Mathematics Unit
Indian Statistical Institute, Delhi

and

T. Bandyopadhyay†

Production & Quantitative Methods Group
Indian Institute of Management, Ahmedabad

and

S. Adhya‡

Department of Statistics
West Bengal State University, Kolkata

Abstract

Logistic regression is an extensively used regression model for binary responses. In many applications, misclassification of binary responses is not uncommon. If the misclassification is ignored, it may severely bias the maximum likelihood estimators (MLE) of regression parameters towards zero. To obviate this difficulty, we propose a pseudo-likelihood method of estimation, that uses data from internal validation study. Under minimal assumptions, we establish rigorous asymptotic results for the maximum pseudo-likelihood estimators. A bootstrapped version of the maximum pseudo likelihood estimators is proposed, and its distributional consistency is proved. It enables us to use bootstrap method for statistical inference. The results of the simulation studies clearly indicate the superiority of the maximum pseudo-likelihood estimators to the maximum full likelihood estimators, and the maximum likelihood estimators based on misclassified binary responses only. Also, inferences on the regression parameters using asymptotic distribution of maximum pseudo-likelihood estimators, and its bootstrap version, are found to be similar.

Keywords: Logistic regression, Response misclassification, Internal validation, Pseudo-likelihood, Bootstrap consistency.

*cha@isid.ac.in

†tathagata@iima.ac.in

‡sumanta.adhya@gmail.com

1 Introduction

Logistic regression is an important and widely used regression model for binary responses, and has found its use in a variety of applied fields (cf. [Hilbe \(2009\)](#), [Hosmer and Lemeshow \(2004\)](#)), especially in epidemiological research, including medical and social sciences (cf. [Jewell \(2003\)](#)). In many applications, it is not uncommon, that the binary responses are subject to classification errors. Misclassified binary responses occur due to various reasons, e.g., faulty data collected through surveys ([Hausman et al. \(1998\)](#), [Hausman \(2001\)](#), [Bollinger and David \(1997\)](#), [Savoca \(2011\)](#)), limited sensitivity and specificity of the diagnostic tests ([Lyles et al. \(2011\)](#), [Edwards et al. \(2013\)](#), [Gilbert et al. \(2014\)](#)), incorrect information gathered from medical and other records ([Spoto et al. \(1992\)](#)) and recall bias in assessing exposure status ([Gordis \(2009\)](#)). For pervasiveness of misclassification of important binary outcomes (*viz.*, program receipt, labor market status, educational attainment, self reported health conditions, physical and mental impairment etc.) in survey data, and its adverse impact on the estimates of the logistic regression parameters we refer to [Meyer and Mittag \(2016\)](#) and [Savoca \(2011\)](#), and the references therein. Suppose, $Y \in \{0, 1\}$ denotes the true response variable and $\tilde{Y} \in \{0, 1\}$ denotes the misclassified version of Y . Then, Y is said to be misclassified, if we observe $\tilde{Y} = 1$, when the corresponding $Y = 0$, or vice-versa. In this article, we investigate the inference problems in logistic regression when the binary responses are subject to classification errors.

Measurement errors in regression models have been widely studied, and excellent text books (cf. [Fuller \(2006\)](#), [Carroll et al. \(2006\)](#) or [Buonaccorsi \(2010\)](#)) are written on it. The effects of measurement error on covariates has been well investigated for simple logistic regression (cf. [Carroll et al. \(1984\)](#), [Stefanski and Carroll \(1985\)](#)), and also for semiparametric logistic regression (cf. [Carroll and Wand \(1991\)](#), [Wang and Wang \(1997\)](#)). In contrast, the study of the effect of misclassified responses in logistic regression has received lesser attention in statistics literature. [Neuhaus \(1999\)](#) studied the bias and efficiency loss due to misclassified responses in binary regression. [Carroll et al. \(2006\)](#) (Section 15.3) discuss misclassification of binary responses in logistic regression, and show that, not accounting for misclassification introduces severe bias in parameter estimates. Throughout, we assume that the information on the underlying covariates is available without any measurement error.

In particular, we assume that the true binary response variable Y is associated with the p -dimensional covariate vector \mathbf{X} by the logistic regression model,

$$\mathbf{P}_0(Y = 1|\mathbf{X} = \mathbf{x}) = \frac{1}{1 + \exp\{-\mathbf{x}'\boldsymbol{\beta}_0\}} \equiv \psi(\mathbf{x}'\boldsymbol{\beta}_0), \quad \text{for all } \mathbf{x} \in \mathbb{R}^p, \quad (1.1)$$

where, $\boldsymbol{\beta}_0 = (\beta_{1,0}, \dots, \beta_{p,0})' \in \mathbb{R}^p$, is an unknown regression parameter, and \mathbf{P}_0 denotes the underlying probability distribution corresponding to $\boldsymbol{\beta}_0$. The binary variable $\tilde{Y} \in \{0, 1\}$, is the misclassified version of the true response Y . We assume the following model for misclassification ([Ekholm and Palmgren \(1987\)](#), [Copas \(1988\)](#), [Magder and Hughes \(1997\)](#), [Bollinger and David \(1997\)](#), [Hausman et al. \(1998\)](#), [Hausman \(2001\)](#), [Roy et al. \(2005\)](#), [Lyles et al. \(2011\)](#), [Gilbert et al. \(2014\)](#), [Savoca \(2011\)](#) and [Meyer and Mittag](#)

(2016)),

$$\mathbf{P}_0(\tilde{Y} = 1|Y = 0) = \theta_{1,0} \quad \text{and} \quad \mathbf{P}_0(\tilde{Y} = 0|Y = 1) = \theta_{2,0}, \quad (1.2)$$

which results in the following regression model for \tilde{Y} ,

$$\mathbf{P}_0(\tilde{Y} = 1|\mathbf{X} = \mathbf{x}) = \theta_{1,0} + (1 - \theta_{1,0} - \theta_{2,0}) \cdot \psi(\mathbf{x}'\boldsymbol{\beta}_0). \quad (1.3)$$

Note $\boldsymbol{\theta}_0 = (\theta_{1,0}, \theta_{2,0})'$ are the unknown misclassification probabilities and \mathbf{P}_0 denotes the probability distribution under $\boldsymbol{\theta}_0$. This misclassification model assumes that, conditional on Y , the data generation process of \tilde{Y} is independent of the covariate \mathbf{X} . The model, as stated above, is simple, but is not unreasonable as a possible description of the actual data generation process. At the end, however, we extend our results to the situations, where $\boldsymbol{\theta}_0$ may depend on \mathbf{x} .

Typically, estimation of the unknown misclassification probabilities $\boldsymbol{\theta}_0$ (cf. (1.2)) requires information on the pair (Y, \tilde{Y}) . Usually, one obtains complete information on $(Y, \tilde{Y}, \mathbf{X})$ for a sub-sample of the main sample known as validation sample. The validation sample is utilized to obtain an improved estimate of $\boldsymbol{\theta}_0$. The rest of the sample has information available only on (\tilde{Y}, \mathbf{X}) , and is known as the non-validation part of the main sample. Use of validation data for adjusting estimates of regression parameters due to misclassification of binary responses is common in many areas of applications. For example, [Lyles et al. \(2011\)](#), [Edwards et al. \(2013\)](#) and [Duffy et al. \(2004\)](#) provide examples in epidemiological studies. [Bollinger and David \(1997\)](#) and [Meyer and Mittag \(2016\)](#) describe such situations in the context of economic and social surveys.

In the applications mentioned above, though the estimation of regression parameter $\boldsymbol{\beta}_0$ is of primary interest, the estimation of parameter $\boldsymbol{\theta}_0$ is also considered to be important as it provides useful information about the level of contamination in the responses. For estimation of $\boldsymbol{\beta}_0$ and $\boldsymbol{\theta}_0$, a possible approach would be to use the full likelihood based on the combined sample comprising validation and non-validation parts of the main sample. Notice that $\boldsymbol{\theta}_0$ shows up in the non-validation part of the full-likelihood due to the use of misclassified responses. The full likelihood, however, as a function of both $\boldsymbol{\beta}_0$ and $\boldsymbol{\theta}_0$ is found to be ill behaved, in the sense that it may often lead to nonsensical estimates of $\boldsymbol{\theta}_0$ even for substantially large sample sizes. In Section 5.1 we will discuss this issue in detail. An alternative approach is to use pseudo-likelihood method proposed by [Gong and Samaniego \(1981\)](#). Recently pseudo-likelihood method has been used in various contexts (cf. [Wang and Zhao \(2007\)](#), [Chen and Liang \(2010\)](#), [Guolo \(2011\)](#), [Ghosh et al. \(2013\)](#), [Lyles and Kupper \(2013\)](#)) especially when, the full likelihood function is ill behaved, yet, the pseudo-likelihood function is well-behaved. In our situation, pseudo-likelihood based approach amounts to replacing the unknown nuisance parameter $\boldsymbol{\theta}_0$ in the likelihood by its estimate from the validation sample, and then assuming as if $\boldsymbol{\theta}_0$ is known, the likelihood is considered as a function of $\boldsymbol{\beta}_0$ only. The pseudo-likelihood as a function of $\boldsymbol{\beta}_0$ is well behaved even for moderately large sample sizes.

Since the validation and non-validation parts of the sample do not have identical distributions, the asymptotic results of [Gong and Samaniego \(1981\)](#) will not be valid in our set up. Also plugging in an estimate of $\boldsymbol{\theta}_0$, introduces dependence among the validation and the non-validation parts of the pseudo-

likelihood function. We develop rigorous asymptotic results under minimal assumptions, and show that the pseudo likelihood estimator of β_0 is asymptotically normal. As shown later in Section 2.2, the asymptotic covariance of the pseudo-likelihood estimator is a complicated function of the unknown parameters β_0 , θ_0 , and the unknown distribution of the covariates. Consequently, the analytical computation of the covariance is unwieldy. Thus, for easy implementation of the pseudo-likelihood based inference, we develop a bootstrap methodology. We prove the distributional consistency of the bootstrapped pseudo-likelihood estimator of β_0 . Consistency of bootstrapped pseudo-likelihood estimators has not been investigated widely in the literature, except in very specific scenarios (cf. [Aerts and Claeskens \(1999\)](#) and [Kawahara and Shimotsu \(2008\)](#)). Existing approaches for proving consistency of bootstrapped M-estimators (cf. [Wellner and Zhan \(1996\)](#), [Chatterjee and Bose \(2005\)](#) and [Cheng and Huang \(2010\)](#)) can not be directly applied in our set-up, precisely due to the reasons mentioned in the beginning of this paragraph. In order to establish the consistency of the bootstrapped pseudo-likelihood estimator of β_0 , we have developed some new bootstrap stochastic equicontinuity results, which may be of independent interest. This is a new contribution to the existing literature on bootstrap consistency for pseudo-likelihood estimators. We discuss the details in Section 3.

We further extend the asymptotic results to the situations where θ_0 may depend on the covariates \mathbf{x} in a way as discussed by [Bollinger and David \(1997\)](#) and [Meyer and Mittag \(2016\)](#). Also, we discuss the asymptotic results when one of the components of θ_0 is zero.

There may be situations when validation data may not be available. In such situations, [Hausman et al. \(1998\)](#) and [Hausman \(2001\)](#) propose maximum likelihood estimators of θ_0 and β_0 based on the contaminated data (\tilde{Y}, \mathbf{x}) . "As a cautionary note [Roy et al. \(2005\)](#) make the important point for practitioners" ([Savoca \(2011\)](#)) that if the sample have a few \mathbf{x} values for which $\psi(\mathbf{x}'\beta_0)$ lie outside the interval $(0.1, 0.9)$ then estimation solely on the basis of the contaminated data may lead to serious identifiability problem. The logit function $\psi(\mathbf{x}'\beta_0)$ in this interval is well approximated by a suitably chosen linear function (cf. [Cox and Snell \(1989\)](#)), and thus it is evident from (1.3), identifiability of β_0 is then a serious issue unless the sample size is large enough to include enough number of \mathbf{x} values for which $\psi(\mathbf{x}'\beta_0)$ lie outside the interval $(0.1, 0.9)$. We consider this issue in detail in Section 5.1.1.

We conduct simulation studies: (i) to compare the performances of the estimators of θ_0 and β_0 based on three different likelihoods discussed above, *viz.*, full likelihood, pseudo-likelihood and likelihood using contaminated data only, and (ii) to investigate the validity of the asymptotic distribution of the pseudo-likelihood estimators and its bootstrap analogues for moderate to large sample sizes. We report the bias and mean squared error (mse) for study (i), and coverage and expected length of the confidence interval of regression parameter for study (ii). Finally, we illustrate our methodology using a real-life data set.

The rest of the article is organized as follows. In Section 2, we state the problem formulation, and describe the pseudo likelihood based estimation methodology. The main theoretical results on the asymptotic properties of the pseudo-likelihood estimator are given in Section 2.2. Theoretical results on the bootstrapped pseudo-likelihood estimator are given in Section 3. In Section 4, we discuss the extension of our results to

more general situations, where misclassification probabilities may depend on the covariates, and also to the situation, where one kind of classification error is absent. Results of numerical study, and the analysis of a real data set are presented in Section 5 and concluding remarks are provided in Section 6. Proofs of the main results are given in Section 7. In the Appendix (cf. Section 8), we provide detailed proofs of auxiliary lemma's needed for proving the main results, we also provide detailed expressions of some matrices used in proving the theorems and prove an useful result on gradients of strictly concave functions.

2 Pseudo-likelihood estimation

Suppose the complete sample is of size n . The validation sample consists of n_1 independent observations on the triplet $(Y, \tilde{Y}, \mathbf{X})$, and is denoted by $\mathcal{X}_{n_1} = \{(Y_i, \tilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq n_1\}$. The non-validation sample consists of $n_2 = (n - n_1)$ independent observations on the pair (\tilde{Y}, \mathbf{X}) , and is denoted by $\mathcal{X}_{n_2} = \{(\tilde{Y}_i, \mathbf{X}_i) : (n_1 + 1) \leq i \leq n\}$. It is assumed that the validation and non-validation samples are independent of each other. The full-likelihood at $(\boldsymbol{\beta}, \boldsymbol{\theta})$ using the complete sample $\mathcal{X}_n = \mathcal{X}_{n_1} \cup \mathcal{X}_{n_2}$ is based on the joint conditional distribution of $[Y, \tilde{Y} | \mathbf{X}]$ and is given by,

$$L_{1,n}(\boldsymbol{\beta}, \boldsymbol{\theta}) \equiv \prod_{i=1}^{n_1} \left\{ (1 - \theta_2)^{\tilde{Y}_i} \theta_2^{(1-\tilde{Y}_i)} \psi(\mathbf{X}'_i \boldsymbol{\beta}) \right\}^{Y_i} \cdot \left\{ \theta_1^{\tilde{Y}_i} (1 - \theta_1)^{(1-\tilde{Y}_i)} [1 - \psi(\mathbf{X}'_i \boldsymbol{\beta})] \right\}^{(1-Y_i)} \\ \times \prod_{i=n_1+1}^n \left\{ \theta_1^{\tilde{Y}_i} (1 - \theta_1)^{(1-\tilde{Y}_i)} (1 - \psi(\mathbf{X}'_i \boldsymbol{\beta})) + (1 - \theta_2)^{\tilde{Y}_i} \theta_2^{(1-\tilde{Y}_i)} \psi(\mathbf{X}'_i \boldsymbol{\beta}) \right\}. \quad (2.1)$$

where, $\psi(\cdot)$ denotes the logistic link function given in (1.1). In order to construct the pseudo likelihood function for $\boldsymbol{\beta}$, we need to plug in an estimator of the unknown $\boldsymbol{\theta}_0$ in (2.1). An estimator of $\boldsymbol{\theta}_0$ arises naturally from the misclassification model (1.2). It is based on the observed cell frequencies of the cells $(0, 0), (0, 1), (1, 0), (1, 1)$, the four possible values of (Y, \tilde{Y}) , obtained from the validation sample. However, for small validation sample sizes n_1 , one or more of these cell frequencies may be zero. In such situations, an adjustment of the cell frequencies often improve the performance of the estimator. One commonly used adjustment (cf. Haldane (1956) and Gart and Zweifel (1967)) is to add 1/2 to each cell frequency, and replace the original cell frequencies in the estimator by the adjusted cell frequencies. We thus define the estimator $\hat{\boldsymbol{\theta}}_n = (\hat{\theta}_{1,n}, \hat{\theta}_{2,n})'$, where

$$\hat{\theta}_{1,n} = \frac{\frac{1}{2} + \sum_{i=1}^{n_1} \mathbf{1}(\tilde{Y}_i = 1, Y_i = 0)}{1 + \sum_{i=1}^{n_1} \mathbf{1}(Y_i = 0)} \quad \text{and} \quad \hat{\theta}_{2,n} = \frac{\frac{1}{2} + \sum_{i=1}^{n_1} \mathbf{1}(\tilde{Y}_i = 0, Y_i = 1)}{1 + \sum_{i=1}^{n_1} \mathbf{1}(Y_i = 1)}. \quad (2.2)$$

The small sample properties of the adjusted estimator of the odds ratio in the context of 2×2 contingency tables have been studied in Parzen et al. (2002). If the validation sample size n_1 is large enough, the estimator $\hat{\boldsymbol{\theta}}_n$ in (2.2) is nearly equivalent to the usual cell frequencies based estimator. Substituting $\hat{\boldsymbol{\theta}}_n$ in the full

likelihood function $L_{1,n}(\boldsymbol{\beta}, \boldsymbol{\theta})$ (cf. (2.1)) the scaled pseudo log-likelihood function at $\boldsymbol{\beta}$ becomes,

$$\begin{aligned} l_n(\boldsymbol{\beta}) &= n^{-1} \cdot \log L_{1,n}(\boldsymbol{\beta}, \hat{\boldsymbol{\theta}}_n) \\ &= \frac{1}{n} \cdot \sum_{1 \leq i \leq n_1, Y_i=1} \left[\tilde{Y}_i \log(1 - \hat{\theta}_{2,n}) + (1 - \tilde{Y}_i) \log \hat{\theta}_{2,n} + \log \psi(\mathbf{X}'_i \boldsymbol{\beta}) \right] \\ &\quad + \frac{1}{n} \cdot \sum_{1 \leq i \leq n_1, Y_i=0} \left[\tilde{Y}_i \log \hat{\theta}_{1,n} + (1 - \tilde{Y}_i) \log(1 - \hat{\theta}_{1,n}) + \log \{1 - \psi(\mathbf{X}'_i \boldsymbol{\beta})\} \right] \\ &\quad + \frac{1}{n} \cdot \sum_{i=n_1+1}^n \log \left[(\hat{\theta}_{1,n})^{\tilde{Y}_i} (1 - \hat{\theta}_{1,n})^{(1-\tilde{Y}_i)} \{1 - \psi(\mathbf{X}'_i \boldsymbol{\beta})\} + (1 - \hat{\theta}_{2,n})^{\tilde{Y}_i} (\hat{\theta}_{2,n})^{(1-\tilde{Y}_i)} \psi(\mathbf{X}'_i \boldsymbol{\beta}) \right]. \end{aligned}$$

In order to write down the estimating equation for $\boldsymbol{\beta}$, we define the following functions:

$$\left. \begin{aligned} h_{1,\boldsymbol{\beta}}(y, \mathbf{x}) &= \mathbf{x} \{y - \psi(\mathbf{x}' \boldsymbol{\beta})\}, \\ h_{2,\boldsymbol{\beta},\boldsymbol{\theta}}(\tilde{y}, \mathbf{x}) &= (1 - \theta_1 - \theta_2) \cdot \frac{\mathbf{x} \cdot \psi(\mathbf{x}' \boldsymbol{\beta}) \{1 - \psi(\mathbf{x}' \boldsymbol{\beta})\} \cdot \{\tilde{y} - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}) \{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x})\}}, \quad \text{and} \\ h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}) &= \theta_1 \cdot \{1 - \psi(\mathbf{x}' \boldsymbol{\beta})\} + (1 - \theta_2) \cdot \psi(\mathbf{x}' \boldsymbol{\beta}), \end{aligned} \right\} \quad (2.3)$$

where, $h_{1,\boldsymbol{\beta}} : \{0, 1\} \times \mathbb{R}^p \mapsto \mathbb{R}^p$, $h_{2,\boldsymbol{\beta},\boldsymbol{\theta}} : \{0, 1\} \times \mathbb{R}^p \mapsto \mathbb{R}^p$ and $h_{3,\boldsymbol{\beta},\boldsymbol{\theta}} : \mathbb{R}^p \mapsto (0, 1)$, for all $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$. It should be noted that $h_{1,\boldsymbol{\beta}}(y, \mathbf{x})$ is the estimating function that arises in simple logistic regression and $h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x})$ is same as the conditional expectation of $[\tilde{Y} | \mathbf{X}]$ as in (1.3), if $(\boldsymbol{\beta}_0, \boldsymbol{\theta}_0)$ is replaced with any arbitrary $(\boldsymbol{\beta}, \boldsymbol{\theta})$.

We also define the sequence of validation sample size fractions,

$$f_n = \frac{n_1}{n} = 1 - \frac{n_2}{n}, \quad n \geq 1, \quad (2.4)$$

and the empirical measures,

$$\mathbb{P}_{n_1} \equiv \frac{1}{n_1} \sum_{i=1}^{n_1} \delta_{(Y_i, \mathbf{X}_i)} \quad \text{and} \quad \mathbb{P}_{n_2} \equiv \frac{1}{n_2} \sum_{i=n_1+1}^{n_1+n_2} \delta_{(\tilde{Y}_i, \mathbf{X}_i)}, \quad n_1, n_2 \geq 1, \quad (2.5)$$

where, $\delta_{(y, \mathbf{x})}(\cdot)$ and $\delta_{(\tilde{y}, \mathbf{x})}(\cdot)$ denote point masses (y, \mathbf{x}) and (\tilde{y}, \mathbf{x}) respectively. For any measurable function $f : (y, \mathbf{x}) \mapsto \mathbb{R}^p$, we define $\mathbb{P}_{n_1} f = \int f(y, \mathbf{x}) d\mathbb{P}_{n_1}(y, \mathbf{x})$ and similarly for \mathbb{P}_{n_2} . With these notations, the score function corresponding to the pseudo log-likelihood function $l_n(\boldsymbol{\beta})$ can be written as,

$$Z_n(\boldsymbol{\beta}) = \frac{d}{d\boldsymbol{\beta}} l_n(\boldsymbol{\beta}) = f_n \cdot Z_{n,1}(\boldsymbol{\beta}) + (1 - f_n) \cdot Z_{n,2}(\boldsymbol{\beta}), \quad \text{where,} \quad \left. \begin{aligned} Z_{n,1}(\boldsymbol{\beta}) &= \mathbb{P}_{n_1} h_{1,\boldsymbol{\beta}} \\ Z_{n,2}(\boldsymbol{\beta}) &= \mathbb{P}_{n_2} h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n} \end{aligned} \right\}, \quad (2.6)$$

Note that, $Z_n : \mathbb{R}^p \mapsto \mathbb{R}^p$ is a random function. The pseudo maximum likelihood estimator (PMLE) of $\boldsymbol{\beta}_0$, which we denote by $\hat{\boldsymbol{\beta}}_n$, satisfies the estimating equation

$$Z_n(\boldsymbol{\beta}) = \mathbf{0}, \quad (2.7)$$

where, $\mathbf{0}$ denotes the $p \times 1$ null vector. As seen above, the estimating function $Z_n(\boldsymbol{\beta})$ in (2.6) is a weighted sum of estimating functions $Z_{n,1}(\boldsymbol{\beta})$ and $Z_{n,2}(\boldsymbol{\beta})$, which arise from the validation and non-validation samples respectively. If $f_n = 1$, then the problem reduces to simple logistic regression based estimation of $\boldsymbol{\beta}_0$. The

presence of misclassified responses \tilde{Y}_i 's in the non-validation sample gives rise to the extra term $Z_{n,2}(\boldsymbol{\beta})$. Since $\hat{\boldsymbol{\theta}}_n$ is based on the validation sample and is involved in $Z_{n,2}(\boldsymbol{\beta})$, the terms $Z_{n,1}(\boldsymbol{\beta})$ and $Z_{n,2}(\boldsymbol{\beta})$ are dependent, even though they arise from independent parts of the complete sample. Using (2.6) and (2.7), the pseudo-likelihood estimation problem can be cast into the Z-estimation framework (cf. Chapter 5 of van der Vaart (1998)). Although $Z_{n,1}(\boldsymbol{\beta})$ is the sample mean of the independent and identically distributed (i.i.d.) summands $\{h_{1,\boldsymbol{\beta}}(Y_i, \mathbf{X}_i) : 1 \leq i \leq n_1\}$, and $Z_{n,2}(\boldsymbol{\beta})$ is the sample mean of identical summands $\{h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n}(\tilde{Y}_i, \mathbf{X}_i) : n_1 + 1 \leq i \leq n\}$, their weighted sum $Z_n(\boldsymbol{\beta})$ can not be represented as a sample mean of i.i.d. summands. Hence, the standard asymptotic theory for Z-estimators based on i.i.d. summands will not be directly applicable for studying the asymptotic properties of $\hat{\boldsymbol{\beta}}_n$.

2.1 Theoretical framework for proving asymptotic properties of the pseudo-likelihood estimator

We describe the theoretical framework for proving the main results. Let \mathbf{P}_0 be the generic notation for the true distribution of (Y, \mathbf{X}) or (\tilde{Y}, \mathbf{X}) or $(Y, \tilde{Y}, \mathbf{X})$ under the true value of the parameter $(\boldsymbol{\beta}_0, \boldsymbol{\theta}_0)$. Although the joint distributions of (Y, \mathbf{X}) , (\tilde{Y}, \mathbf{X}) and $(Y, \tilde{Y}, \mathbf{X})$ are different, we use the same notation for simplicity. There does not seem to be any confusion, since the underlying random vector is evident from the notation, and the context.

We use the notation \mathbf{E}_0^* , $o_{\mathbf{P}_0^*}(1)$ and $O_{\mathbf{P}_0^*}(1)$ to denote outer expectation, convergence to zero and bounded in outer probability with respect to the probability measure \mathbf{P}_0 (cf. van der Vaart and Wellner (1996)), respectively. Usual expectation, variance and covariance with respect to \mathbf{P}_0 will be denoted by \mathbf{E}_0 , \mathbf{Var}_0 and \mathbf{cov}_0 , respectively. The symbol \xrightarrow{d} denotes convergence in distribution. Unless stated otherwise, we use the symbols, \mathbf{E} and \mathbf{Var} to denote usual expectation and variance of a random quantity (with respect to the underlying probability distribution).

2.1.1 Technical assumptions

(A1) The true regression parameter $\boldsymbol{\beta}_0 = (\beta_{1,0}, \dots, \beta_{p,0})' \in \mathbb{R}^p$.

(A2) The true misclassification probabilities, $\boldsymbol{\theta}_0 = (\theta_{1,0}, \theta_{2,0})'$ satisfy the following conditions:

- (i) there exist constants, $0 < \delta_1 < \delta_2 < 1$, such that, $0 < \delta_1 < \theta_{1,0}, \theta_{2,0} < \delta_2 < 1$.
- (ii) $\theta_{1,0} + \theta_{2,0} \neq 1$.

The parameter space of $\boldsymbol{\theta}$ satisfying the above restrictions is denoted by Θ .

(A3) Let the marginal distribution of the p -dimensional covariate vector $\mathbf{X} = (X_1, \dots, X_p)'$ be denoted by $Q(\mathbf{x})$. We assume that Q is such that,

$$\mathbf{Var}(\mathbf{X}) = \mathbf{Var}\left((X_1, \dots, X_p)'\right) \text{ exists, and is positive definite.}$$

(A4) The validation sampling fractions $\{f_n : n \geq 1\}$, defined in (2.4), satisfy the following conditions,

- (i) $\lim_{n \rightarrow \infty} f_n = f \in (0, 1)$.
- (ii) $|f_n - f| = o(n^{-1/2})$.

Assumption (A2)(ii) is required to avoid a non-identifiability problem. For a discussion on the implication of this assumption, we refer to Hausman et al. (1998). Assumption (A2)(i) ensures that the true misclassification probabilities are bounded away from 0 and 1. It is an important technical assumption, and is used throughout the proofs for obtaining upper bounds on the estimating function $h_{2,\beta,\theta}(\tilde{y}, \mathbf{x})$ (cf. (2.3)).

Assumption (A1) states the underlying true parameter β_0 may be unbounded. Assumption (A3) implies the existence of first and second moments. Hence, $\mathbf{E}X_j^2 \in (0, \infty)$ and $\mathbf{E}|X_j| < \infty$, for all $j = 1, \dots, p$. The positive definiteness assumption ensures that the components of \mathbf{X} are not linearly dependent among themselves, and it is an essential condition to ensure identifiability of the model (1.1). Assumptions (A1) and (A3) can be compared with some of the classical assumptions used for studying asymptotic properties of MLE's in simple logistic regression. For example, Amemiya (1985) assumes boundedness of β_0 and Gouriéroux and Monfort (1981) assumes that the covariates are bounded. Fahrmeir and Kaufmann (1985) studied MLE's in generalized linear models, and do not directly assume boundedness of β_0 or the covariates. However, they use other assumptions on the observed Fisher information matrix which are hard to verify, and are dependent on secondary sufficient conditions, among which one of conditions is a boundedness assumption on the covariates (cf. page 355 of Fahrmeir and Kaufmann (1985)). Compared to these restrictive assumptions, assumptions (A1) and (A3) are much weaker and easy to justify, but this leads to substantial technical complications in handling the proofs, and necessitates the use of empirical process tools. Assumptions (A1), (A2) and (A3) are related to each other, and as we will see later in Section 4.2, if one of the misclassification probabilities is set to zero, the unboundedness assumption on β_0 and the covariates has to be modified.

In assumption (A4), the first condition ensures that the limiting validation sampling fraction f is bounded away from 0 and 1, which in turn implies that both validation and non-validation sample sizes increase with the total sample size, and $\min\{n_1, n_2\} \rightarrow \infty$, as $n \rightarrow \infty$. The second condition in (A4) provides a convergence rate for $|f_n - f|$, as $n \rightarrow \infty$. Effectively, this ensures that $\{f_n : n \geq 1\}$ converges to f sufficiently quickly as the sample size increases.

Before stating the next assumption, we introduce the following notations. For any measurable function $h(y, \mathbf{x})$, we write, $\mathbf{P}_0 h(Y, \mathbf{X}) = \int h(y, \mathbf{x}) d\mathbf{P}_0(y, \mathbf{x})$, and similarly for a function $h(\tilde{y}, \mathbf{x})$. Consider the following nonrandom maps, $Z_1, Z_2, Z : \mathbb{R}^p \mapsto \mathbb{R}^p$, defined as,

$$\left. \begin{aligned} Z_1(\beta) &= \mathbf{P}_0 h_{1,\beta}(Y, \mathbf{X}) = \mathbf{E}[\mathbf{X} \cdot \{\psi(\mathbf{X}'\beta_0) - \psi(\mathbf{X}'\beta)\}], \\ Z_2(\beta) &= \mathbf{P}_0 h_{2,\beta,\theta_0}(\tilde{Y}, \mathbf{X}) \\ &= (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot \mathbf{E}\left[\mathbf{X} \cdot \frac{\psi(\mathbf{X}'\beta)\{1 - \psi(\mathbf{X}'\beta)\}\{\psi(\mathbf{X}'\beta_0) - \psi(\mathbf{X}'\beta)\}}{h_{3,\beta,\theta_0}(\mathbf{X})\{1 - h_{3,\beta,\theta_0}(\mathbf{X})\}}\right], \quad \text{and} \\ Z(\beta) &= f \cdot Z_1(\beta) + (1 - f) \cdot Z_2(\beta), \end{aligned} \right\} \quad (2.8)$$

where, $h_{1,\beta}$, $h_{2,\beta,\theta}$, $h_{3,\beta,\theta}$ are defined in (2.3), and f is the limiting validation sampling fraction defined in assumption (A3). Note that, in (2.8), \mathbf{E} denotes expectation with respect to the distribution function $Q(\mathbf{x})$ of the covariate \mathbf{X} . Let us denote the $p \times p$ matrix of partial derivatives of $Z(\beta)$ as $\dot{Z}(\beta)$.

(A5) Assume that $\dot{Z}(\beta_0)$ exists and is nonsingular.

This assumption ensures that the limiting score function $Z(\beta)$ has a non-singular derivative at β_0 , which is a commonly used assumption, similar to the non-singularity assumption of the Fisher information matrix used in maximum likelihood estimation. The matrix $\dot{Z}(\beta_0)$ is given in (8.41).

2.1.2 Modifications of technical assumptions in presence of an intercept term

The true model (1.1) does not explicitly include an intercept term. In case (1.1) includes an intercept term, we can write

$$\mathbf{P}_0(Y = 1 | (X_2, \dots, X_p)' = (x_2, \dots, x_p)') = \psi \left(\beta_{1,0} + \sum_{j=2}^p \beta_{j,0} x_j \right), \quad \text{for all } \mathbf{x} \in \mathbb{R}^{p-1}, \quad (2.9)$$

where, $\beta_0 = (\beta_{1,0}, \beta_{2,0}, \dots, \beta_{p,0})'$ is the unknown regression coefficient with intercept term $\beta_{1,0}$. To handle this model, we require a simple modification of assumption (A3).

(A3[†]) The first component of the p -dimensional covariate vector \mathbf{X} is equal to 1, and the remaining $(p - 1)$ components $(X_2, \dots, X_p)'$ satisfy the following assumptions:

- (i) $\mathbf{E}(X_j) = 0$, for all $j = 2, \dots, p$.
- (ii) $\mathbf{Var} \left((X_2, \dots, X_p)' \right)$ exists and is positive definite.

In Lemma 8.2, we make use of the assumption (A3[†]) to prove identifiability of the model with an intercept term. Assumption (A3[†])(ii) is similar to assumption (A3), and the additional condition in (A3[†])(i), which is required for centering of the covariates, is to ensure identifiability of the intercept term in (2.9). In the rest of the article, all theorems will be stated assuming the model (1.1) without an intercept term and using assumption (A3). For model (2.9) with an intercept term, one simply needs to replace assumption (A3) with (A3[†]).

2.2 Consistency and asymptotic normality of the pseudo-likelihood estimator

In order to study the limiting behavior of $\hat{\beta}_n$, the first step is to study the limit distribution of the estimated misclassification probabilities $\hat{\theta}_n$ (cf. (2.2)), which is given in Lemma 2.1. In order to describe the results, we introduce the following notations. Following the assumptions stated in Section 2.1.1, we define the following:

$$\mathbf{B}_0 = \left. \begin{array}{l} a_0 = \int \psi(\mathbf{x}'\beta_0) dQ(\mathbf{x}), \quad \pi_{2,0} = \theta_{2,0}a_0, \quad \pi_{3,0} = \theta_{1,0}(1 - a_0), \quad \text{and} \\ \begin{pmatrix} \pi_{3,0}/(1 - a_0)^2 & 0 & (1 - a_0 - \pi_{3,0})/(1 - a_0)^2 \\ \pi_{2,0}/a_0^2 & 1/a_0 & \pi_{2,0}/a_0^2 \end{pmatrix} \end{array} \right\}. \quad (2.10)$$

It should be noted that, $a_0 \in (0, 1)$. This follows by considering the properties of the logistic link function $\psi(\cdot)$, and combining assumptions (A1) and (A3), which state that the components of $\boldsymbol{\beta}_0$ are finite, and \mathbf{X} is non-degenerate and tight in \mathbb{R}^p . Along with assumption (A2), this implies $\pi_{2,0}, \pi_{3,0} \in (0, 1)$, and all elements of \mathbf{B}_0 in (2.10) are finite. We also define the matrix,

$$\boldsymbol{\Sigma}_{2,2} = \begin{pmatrix} \{1 - (a_0 + \pi_{3,0})\}(a_0 + \pi_{3,0}) & -\{1 - (a_0 + \pi_{3,0})\}\pi_{2,0} & -\{1 - (a_0 + \pi_{3,0})\}\pi_{3,0} \\ & \pi_{2,0}(1 - \pi_{2,0}) & -\pi_{2,0}\pi_{3,0} \\ & & \pi_{3,0}(1 - \pi_{3,0}) \end{pmatrix}. \quad (2.11)$$

Lemma 2.1. *Suppose, assumptions (A1) - (A4) hold. Then,*

$$\sqrt{n_1} (\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \xrightarrow{d} N_2(\mathbf{0}, \mathbf{B}_0 \boldsymbol{\Sigma}_{2,2} \mathbf{B}'_0), \quad (2.12)$$

where, $\boldsymbol{\Sigma}_{2,2}$ and \mathbf{B}_0 are defined in (2.11) and (2.10) respectively.

This result can be used for statistical inference on $\boldsymbol{\theta}_0$. The limit distribution of $\hat{\boldsymbol{\beta}}_n$, however, is affected by the asymptotic covariance of $\hat{\boldsymbol{\theta}}_n$.

Now we state the main result on the asymptotic behavior of $\hat{\boldsymbol{\beta}}_n$. In order to state our results, we need to invoke the definitions of the matrices $\boldsymbol{\Sigma}_{11}$, $\boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}'_{21}$, $\boldsymbol{\Sigma}_{22}$, $\boldsymbol{\Gamma}$, \mathbf{B}_0 and \mathbf{A}_0 , which are given in (8.37), (2.11), (8.38), (2.10) and (8.40) respectively. The explicit forms of these matrices are given in Section 8.

Theorem 2.2. *Suppose, assumptions (A1) - (A5) hold. Then the following statements are true:*

(i) *The pseudo-maximum likelihood estimator is consistent, i.e.,*

$$\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\| = o_{\mathbf{P}_0^*}(1).$$

(ii) *The pseudo-maximum likelihood estimator is asymptotically normal, i.e.,*

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) \xrightarrow{d} N_p\left(\mathbf{0}, [\dot{Z}(\boldsymbol{\beta}_0)]^{-1} \boldsymbol{\Sigma}_0 ([\dot{Z}(\boldsymbol{\beta}_0)]^{-1})'\right), \quad (2.13)$$

where,

$$\boldsymbol{\Sigma}_0 = f \cdot \boldsymbol{\Sigma}_{11} + (1 - f) \cdot \left\{ \mathbf{A}_0 \mathbf{B}_0 \boldsymbol{\Sigma}_{21} + (\mathbf{A}_0 \mathbf{B}_0 \boldsymbol{\Sigma}_{21})' \right\} + \frac{(1 - f)^2}{f} \cdot \mathbf{A}_0 \mathbf{B}_0 \boldsymbol{\Sigma}_{22} \mathbf{B}'_0 \mathbf{A}'_0 + (1 - f) \cdot \boldsymbol{\Gamma}, \quad (2.14)$$

and $\dot{Z}(\boldsymbol{\beta}_0)$ as defined in (8.41).

As stated earlier, asymptotic properties of $\hat{\boldsymbol{\beta}}_n$ cannot be studied by using results for Z-estimators based on i.i.d. summands. We apply the results for general Z-estimators based on arbitrary estimating equations (cf. Theorems 2.10 and 2.11 of Kosorok (2008)), which necessitates the verification of the required conditions. While verifying these conditions in our case, we face two major difficulties.

First, for proving the consistency of $\hat{\beta}_n$ a key step is to show that β_0 is a unique and well-separated zero of $Z(\beta)$ (cf. (2.8)). The difficulty arises because, in our case, the parameter space for β_0 is not compact and unlike the case of simple logistic regression, the limiting score function $Z(\beta)$ is not the gradient of a strictly concave function. If the limiting score function has a unique zero, then either of these two conditions are helpful in proving well-separatedness. However, in Lemmas 8.2 and 8.3, we show that it is possible to prove the uniqueness and well-separatedness of β_0 under the assumptions stated in Section 2.1.1.

The next difficulty is to properly handle the random indexing term $\hat{\theta}_n$ in the classes of functions, $\{h_{2,\beta,\hat{\theta}_n} : \beta \in \mathbb{R}^p\}$, $n \geq 1$. To be more precise, define the centered and scaled empirical processes,

$$\mathbb{G}_{n_i} = \sqrt{n_i} \cdot (\mathbb{P}_{n_i} - \mathbf{P}_0), \quad \text{for } i = 1, 2, \quad (2.15)$$

where, the meaning of \mathbf{P}_0 is clear from the context. In our case, verifying the Donsker property for the classes $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ and $\{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$ is not enough. We need to show that it is possible to replace the functions $h_{2,\beta,\hat{\theta}_n}$ with h_{2,β,θ_0} , with appropriate scaling and uniformly over β , as $\hat{\theta}_n$ converges to θ_0 . In Lemma 8.6 we have shown that, $\sup_{\beta} \|\mathbb{G}_{n_2}(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0})\|$ converges in probability to zero. This crucial uniform convergence result has been used in (a) verifying uniform convergence of Z_n while proving consistency of $\hat{\beta}_n$, and (b) verifying a stochastic equicontinuity condition about the process $\sqrt{n}(Z_n - Z)$ while proving asymptotic normality of $\hat{\beta}_n$.

Since we consider a specific model (cf. (1.1) and (1.2)) for the proofs, we could avoid the high-level stochastic equicontinuity assumptions on the empirical processes \mathbb{G}_{n_i} , $i = 1, 2$, and are able to handle the dependence between $Z_{n,1}(\beta)$ and $Z_{n,2}(\beta)$. For the same reason, we do not need a compactness assumption on the parameter space for β_0 or a boundedness assumption on the covariates. This can be compared with conditions used in articles with a much wider scope (cf. condition (S1) of Cheng and Huang (2010) or Assumption A.2 of Wellner and Zhan (1996)), where high-level assumptions are used.

The terms in the expression (2.14) for the covariance matrix Σ_0 clearly show the distinct contributions of the validation and non-validation parts of the sample. The first term is the contribution of the validation sample $\{(Y_i, \mathbf{X}_i) : 1 \leq i \leq n_1\}$, which follows immediately from the definition of $Z_{n,1}(\beta)$ (cf. (2.6)). The fourth term is due to the non-validation sample $\{(\tilde{Y}_i, \mathbf{X}_i) : n_1 + 1 \leq i \leq n\}$, which can be shown to be true by constructing a process similar to $\{Z_{n,2}(\beta) : \beta \in \mathbb{R}^p\}$, by replacing $\hat{\theta}_n$ with θ_0 . Finally, the second and third terms arise due to replacing the unknown θ_0 with its estimator $\hat{\theta}_n$ in the non-validation part of the estimating equation (2.7). Note that, the asymptotic covariance of $\hat{\theta}_n$ is embedded in the third term. The details are provided in the proofs.

The scaling factors in some of the terms in (2.14) are intriguing. The factors in the first and fourth terms are clearly due to the fact, that validation and non-validation samples comprise f_n and $(1 - f_n)$ proportions of the total sample, respectively. The scaling factor $(1 - f)$ in the second term of Σ_0 arises from the joint distribution of $\mathbb{P}_{n_1} h_{1,\beta}(Y, \mathbf{X})$ and $\hat{\theta}_n$, which depends on $\{(Y_i, \mathbf{X}_i) : 1 \leq i \leq n_1\}$. The scaling factor $(1 - f)^2/f$ in the third term is hard to anticipate without a theoretical derivation. It arises from the joint distribution

of $\mathbb{P}_{n_2} h_{2,\beta,\hat{\theta}_n}(Y, \mathbf{X})$ and $\hat{\boldsymbol{\theta}}_n$, and shows that the effect of increasing or decreasing f is asymmetric, and the scaling factor is non-linear in f .

Finally, notice that for drawing inference on $\boldsymbol{\beta}_0$, the asymptotic distribution in (2.13) is to be estimated. The elements of $\boldsymbol{\Sigma}_0$ and $\dot{Z}(\boldsymbol{\beta}_0)$ are complicated functions of the unknown parameters $\boldsymbol{\beta}_0, \boldsymbol{\theta}_0$, and are defined in terms of expected values with respect to the unknown distribution of \mathbf{X} . For estimation of $\boldsymbol{\Sigma}_0$, it is possible to plug-in $\hat{\boldsymbol{\beta}}_n$ and $\hat{\boldsymbol{\theta}}_n$ for $\boldsymbol{\beta}_0$ and $\boldsymbol{\theta}_0$, and replace the expectation with respect to $Q(\mathbf{x})$ by expectation with respect to the empirical version of $Q(\mathbf{x})$ based on the observed covariate values. However, a more convenient approach, especially from the point of view of implementation, is bootstrapping, which we describe in the next section.

3 Bootstrapped pseudo-likelihood estimator

In this section, we establish the distributional consistency of the bootstrapped PMLE in order to enable us to use the bootstrap approximation to the distribution of $\hat{\boldsymbol{\beta}}_n$ for statistical inference on $\boldsymbol{\beta}_0$. For the finite dimensional parameter $(\boldsymbol{\beta}_0, \boldsymbol{\theta}_0)$, we are interested in inference unconditional on \mathbf{X} . There can be several approaches to obtain a bootstrap sample in our setup. Since the distribution function $Q(\mathbf{x})$ is typically unknown (except for the assumptions made in (A3)), to obtain the results unconditional on \mathbf{X} , we would prefer to use the non-parametric bootstrap or the simple (Efron's) bootstrap (Efron (1979)). For bootstrapping, we select independent random samples with replacement of sizes n_1 and n_2 from the validation, and the non-validation samples, respectively. We denote these bootstrap samples as,

$$\hat{\mathcal{X}}_{n_1} = \{(\hat{Y}_i, \hat{Y}_i, \hat{\mathbf{X}}_i) : 1 \leq i \leq n_1\} \quad \text{and} \quad \hat{\mathcal{X}}_{n_2} = \{(\hat{Y}_i, \hat{\mathbf{X}}_i) : (n_1 + 1) \leq i \leq n\}.$$

Based on $\hat{\mathcal{X}}_{n_1}$, we define the bootstrap estimates of the misclassification probabilities $\hat{\boldsymbol{\theta}}_n = (\hat{\theta}_{1,n}, \hat{\theta}_{2,n})'$, similar to $\hat{\boldsymbol{\theta}}_n$ in (2.2), where

$$\hat{\theta}_{1,n} = \frac{\frac{1}{2} + \sum_{i=1}^{n_1} \mathbf{1}(\hat{Y}_i = 1, \hat{Y}_i = 0)}{1 + \sum_{i=1}^{n_1} \mathbf{1}(\hat{Y}_i = 0)} \quad \text{and} \quad \hat{\theta}_{2,n} = \frac{\frac{1}{2} + \sum_{i=1}^{n_1} \mathbf{1}(\hat{Y}_i = 0, \hat{Y}_i = 1)}{1 + \sum_{i=1}^{n_1} \mathbf{1}(\hat{Y}_i = 1)}, \quad n \geq 1. \quad (3.1)$$

The usual with replacement n -out-of- n bootstrap is a special case of the exchangeably weighted bootstrap (cf. Præstgaard and Wellner (1993)) with multinomially distributed weight vectors. For the bootstrap samples $\hat{\mathcal{X}}_{n_i}$, $i = 1, 2$, we define two independent multinomial random vectors \mathbf{M}_{n_i} , $i = 1, 2$, where

$$\mathbf{M}_{n_i} = (M_{n_i,1}, \dots, M_{n_i,n_i})' \sim \text{Multinomial} \left(n_i, \frac{1}{n_i}, \dots, \frac{1}{n_i} \right), \quad n_i \geq 1, \quad i = 1, 2,$$

where, $M_{n_i,k}$ denotes the frequency of occurrence of the k -th sample unit in the with replacement sample of size n_i , $i = 1, 2$. Clearly, \mathbf{M}_{n_1} and \mathbf{M}_{n_2} are independent of \mathcal{X}_n . Now, using the multinomial weights \mathbf{M}_{n_1} and \mathbf{M}_{n_2} , we define the weighted empirical measures

$$\hat{\mathbb{P}}_{n_1} = \frac{1}{n_1} \sum_{i=1}^{n_1} M_{n_1,i} \cdot \delta_{(Y_i, \mathbf{X}_i)} \quad \text{and} \quad \hat{\mathbb{P}}_{n_2} = \frac{1}{n_2} \sum_{i=n_1+1}^n M_{n_2,i} \cdot \delta_{(\tilde{Y}_i, \mathbf{X}_i)}, \quad n \geq 1,$$

where, $\delta_{(Y_i, \mathbf{X}_i)}$ and $\delta_{(\tilde{Y}_i, \mathbf{X}_i)}$ denote point masses at (Y_i, \mathbf{X}_i) and $(\tilde{Y}_i, \mathbf{X}_i)$, respectively. Following the definition of $Z_n(\boldsymbol{\beta})$ in (2.6), we define the corresponding bootstrap version,

$$\widehat{Z}_n(\boldsymbol{\beta}) = f_n \cdot \widehat{Z}_{n,1}(\boldsymbol{\beta}) + (1 - f_n) \cdot \widehat{Z}_{n,2}(\boldsymbol{\beta}), \quad \text{where, } \left. \begin{aligned} \widehat{Z}_{n,1}(\boldsymbol{\beta}) &= \widehat{\mathbb{P}}_{n_1} h_{1,\boldsymbol{\beta}} \\ \widehat{Z}_{n,2}(\boldsymbol{\beta}) &= \widehat{\mathbb{P}}_{n_2} h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n} \end{aligned} \right\}, \quad (3.2)$$

where, $\widehat{\boldsymbol{\theta}}_n$ is defined in (3.1). The bootstrapped PMLE, denoted by $\widehat{\boldsymbol{\beta}}_n$, satisfies the estimating equation:

$$\widehat{Z}_n(\widehat{\boldsymbol{\beta}}_n) = \mathbf{0}. \quad (3.3)$$

It should be noted that there is an extensive literature on the asymptotic properties of bootstrap methods in the context of general M-estimation. Arcones and Giné (1992) study almost sure convergence results for bootstrapped M-estimators under strong assumptions on the underlying parameter space and estimating functions. Chatterjee and Bose (2005) investigate generalized bootstrap methods for estimating equations, and their results have a wide scope. However, their results cannot handle different exchangeable weights like \mathbf{M}_{n_1} and \mathbf{M}_{n_2} , and non-identical random variables (Y, \mathbf{X}) and (\tilde{Y}, \mathbf{X}) . Wellner and Zhan (1996) study consistency of exchangeably weighted bootstrap for Z-estimators, but their results are valid for i.i.d. observations, and also require specific measurability assumptions on the class of estimating functions (cf. Giné and Zinn (1990)). In an important piece of work, Cheng and Huang (2010) study consistency of exchangeably weighted bootstrap for Z-estimators in general semiparametric setup in the presence of infinite dimensional nuisance parameters using estimating equations based on i.i.d. observations. The generality and wide scope of these results require assumptions like compactness of the finite dimensional parameter space, stochastic equicontinuity of the underlying empirical process, measurability assumptions on the class of estimating functions, and smoothness conditions on the limiting score functions. Besides, while proving bootstrap distributional consistency, in most cases consistency of the estimators of finite dimensional parameters, and their bootstrapped version is assumed.

However, existing bootstrap consistency results for Z-estimators, including those mentioned above, cannot be directly used for finding the asymptotic distribution of the bootstrapped pseudo MLE $\widehat{\boldsymbol{\beta}}_n$. The primary difficulty arises from the dependence between $\widehat{Z}_{n,1}(\boldsymbol{\beta})$ and $\widehat{Z}_{n,2}(\boldsymbol{\beta})$ because of the presence of $\widehat{\boldsymbol{\theta}}_n$ in the latter. Also, usual empirical process results are not directly applicable to the estimated (random) class of functions $\{h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n}(\tilde{y}, \mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$. Finally, the processes $\{\widehat{Z}_{n,i}(\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$, $i = 1, 2$, are not based on the same set of multinomial weights and involve different sets of random variables (either (Y, \mathbf{X}) or (\tilde{Y}, \mathbf{X})), unlike standard with replacement bootstrap procedure. In order to deal with these additional difficulties, the asymptotic behavior of each of these two processes needs to be studied separately.

In order to study the asymptotic behavior of the bootstrapped PMLE $\widehat{\boldsymbol{\beta}}_n$, we need to prove two new equicontinuity results. In particular, we show that

$$\sup_{\boldsymbol{\beta}} \|\widehat{\mathbb{G}}_{n_2}(h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n})\| \quad \text{and} \quad \sup_{\boldsymbol{\beta}} \|\mathbb{G}_{n_2}(h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\widehat{\boldsymbol{\theta}}_n})\|,$$

converge to zero in probability in an appropriate sense. Unlike the equicontinuity result proved for the PMLE $\widehat{\beta}_n$, for proving these new results, the results of [van der Vaart and Wellner \(2007\)](#) can not be directly used. We show that their approach can be extended for handling both these cases.

As evident from the above discussion, stochastic convergence of the above mentioned supremums has to be carefully defined, since the random quantities are defined on different probability spaces. The details are given in the proofs. This type of stochastic equicontinuity results for bootstrapped empirical processes indexed by a fixed and a random parameter are new, and the line of arguments can be possibly extended to develop bootstrap consistency results for similar problems.

3.1 Asymptotic results for the bootstrapped pseudo-likelihood estimator

In order to state the bootstrap consistency results, we introduce some concepts and definitions following [Wellner and Zhan \(1996\)](#) and [Cheng and Huang \(2010\)](#).

Note that, $\{(Y_i, \tilde{Y}_i, \mathbf{X}_i) : i \geq 1\}$ are i.i.d. observations from a distribution \mathbf{P}_0 on a probability space $(\mathcal{X}, \mathcal{A})$. To deal with measurability issues, we view $(Y_i, \tilde{Y}_i, \mathbf{X}_i)$ as the i -th coordinate projection from the underlying canonical product probability space $(\mathcal{X}^\infty, \mathcal{A}^\infty, \mathbf{P}_0^\infty)$ into the i -th copy of \mathcal{X} . We assume that the multinomial weight vectors $\{\mathbf{M}_{n_1} : n_1 \geq 1\}$ and $\{\mathbf{M}_{n_2} : n_2 \geq 1\}$ are independent of $\{(Y_i, \tilde{Y}_i, \mathbf{X}_i) : i \geq 1\}$, and form a triangular array which is defined on some probability space $(\mathcal{Z}, \mathcal{C}, \mathbf{P}_M)$, for all $n_1, n_2 \geq 1$. To handle the joint randomness of $\{(Y_i, \tilde{Y}_i, \mathbf{X}_i) : i \geq 1\}$ and $(\mathbf{M}_{n_1}, \mathbf{M}_{n_2})$, we define the product probability space

$$(\mathcal{X}^\infty, \mathcal{A}^\infty, \mathbf{P}_0^\infty) \times (\mathcal{Z}, \mathcal{C}, \mathbf{P}_M) = (\mathcal{X}^\infty \times \mathcal{Z}, \mathcal{A}^\infty \times \mathcal{C}, \mathbf{Pr}),$$

where, $\mathbf{Pr} = \mathbf{P}_0^\infty \times \mathbf{P}_M$, is a product probability measure. For simplicity, we will denote the product measure \mathbf{P}_0^∞ as \mathbf{P}_0 . We will use the symbols \mathbf{Pr}^* (\mathbf{Er}^*), \mathbf{P}_0^* (\mathbf{E}_0^*) to denote outer probabilities (expectations) corresponding to \mathbf{Pr} and \mathbf{P}_0 , respectively. Usual expectation with respect to \mathbf{P}_M and \mathbf{Pr} will be denoted by \mathbf{E}_M and \mathbf{Er} respectively.

We now define stochastic orders based on the above probability measures. The definitions would be repeatedly used in the proofs, and are crucial for a clear understanding of the results. Consider a sequence of real valued functions $\{\Delta_n : n \geq 1\}$, defined on the product probability space $(\mathcal{X}^\infty \times \mathcal{Z}, \mathcal{A}^\infty \times \mathcal{C}, \mathbf{Pr})$. We say that, $\Delta_n = o_{\mathbf{P}_M^*}^*(1)$ in \mathbf{P}_0^* -probability, if for any $\epsilon, \eta > 0$,

$$\mathbf{P}_0^* (\mathbf{P}_M^* (|\Delta_n| > \epsilon) > \eta) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Similarly, $\Delta_n = O_{\mathbf{P}_M^*}^*(1)$ in \mathbf{P}_0^* -probability if, for any $\delta > 0$ and $\epsilon > 0$, there exists a $M = M(\epsilon, \delta) \in (0, \infty)$ and an integer $n_0 = n_0(\epsilon, \delta) \in \mathbb{N}$, such that

$$\mathbf{P}_0^* (\mathbf{P}_M^* (|\Delta_n| > M) > \delta) < \epsilon, \quad \text{for all } n \geq n_0(\epsilon, \delta).$$

[Lemma 8.11](#) describes some of the relations among the stochastic orders in terms of different probability measures, and will be heavily used in clearly describing bootstrap convergence results. For a few other

similar results, we refer to Lemma 3 of [Cheng and Huang \(2010\)](#).

Lemma 3.1 describes the limiting behavior of $\widehat{\boldsymbol{\theta}}_n$.

Lemma 3.1. *Suppose, assumptions (A1)-(A4) hold. Then,*

- (i) $\sqrt{n_1} \left(\widehat{\boldsymbol{\theta}}_n - \widehat{\boldsymbol{\theta}}_n \right) \xrightarrow{d} N_2(\mathbf{0}, \mathbf{B}_0 \boldsymbol{\Sigma}_{2,2} \mathbf{B}'_0)$, conditionally almost surely (\mathbf{P}_0), where \mathbf{B}_0 is defined in (2.10).
- (ii) $\sqrt{n_1} \left(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0 \right) \xrightarrow{d} N_2(\mathbf{0}, 2\mathbf{B}_0 \boldsymbol{\Sigma}_{2,2} \mathbf{B}'_0)$, unconditionally. The unconditional distribution of $\widehat{\boldsymbol{\theta}}_n$ is with respect to the product probability measure \mathbf{Pr} .

Part (i) of Lemma 3.1 states the bootstrap estimator $\widehat{\boldsymbol{\theta}}_n$ is distributionally consistent, and this can be used to carry out bootstrap based inference on $\boldsymbol{\theta}_0$. Also, we require this result for proving the bootstrap consistency of $\widehat{\boldsymbol{\beta}}_n$. The unconditional convergence result stated in part (ii) of Lemma 3.1 is a major technical tool which is used in developing some of the new bootstrap stochastic equicontinuity results. Now we can state the main result on bootstrap consistency.

Theorem 3.2. *Suppose, assumptions (A1) - (A5) hold. Then, the following statements are true:*

- (i) *The bootstrapped PMLE $\widehat{\boldsymbol{\beta}}_n$ is consistent in \mathbf{P}_0^* -probability, i.e.,*

$$\|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\| = o_{\mathbf{P}_0^*}(1), \quad \text{in } \mathbf{P}_0^*\text{-probability.}$$

- (ii) *The conditional distribution of $\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \widehat{\boldsymbol{\beta}}_n)$, given the sample \mathcal{X}_n , consistently estimates the distribution of $\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)$ in the following sense,*

$$\sup_{\mathbf{t} \in \mathbb{R}^p} \left| \mathbf{P}_M^* \left(\sqrt{n} \left(\widehat{\boldsymbol{\beta}}_n - \widehat{\boldsymbol{\beta}}_n \right) \leq \mathbf{t} \right) - \mathbf{P}_0 \left(\sqrt{n} \left(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0 \right) \leq \mathbf{t} \right) \right| = o_{\mathbf{P}_0^*}(1),$$

where, the inequality “ \leq ” is understood to be componentwise.

Note that, the bootstrap consistency result in Theorem 3.2 does not require any additional technical conditions beyond what was needed while proving asymptotic normality of the PMLE $\widehat{\boldsymbol{\beta}}_n$. The uniform convergence result, in part (ii) above, holds for classes of sets, larger than all p -dimensional rectangles of the form $(-\infty, t_1] \times \cdots \times (-\infty, t_p]$, where $\mathbf{t} = (t_1, \dots, t_p)' \in \mathbb{R}^p$. The proof of part (ii) shows that, $\sqrt{n}(\widehat{\boldsymbol{\beta}}_n - \widehat{\boldsymbol{\beta}}_n)$ converges conditionally in distribution (in outer probability) to the same limiting Gaussian distribution, which is described in Theorem 2.2(ii). Using the absolute continuity of the Gaussian distribution, and the results on uniformity classes (cf. Section 1.2 of [Bhattacharya and Ranga Rao \(1986\)](#)), it can be shown that uniform convergence in part (ii) is valid over all Borel-measurable convex subsets of \mathbb{R}^p (cf. page 55 of [Lahiri \(2003\)](#)).

3.2 Applications of bootstrap method in some inference problems

Theorem 3.2 enables us to use the bootstrap method for a variety of inference problems on β_0 . More generally, the delta method for bootstrapped estimators can be used for inference on appropriate functions of β_0 . Inference on regression coefficient β_0 itself is of interest to assess the effect of the covariates on the probability of occurrence of the event under study. Specifically, in epidemiology, β_{0j} , the j -th component of β_0 , represents the log-odds ratio of a disease associated with a unit increase in the scale of x_j , holding all other variables in \mathbf{x} fixed. If x_j is a binary variable, then β_{0j} represents the adjusted log-odds ratio adjusted for the covariates in the regression model other than x_j . Inference on a linear parametric function $\mathbf{c}'\beta_0$ may be of interest in many applications, e.g., in assessing credit worthiness of a person by a credit bureau. Also estimating $\psi(\mathbf{x}'\beta_0)$ for a given value of \mathbf{x} , say, \mathbf{x}_0 is of interest as it represents the probability of the event under study to happen if $\mathbf{x} = \mathbf{x}_0$. Especially, we may be interested in the estimation of $\psi(\mathbf{x}'\beta_0)$, wherein it represents the risk of an adverse event. For example, in epidemiological studies, it may represent the risk of a disease associated with the risk factors \mathbf{x} ; in financial applications, it may represent the risk of a default associated with an individual's profile \mathbf{x} . In the following, we illustrate two common applications of Theorem 3.2.

3.2.1 Interval estimation of a linear parametric function

We can use Theorem 3.2 to construct an asymptotically consistent confidence interval for $\mathbf{c}'\beta_0$. The confidence interval is,

$$I_{n,\eta} = \left(\mathbf{c}'\hat{\beta}_n - n^{-1/2} \cdot \hat{\xi}_{n,1-\eta}(\mathbf{c}), \mathbf{c}'\hat{\beta}_n - n^{-1/2} \cdot \hat{\xi}_{n,\eta}(\mathbf{c}) \right), \quad (3.4)$$

where, $\hat{\xi}_{n,\eta}(\mathbf{c})$ is the η -quantile of $\sqrt{n}(\mathbf{c}'\hat{\beta}_n - \mathbf{c}'\beta_0)$, for any $\eta \in (0, 1)$. As stated in the next corollary, $I_{n,\eta}$ is an asymptotically consistent level $(1 - 2\eta)$ confidence interval for $\mathbf{c}'\beta_0$.

Corollary 3.3. *Under the conditions of Theorem 3.2, for any $\eta \in (0, 1/2)$ and any fixed $\mathbf{c} \in \mathbb{R}^p$,*

$$\mathbf{P}_0^* (\mathbf{c}'\beta_0 \in I_{n,\eta}) \rightarrow (1 - 2\eta).$$

3.2.2 Interval estimation of Risk

The risk at \mathbf{x}_0 is given by

$$\pi_{\mathbf{x}_0}(\beta_0) \equiv \psi(\mathbf{x}_0'\beta_0).$$

To avoid trivialities, we assume that $\mathbf{x}_0 \neq \mathbf{0}$. The naive estimator of risk at \mathbf{x}_0 is $\pi_{\mathbf{x}_0}(\hat{\beta}_n) = \psi(\mathbf{x}_0'\hat{\beta}_n)$, and the bootstrap estimator is $\pi_{\mathbf{x}_0}(\hat{\hat{\beta}}_n) = \psi(\mathbf{x}_0'\hat{\hat{\beta}}_n)$. For any $\eta \in (0, 1/2)$, let $\hat{\kappa}_{n,\eta}$ denote the η -quantile of $\sqrt{n}(\pi_{\mathbf{x}_0}(\hat{\hat{\beta}}_n) - \pi_{\mathbf{x}_0}(\hat{\beta}_n))$. Define the following confidence interval for $\pi_{\mathbf{x}_0}(\beta_0)$,

$$J_{n,\eta} = \left(\pi_{\mathbf{x}_0}(\hat{\beta}_n) - n^{-1/2} \cdot \hat{\kappa}_{n,1-\eta}, \pi_{\mathbf{x}_0}(\hat{\beta}_n) - n^{-1/2} \cdot \hat{\kappa}_{n,\eta} \right). \quad (3.5)$$

Corollary 3.4. *Under the conditions of Theorem 3.2, for any $\eta \in (0, 1/2)$,*

$$\mathbf{P}_0^* (\pi_{\mathbf{x}_0}(\boldsymbol{\beta}_0) \in J_{n,\eta}) \rightarrow (1 - 2\eta).$$

4 Some additional results

Motivated by the practical applications, we consider two specific situations, wherein the asymptotic results proved in Sections 2 and 3 are not directly applicable. In this section, we show that, our results continue to hold with appropriate modifications.

4.1 Differential classification error

Until now, we have developed the asymptotic theory, assuming that, the probabilities of misclassification are independent of the covariates \mathbf{x} . In other words, we have assumed a non-differential classification error. However, as mentioned by Meyer and Mittag (2016), "there is little ex ante reason to believe that misclassification is independent of the covariates." In a HIV related study, Lyles et al. (2011) consider a regression model for the misclassification probabilities with subject specific binary covariates. In social and economic surveys (cf. Bollinger and David (1997) and Abrevaya and Hausman (1999)), it has been observed that, the misclassification probabilities of the binary response (like, whether beneficiary of a program or not) depend on the covariates, like union membership of a worker (a member or not), income (above or below median income), age (above the median age or not) etc.

We now consider extension of our theory to differential classification errors. For simplicity, we assume that, the sample space for \mathbf{x} can be divided into K non-overlapping subsets, each representing the profile of a distinct group, possibly with different misclassification probabilities. Denote these K groups by G_1, \dots, G_K . We assume that, a random sample of size n_k is obtained from G_k , with K remaining fixed, and $\min\{n_k : 1 \leq k \leq K\} \rightarrow \infty$. The underlying true model is,

$$\left. \begin{aligned} \mathbf{P}_0(Y = 1 \mid \mathbf{X} = \mathbf{x}) &= \psi(\mathbf{x}'\boldsymbol{\beta}_0), \quad \text{and} \\ \mathbf{P}_0(\tilde{Y} = 1 \mid Y = 0) &= \theta_{1,0}^{(k)}, \quad \mathbf{P}_0(\tilde{Y} = 0 \mid Y = 1) = \theta_{2,0}^{(k)}, \end{aligned} \right\} \text{ for all } \mathbf{x} \in G_k, k \in \{1, \dots, K\}, \quad (4.1)$$

where, $\boldsymbol{\theta}_0^{(k)} = (\theta_{1,0}^{(k)}, \theta_{2,0}^{(k)})'$ are the misclassification probabilities for the group G_k . Writing $f_{n,k} = n_{1,k}/n_k$, where $n_{1,k}$ and $n_{2,k} = (n_k - n_{1,k})$ are the validation and the non-validation sample sizes, respectively, we assume that assumption (A2) remains valid with $\boldsymbol{\theta}_0$ replaced by $\boldsymbol{\theta}_0^{(k)}$. Similarly, we assume that assumption (A4) remains true in the following sense: $f_{n,k} \rightarrow f_k \in (0, 1)$ and $\sqrt{n_k} \cdot |f_{n,k} - f_k| \rightarrow 0$ as $n \rightarrow \infty$. Assumptions (A1) and (A3) or (A3[†]) remain unchanged. Based on the observations from G_k , one can estimate $\hat{\boldsymbol{\theta}}_n^{(k)}$ using

(2.2). Following similar arguments as above, the estimating function can be written as,

$$Z_{n,K}(\boldsymbol{\beta}) = \sum_{k=1}^K \left[f_{n,k} \cdot \mathbb{P}_{n_{1,k}} h_{1,\boldsymbol{\beta}}(Y, \mathbf{X}) + (1 - f_{n,k}) \cdot \mathbb{P}_{n_{2,k}} h_{2,\boldsymbol{\beta}, \hat{\boldsymbol{\theta}}_n^{(k)}}(\tilde{Y}, \mathbf{X}) \right],$$

where, $\mathbb{P}_{n_{1,k}}$ and $\mathbb{P}_{n_{2,k}}$ denote the empirical measures corresponding to validation and non-validation samples in group G_k , respectively. The PMLE $\hat{\boldsymbol{\beta}}_n$ satisfies the estimating equation,

$$Z_{n,K}(\boldsymbol{\beta}) = \mathbf{0}. \quad (4.2)$$

Following the arguments given in Lemma 2.1, and using the above mentioned assumptions on $\boldsymbol{\theta}_0^{(k)}$ and $f_{n,k}$, it can be shown that $\hat{\boldsymbol{\theta}}_n^{(k)}$ will be asymptotically normal. Also, using the arguments similar to those given in Lemma 8.7, it will follow that, $\sup_{\boldsymbol{\beta}} \|Z_{n,K}(\boldsymbol{\beta}) - Z_K(\boldsymbol{\beta})\| = o_{\mathbf{P}_0^*}(1)$, where,

$$Z_K(\boldsymbol{\beta}) = \sum_{k=1}^K \left[f_k \cdot \mathbf{P}_0 h_{1,\boldsymbol{\beta}}(Y, \mathbf{X}) + (1 - f_k) \cdot \mathbf{P}_0 h_{2,\boldsymbol{\beta}, \boldsymbol{\theta}_0^{(k)}}(\tilde{Y}, \mathbf{X}) \right].$$

Further, following arguments similar to that used in Lemma 8.2 and 8.3, it can be shown that $Z_K(\boldsymbol{\beta})$ has a unique and well-separated root at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$. Specifically, note that $\sum_k f_k \cdot \mathbf{P}_0 h_{1,\boldsymbol{\beta}} = f \cdot \mathbf{P}_0 h_{1,\boldsymbol{\beta}} = f \cdot Z_1(\boldsymbol{\beta})$, where f is defined in (A4), and $Z_1(\boldsymbol{\beta})$ is defined in (2.8). So, the arguments in Lemma 8.2 can be repeated in this case for $Z_1(\boldsymbol{\beta})$, and the rest of the proof follows by sandwiching $Z_K(\boldsymbol{\beta})$ within constant positive multiples of $Z_1(\boldsymbol{\beta})$, as shown in Lemma 8.3. The proof of Donsker and Glivenko-Cantelli properties can be carried out similarly, and this will imply consistency of $\hat{\boldsymbol{\beta}}_n$ (cf. (4.2)). Next, we will assume that assumption (A5) holds if $Z(\boldsymbol{\beta})$ (cf. (2.8)) is replaced by $Z_K(\boldsymbol{\beta})$. Then, following the arguments in case of a single group, it can be shown that $\hat{\boldsymbol{\beta}}_n$ will be asymptotically normal, and the bootstrap version of $\hat{\boldsymbol{\beta}}_n$, which can be constructed in the same manner by independently resampling within each group, will also be distributionally consistent in probability. The asymptotic variance of $\hat{\boldsymbol{\beta}}_n$ defined in (4.2) can be obtained using the variance expression given in (2.13). The primary differences will be, $\dot{Z}(\boldsymbol{\beta}_0)$ will be replaced by $\dot{Z}_K(\boldsymbol{\beta}_0)$ and $\boldsymbol{\Sigma}_0$ in (2.14) will be replaced by,

$$\begin{aligned} \boldsymbol{\Sigma}_{0,K} = & f \cdot \boldsymbol{\Sigma}_{11} + \sum_{k=1}^K \left[(1 - f_k) \cdot \left\{ \mathbf{A}_{0,k} \mathbf{B}_{0,k} \boldsymbol{\Sigma}_{21,k} + (\mathbf{A}_{0,k} \mathbf{B}_{0,k} \boldsymbol{\Sigma}_{21,k})' \right\} \right. \\ & \left. + \frac{(1 - f_k)^2}{f_k} \cdot \mathbf{A}_{0,k} \mathbf{B}_{0,k} \boldsymbol{\Sigma}_{22,k} \mathbf{B}'_{0,k} \mathbf{A}'_{0,k} + (1 - f_k) \cdot \boldsymbol{\Gamma}_k \right], \end{aligned}$$

where, the matrices $\boldsymbol{\Sigma}_{12,k} = \boldsymbol{\Sigma}'_{21,k}$, $\boldsymbol{\Sigma}_{22,k}$, $\boldsymbol{\Gamma}_k$, $\mathbf{A}_{0,k}$ and $\mathbf{B}_{0,k}$, are similarly defined as in (8.37), (2.11), (8.38), (8.40) and (2.10) respectively, by replacing $\boldsymbol{\theta}_0$ with $\boldsymbol{\theta}_0^{(k)}$.

4.2 One of the misclassification probabilities equal to zero

The classification errors are usually asymmetric. An extreme case of asymmetry could be, one of the misclassification probabilities is equal to zero, while the other is non-zero. For example, it may be that

$\mathbf{P}_0(\tilde{Y} = 0 \mid Y = 1) = \theta_{2,0} = 0$, while $\mathbf{P}_0(\tilde{Y} = 1 \mid Y = 0) = \theta_{1,0}$ is non-zero. In this case, the parameter $\boldsymbol{\theta}_0$ reduces to a scalar parameter $\theta_{1,0}$. While discussing about literacy data, collected through census in India, [Kothari and Bandyopadhyay \(2011\)](#) mention that the indirect method of determining the literacy status of an individual leads to misclassification. In this case, the chances of misclassifying a literate person ($Y = 1$) as an illiterate ($\tilde{Y} = 0$) are negligible and can be practically considered equal to zero. However, the chance of other type of misclassification is significantly high. The detailed discussion is deferred to Section 5.4, where this particular example has been analysed. This type of extreme asymmetrical misclassification may arise in medical diagnostic studies, [Demidenko \(2004\)](#) (pp. 512) provides such an example related to cancer detection. The unknown non-zero misclassification probability $\theta_{1,0}$ can be estimated using (2.2) and for purposes of computation, we can define, $\hat{\theta}_{2,n} = 0$ with probability 1. The pseudo-likelihood based estimating equation in (2.7) can be used with this modified definition of $\hat{\boldsymbol{\theta}}_n = (\hat{\theta}_{1,n}, 0)'$. In this case, under assumption (A2)(i), we have

$$\delta_1 < h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}) = h_{3,\boldsymbol{\beta},\theta_1}(\mathbf{x}) = \theta_{1,0} + (1 - \theta_{1,0}) \cdot \psi(\mathbf{x}'\boldsymbol{\beta}) < 1, \quad (\text{cf. (2.3)}), \text{ for all } \mathbf{x}, \boldsymbol{\beta} \text{ and } \theta_1 \in (\delta_1, \delta_2).$$

As a result, $h_{3,\boldsymbol{\beta},\theta_1}(\mathbf{x})$ remains bounded away from 0, but not bounded away from 1. This implies, $(1 - h_{3,\boldsymbol{\beta},\theta_1}(\mathbf{x}))^{-1}$ can become arbitrarily large. This is unlike the case when both misclassification probabilities are present in the model, and are bounded away from 0 and 1. As a result, all technical arguments used in the general case will fail in this situation.

In order to extend the theoretical results to this case, we need to modify some of the original technical assumptions stated in Section 2.1.1. Specifically, we strengthen assumption (A1) to ensure boundedness of the coefficient $\boldsymbol{\beta}_0$, assumption (A2)(i) is only applicable on the non-zero misclassification probability $\theta_{1,0}$ and it remains bounded away from 0 and 1, assumption (A3) is also strengthened and we assume boundedness of the covariate \mathbf{X} along with positive definiteness of $\mathbf{Var}(\mathbf{X})$. Assumption (A2)(ii) is redundant in this case, and assumptions (A4) and (A5) remain unchanged. For the model with intercept (2.9), the same modifications are applicable along with the changes described in assumption (A3[†]).

With these new set of assumptions, $\psi(\mathbf{x}'\boldsymbol{\beta}_0)$ remains bounded away from 0 and 1, and as a result, $h_{3,\boldsymbol{\beta},\theta_1}(\mathbf{x})$ is also bounded away from 0 and 1, for all \mathbf{x} , $\theta_{1,0} \in (\delta_1, \delta_2)$ and $\boldsymbol{\beta} \in \mathbb{R}^p$. It is now possible to prove the theoretical results about the PMLE $\hat{\boldsymbol{\beta}}_n$, and it's bootstrap version, by retracing the technical arguments for the general case, replacing $\hat{\boldsymbol{\theta}}_n$ by $\hat{\theta}_{1,n}$ and $\hat{\boldsymbol{\theta}}_n$ by $\hat{\theta}_{1,n}$. The boundedness of $\boldsymbol{\beta}_0$ will simplify some of the convergence arguments used in the proofs. We skip the details. There will be some changes while computing the expression for asymptotic covariance of $\hat{\boldsymbol{\beta}}_n$ (cf. (2.13) and (2.14)), which depend on $\boldsymbol{\theta}_0$. For all matrices, the unknown $\theta_{2,0}$ is to be replaced by $\theta_{2,0} = 0$. Since $\theta_{2,0}$ is not estimated, the asymptotic covariance expression of $\hat{\boldsymbol{\theta}}_n$ (cf. Lemma 2.1), which is used in the third term of (2.14) will be slightly changed. The matrix $\mathbf{B}_0 \boldsymbol{\Sigma}_{22} \mathbf{B}'_0$ should be replaced by a 2×2 matrix, with (1,1)-th element equal to the corresponding element of $\mathbf{B}_0 \boldsymbol{\Sigma}_{22} \mathbf{B}'_0$ and all other elements equal to zero.

Apparently a simplification of the model (1.2), with one of the misclassification probabilities being equal

to zero, leads to the use of stronger assumptions requiring boundedness of the covariates and the true regression coefficient. At present, we do not know of an alternative proof under weaker assumptions.

5 Simulation results and real data analysis

5.1 Comparisons with other likelihood based estimation methods

As mentioned in Section 1, we carry out extensive simulation studies to compare the performances of PMLE with other likelihood based estimators, one based only on contaminated (misclassified) responses, considered by Hausman et al. (1998) and Hausman (2001) in the absence of any validation sample, and the other, based on the joint likelihood of both validation and non-validation samples considered by Lyles et al. (2011). We refer to them as CMLE, and JMLE, respectively. In the following, we briefly discuss the likelihood functions, and the associated score equations for finding CMLE and JMLE. Computation of the likelihood estimates, including PMLE involves solving p or $(p + 2)$ (in case of JMLE and CMLE) nonlinear equations with the same number of variables. For solving the nonlinear equations, we have used R-package BB, developed by Varadhan and Gilbert (2009).

5.1.1 CMLE: Likelihood & score equations

In the absence of a validation sample, information on the true responses is not available. In this case, Hausman et al. (1998) suggested jointly estimating $(\boldsymbol{\beta}, \boldsymbol{\theta})$ using the likelihood based on the set of observations $\{(\tilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq n\}$. However, this approach has some serious drawbacks, which we discuss briefly. As noted earlier (cf. (1.3)), $[\tilde{Y}_i | \mathbf{X}_i = \mathbf{x}_i] \sim \text{Bernoulli}(h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}_i))$, for all $i = 1, \dots, n$, and are independent, where $h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x})$ is defined in (2.3). The likelihood function will be,

$$L_{n,C}(\boldsymbol{\beta}, \boldsymbol{\theta}) = \prod_{i=1}^n (h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}_i))^{y_i} \cdot (1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{x}_i))^{1-y_i}.$$

The estimating equations for $(\boldsymbol{\beta}, \boldsymbol{\theta})$ are,

$$\left. \begin{aligned} \frac{\partial}{\partial \boldsymbol{\beta}} n^{-1} \log L_{n,C}(\boldsymbol{\beta}, \boldsymbol{\theta}) &= \mathbb{P}_n \left[\mathbf{X} \cdot \frac{\psi(\mathbf{X}'\boldsymbol{\beta}) \cdot \{1 - \psi(\mathbf{X}'\boldsymbol{\beta})\} \cdot \{\tilde{Y} - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X}) \cdot \{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}} \right] = \mathbf{0}, \\ \frac{\partial}{\partial \boldsymbol{\theta}_1} n^{-1} \log L_{n,C}(\boldsymbol{\beta}, \boldsymbol{\theta}) &= \mathbb{P}_n \left[\frac{\{1 - \psi(\mathbf{X}'\boldsymbol{\beta})\} \cdot \{\tilde{Y} - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X}) \cdot \{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}} \right] = 0, \\ \frac{\partial}{\partial \boldsymbol{\theta}_2} n^{-1} \log L_{n,C}(\boldsymbol{\beta}, \boldsymbol{\theta}) &= \mathbb{P}_n \left[\frac{\psi(\mathbf{X}'\boldsymbol{\beta}) \cdot \{\tilde{Y} - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X}) \cdot \{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}}(\mathbf{X})\}} \right] = 0, \end{aligned} \right\} \quad (5.1)$$

where, $\mathbb{P}_n = n^{-1} \sum_{i=1}^n \delta_{(\tilde{Y}_i, \mathbf{X}_i)}$, denotes the empirical measure based on $\{(\tilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq n\}$. Denote these estimates by $(\hat{\boldsymbol{\beta}}_{n,C}, \hat{\boldsymbol{\theta}}_{n,C})$. Often, the estimating equations in (5.1) become nearly non-identifiable, and thus, lead to nonsensical estimates $(\hat{\boldsymbol{\beta}}_{n,C}, \hat{\boldsymbol{\theta}}_{n,C})$ for the following reasons.

It has been noted by [Cox and Snell \(1989\)](#), that for every β , there exists a β^* , such that $\psi(\mathbf{x}'\beta) \approx \mathbf{x}'\beta^*$ for $\psi(\mathbf{x}'\beta) \in (0.1, 0.9)$. Thus, the logistic and linear functions are almost identical on a large part of their range, except for the tails. Hence, for $\psi(\mathbf{x}'\beta) \in (0.1, 0.9)$, we have,

$$h_{3,\beta,\theta}(\mathbf{x}) = \theta_1 + (1 - \theta_1 - \theta_2) \cdot \psi(\mathbf{x}'\beta) \approx \theta_1 + (1 - \theta_1 - \theta_2) \cdot \mathbf{x}'\beta^* = \theta_1 + \mathbf{x}'\beta_*, \quad (5.2)$$

where, $\beta_* = (1 - \theta_1 - \theta_2) \cdot \beta^*$. Thus, if (θ_1, θ_2) are unknown, the likelihood estimates based on $L_{n,C}$ fails to recover the estimates of the model parameters, unless there are enough number of observations in the tails of $\psi(\mathbf{x}'\beta)$. Thus, if the covariate \mathbf{X} is such that, the probability

$$Q(\mathbf{X} : \psi(\mathbf{X}'\beta_0) \in (0.1, 0.9)) \quad \text{is high,}$$

then recovering estimates of the model parameters from $L_{n,C}$ is nearly impossible, unless the sample size is very large.

5.1.2 JMLE: Likelihood & score equations

If the validation sample is available, then a natural approach is to consider the joint likelihood function $L_{1,n}(\beta, \theta)$, given in (2.1). Following earlier calculations, the estimating equations for (β, θ) will be,

$$\left. \begin{aligned} \frac{\partial}{\partial \beta} n^{-1} l_n(\beta, \theta) &= f_n \cdot \mathbb{P}_{n_1} h_{1,\beta}(Y, \mathbf{X}) + (1 - f_n) \cdot \mathbb{P}_{n_2} h_{2,\beta,\theta}(\tilde{Y}, \mathbf{X}) = \mathbf{0}, \\ \frac{\partial}{\partial \theta_1} n^{-1} l_n(\beta, \theta) &= f_n \cdot \mathbb{P}_{n_1} \left[\frac{(1 - Y)(\tilde{Y} - \theta_1)}{\theta_1(1 - \theta_1)} \right] + (1 - f_n) \cdot \mathbb{P}_{n_2} \left[\frac{\{1 - \psi(\mathbf{X}'\beta)\} \cdot \{\tilde{Y} - h_{3,\beta,\theta}(\mathbf{X})\}}{h_{3,\beta,\theta}(\mathbf{X}) \cdot \{1 - h_{3,\beta,\theta}(\mathbf{X})\}} \right] = 0, \\ \frac{\partial}{\partial \theta_2} n^{-1} l_n(\beta, \theta) &= f_n \cdot \mathbb{P}_{n_1} \left[\frac{Y(1 - \tilde{Y} - \theta_2)}{\theta_2(1 - \theta_2)} \right] - (1 - f_n) \cdot \mathbb{P}_{n_2} \left[\frac{\psi(\mathbf{X}'\beta) \cdot \{\tilde{Y} - h_{3,\beta,\theta}(\mathbf{X})\}}{h_{3,\beta,\theta}(\mathbf{X}) \cdot \{1 - h_{3,\beta,\theta}(\mathbf{X})\}} \right] = 0. \end{aligned} \right\} \quad (5.3)$$

Let us denote the joint maximum likelihood estimates as $(\hat{\beta}_{n,J}, \hat{\theta}_{n,J})$. Note that (2.1) has two components representing the contributions of the validation and non-validation data to the likelihood, respectively. If $\psi(\mathbf{x}'\beta_0) \in (0.1, 0.9)$, the non-identifiability problem discussed in the context of (5.2) will still be present, but plausibly to a lesser extent, especially when the validation sample size is small.

Further, note that the first and second components of the joint log-likelihood mentioned above, are sums of n_1 , and n_2 i.i.d. random variables, respectively. But the random variables in the two sums are not identically distributed. The first involves the random vector (Y, \mathbf{X}) , and the second, (\tilde{Y}, \mathbf{X}) . Since, $\min\{n_1, n_2\} \rightarrow \infty$, none of the two sums can be neglected asymptotically. Consequently, the log-likelihood cannot be written as a sum of i.i.d. random variables, even asymptotically, and hence, the classical asymptotic results on maximum likelihood estimation can not be used. There are other difficulties too, in dealing with the score functions obtained in (5.3). First, the restrictions on the parameter space of Θ (cf. assumption (A2)), and the non-compactness of the parameter space for β_0 , make the study of the joint uniform convergence of the score functions (over β and θ) difficult. Second, one has to show that, the limiting score functions have an unique and well-separated zero at (β_0, θ_0) . The techniques, used in proving the asymptotic results for PMLE, cannot be used for JMLE. Developing an asymptotic theory for the JMLE seems to be difficult.

5.1.3 A comparison between PMLE, CMLE, JMLE and naive logistic regression

In order to provide a clear motivation for using pseudo-likelihood, we design a simulation study. We choose, $p = 2$ with independent covariates $X_1 \sim N(0, 1)$ and $X_2 \sim N(0, \sigma^2)$, where σ^2 will be chosen later on. We fix, $\boldsymbol{\beta}_0 = (1, 2)$ and $\boldsymbol{\theta}_0 = (0.1, 0.3)$. Given any $\eta \in (0, 1)$, we choose $\sigma^2 = \sigma^2(\eta)$, such that

$$\eta = Q(\psi(\mathbf{X}'\boldsymbol{\beta}_0) \in (0.1, 0.9)) = Q(0.1 < \psi(X_1 + 2X_2) < 0.9) = Q(-\log 9 < X_1 + 2X_2 < \log 9).$$

This leads to,

$$\sigma^2(\eta) = \frac{1}{4} \cdot \left(\frac{(\log 9)^2}{(\Phi^{-1}(\frac{1+\eta}{2}))^2} - 1 \right), \quad (5.4)$$

where, Φ^{-1} denote the standard normal quantile function. In order to make sure that $\sigma^2(\eta) > 0$, we need to ensure that, $\eta \in (0, 0.97)$. By changing η , and consequently $\sigma^2(\eta)$, we could control the probability of $\psi(\mathbf{X}'\boldsymbol{\beta}_0)$ lying inside $(0.1, 0.9)$. Higher is the value of η , lower is the chance of getting extreme observations (i.e., observations whose $\psi(\mathbf{X}'\boldsymbol{\beta}_0)$ values are near 0 or 1) in a sample and more exacerbated would be the effect of non-identifiability. The primary goal of this study is to investigate the above phenomena.

In the simulation study, besides PMLE, CMLE and JMLE, we have also included naive logistic regression based estimator of $\boldsymbol{\beta}_0$, denoted by $\hat{\boldsymbol{\beta}}_{NL}$, which is obtained on the basis of $\{(\tilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq n\}$ ignoring the misclassification errors. The sample size for the simulation study is fixed at $n = 300$, with $n_1 = 60$ and $n_2 = 240$. Notice that, for computation of CMLE, and $\hat{\boldsymbol{\beta}}_{NL}$, we consider only the values of (\tilde{Y}, \mathbf{X}) for the entire sample.

In Table 1, we present the results of the simulation study. The simulation set-up is described at the top of the table. For simplicity of presentation, we report the average bias and mean-squared error (MSE) of different likelihood estimators of $(\beta_{1,0}, \theta_{1,0})$ using 250 simulated data-sets. For naive logistic regression it is for $\beta_{1,0}$ only.

Results from Table 1 clearly show that PMLE of $\beta_{1,0}$ has the best performance in terms of MSE. The same conclusion can be drawn from the frequency plots in Figures 2 and 3. For estimating $\theta_{1,0}$, PMLE is superior to others both in terms of bias and MSE. As expected, the CMLE estimators' performance is miserable in terms of both bias and MSE. In fact, often the score equations fail to yield stable solutions. Unless possibly the sample size is very large (Hausman et al. (1998) considered $n = 5000$), the non-identifiability problem persists. Moreover, the estimator $\hat{\theta}_{1,C}$ takes negative values and values greater than one (cf. Figure 1). Similar problem also arises, but to a lesser extent, for joint likelihood based estimator $\hat{\theta}_{1,J}$, as seen from the frequency plots in Figures 1. Evidently, the joint likelihood estimate of $\boldsymbol{\theta}$ could very well be nonsensical. However, PMLE does not suffer from this drawback and seems to perform better than the joint likelihood estimator at least for moderately large sample sizes which is often the case in many real-life applications. The naive logistic regression based estimator $\hat{\boldsymbol{\beta}}_{1,NL}$ suffers from serious underestimation and has larger MSE than PMLE or JMLE. Overall, PMLE based estimates seem to be the best choice among alternative likelihood based procedures.

Finally, in terms of computational speed, PMLE takes the least time to converge followed by the JMLE, which is, of course, considerably slower. The CMLE is the worst, and converges very slowly. As expected, often, the solutions converge to non-nonsensical values because of the non-identifiability problem.

Table 1: Average bias and (MSE) (in parenthesis below) comparison of estimates obtained using JMLE ($\hat{\beta}_{1,J}, \hat{\theta}_{1,J}$), PMLE ($\hat{\beta}_{1,n}, \hat{\theta}_{1,n}$), CMLE ($\hat{\beta}_{1,C}, \hat{\theta}_{1,C}$) and naive logistic regression $\hat{\beta}_{1,NL}$. Also shown are the different values of η and $\sigma(\eta) = \sqrt{\text{Var}(X_2)}$ (cf. (5.4)). Here, $n = 300$ and $n_1 = 60$. The true parameter values are $\beta_{1,0} = 1$ and $\theta_{1,0} = 0.1$.

		Bias and (MSE)						
		Estimates of $\beta_{1,0}$				Estimates of $\theta_{1,0}$		
η	$\sigma(\eta)$	$\hat{\beta}_{1,J}$	$\hat{\beta}_{1,n}$	$\hat{\beta}_{1,C}$	$\hat{\beta}_{1,NL}$	$\hat{\theta}_{1,J}$	$\hat{\theta}_{1,n}$	$\hat{\theta}_{1,C}$
0.6	1.21	-0.0152 (0.1829)	0.0761 (0.1487)	914.74 (1.53×10^8)	-0.6270 (0.4116)	-0.0415 (0.0135)	0.0094 (0.0026)	444.1 (5.35×10^{11})
0.7	0.935	-0.0162 (0.1682)	0.0767 (0.1361)	13.6 (5.74×10^6)	-0.584 (0.3604)	-0.0442 (0.0131)	0.0115 (0.0029)	3.05×10^4 (2.24×10^{11})
0.8	0.696	-0.0169 (0.1020)	0.0294 (0.0881)	474.0 (3×10^8)	-0.538 (0.306)	-0.0305 (0.01)	0.0127 (0.0026)	-3.30×10^4 (1.38×10^{12})
0.9	0.443	-0.0154 (0.102)	0.0178 (0.0842)	-1084.05 (2.47×10^8)	-0.5167 (0.285)	-0.0231 (0.0093)	0.0113 (0.0029)	1.26×10^5 (1.2×10^{12})

5.2 Inference: Asymptotic versus bootstrap

In this section, we design simulation studies to compare the performances of the penalized likelihood based bootstrap percentile intervals of the regression parameters, with the corresponding intervals based on the asymptotic distribution of the PMLE of β_0 . For evaluating the performances, we compute the empirical coverage and average length of each such interval. We consider two regression models for the simulation studies, one with three covariates and the other with nine covariates. The latter is considered to demonstrate that the PMLE based inference is computationally feasible even with a reasonable number of covariates. In the following we discuss the results.

5.2.1 Simulation study: Model with three covariates

We consider the following model (cf. Hausman et al. (1998)),

$$\mathbf{P}_0(Y = 1 | (X_2, X_3, X_4) = (x_2, x_3, x_4)) = \psi(\beta_{1,0} + \beta_{2,0}x_2 + \beta_{3,0}x_3 + \beta_{4,0}x_4), \quad \text{for all } (x_2, x_3, x_4). \quad (5.5)$$

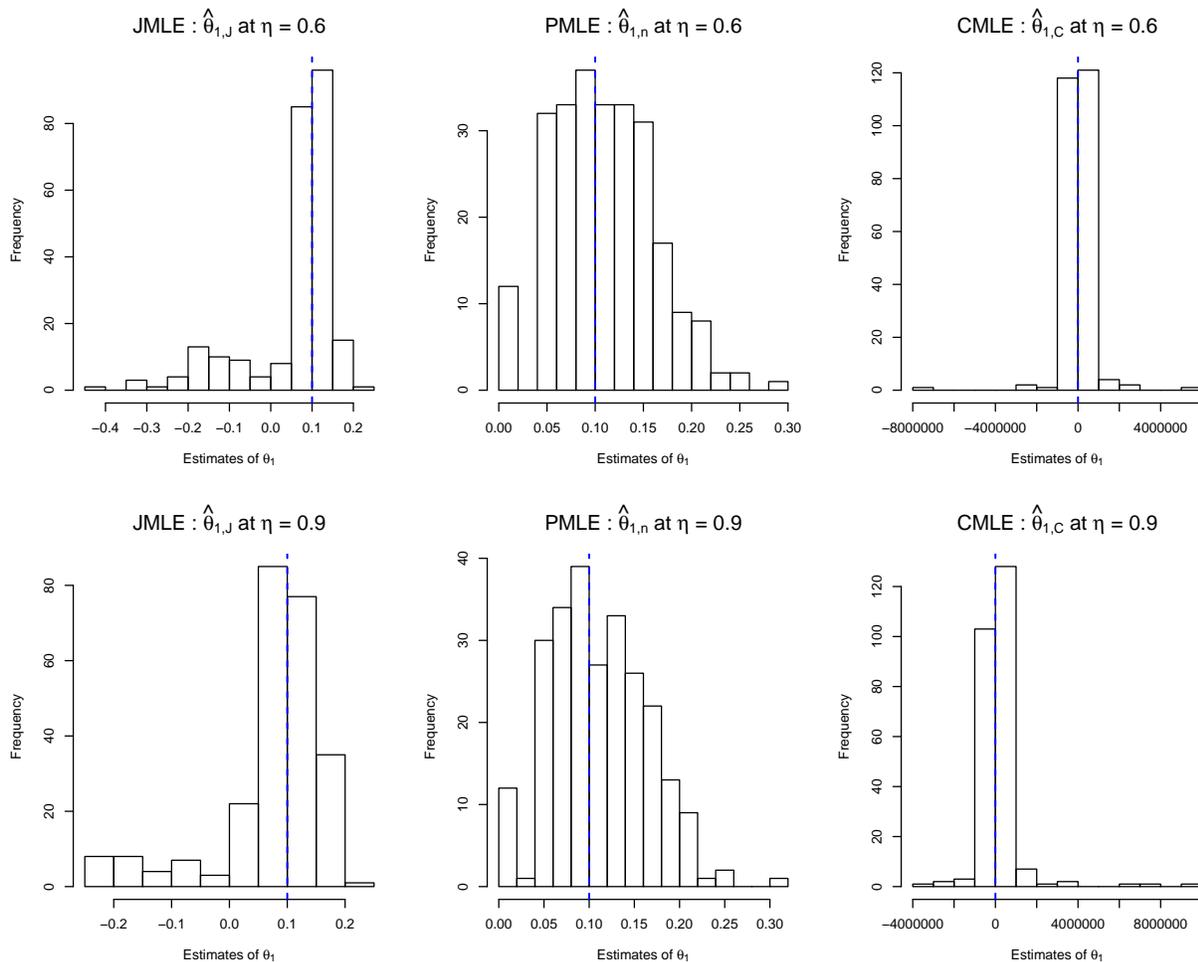


Figure 1: Frequency plots in the first and second rows show the distributions of $\hat{\theta}_{1,J}$, $\hat{\theta}_{1,n}$ and $\hat{\theta}_{1,C}$ at $\eta = 0.6$ and $\eta = 0.9$, respectively. Here, $n = 300$, $n_1 = 60$ and the frequency plots were based on 250 simulations. Blue vertical dotted line shows the true value $\theta_{1,0} = 0.1$. Both $\hat{\theta}_{1,J}$ and $\hat{\theta}_{1,C}$ can lead to negative estimates.

The covariates X_2 , X_3 and X_4 are independent, with $X_2 \sim \text{Log-normal}(0, 1) - e^{1/2}$, $X_3 \sim \text{Bernoulli}(1/3) - 1/3$ and $X_4 \sim \text{Uniform}(0, 1) - 1/2$. All covariates are centered. We study three different models:

(a) $\beta_0 = (\beta_{1,0}, \beta_{2,0}, \beta_{3,0}, \beta_{4,0}) = (0, 0.7, 1.5, -0.6)'$, $\theta_0 = (0.1, 0.3)'$.

(b) Same choice of β_0 as in model (a), with $\theta_0 = (0.1, 0.1)'$.

(c) $\beta_0 = (-1, 0.7, 1.5, -0.6)'$ and θ_0 is same as in model (a).

Models (a) and (c) differ only with respect to the intercept term, while models (a) and (b) with respect to the misclassification probabilities only.

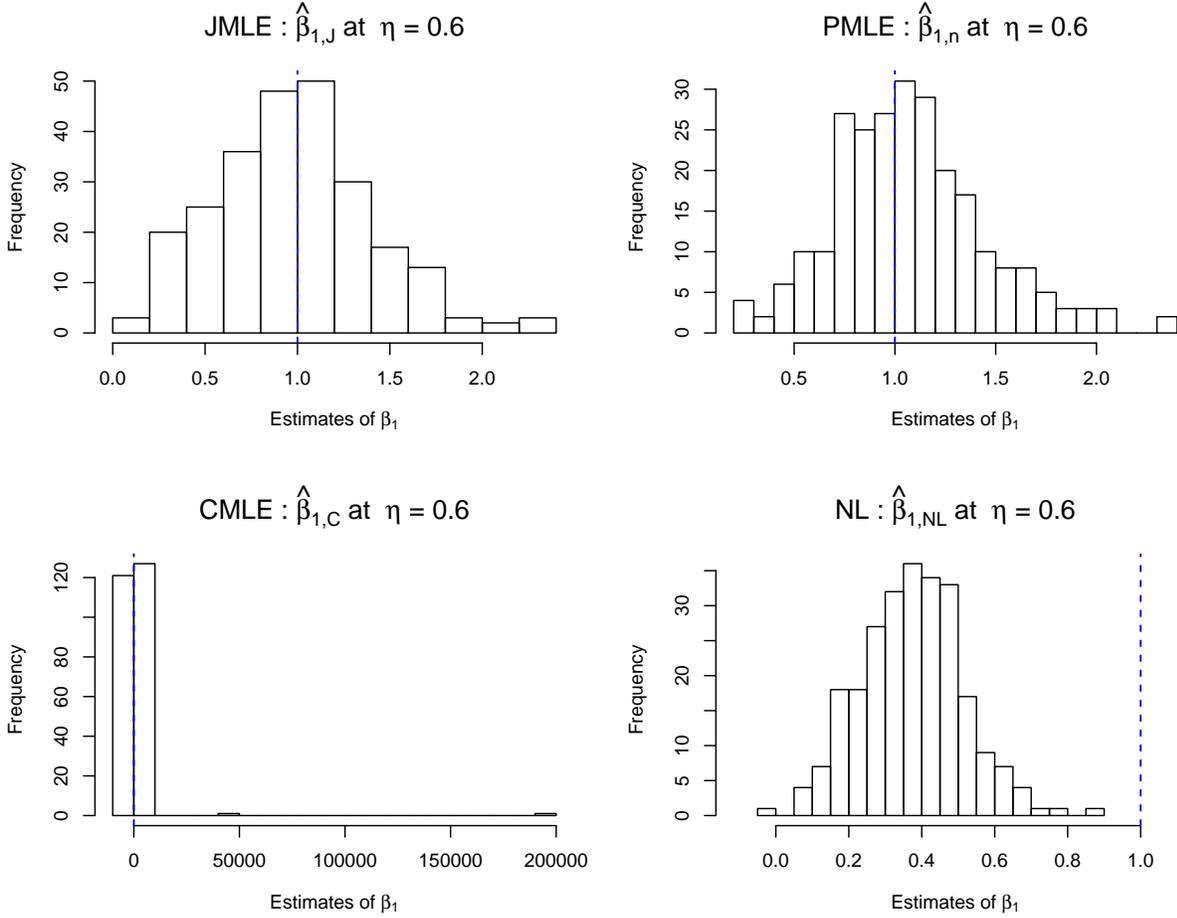


Figure 2: Frequency plots showing distributions of various estimators: JMLE $\hat{\beta}_{1,J}$, PMLE $\hat{\beta}_{1,n}$, CMLE $\hat{\beta}_{1,C}$ and NL $\hat{\beta}_{1,NL}$ at $\eta = 0.6$. Here, $n = 300$ and $n_1 = 60$. Blue vertical dotted line shows the true value $\beta_{1,0} = 1$. JMLE has slightly larger variability than PMLE, CMLE is the worst and NL has underestimation and higher variability.

For each of the three models considered above, we consider three different sample sizes $n \in \{300, 600, 1000\}$, and also three different validation sample fractions $f_n \in \{0.1, 0.2, 0.3\}$. We compute the empirical coverage probabilities and the average lengths of the bootstrap percentile interval, and the asymptotic confidence intervals for the parameters $\beta_{j,0}$, $j = 1, 2, 3, 4$, for each of these models. The bootstrap percentile intervals are obtained from (3.4). The asymptotic confidence intervals are obtained from (2.13) by plugging in the estimated values of β_0 and θ_0 , and replacing the expectation with respect to the unknown distribution $Q(\mathbf{x})$ of $\mathbf{X} = (X_2, X_3, X_4)$ by the same with respect to the empirical version of $Q(\mathbf{x})$ in the expression for the asymptotic covariance matrix. Finally, the limiting validation sampling fraction f is replaced by f_n . The empirical coverage probabilities, and the average lengths are computed on the basis of 250 simulated

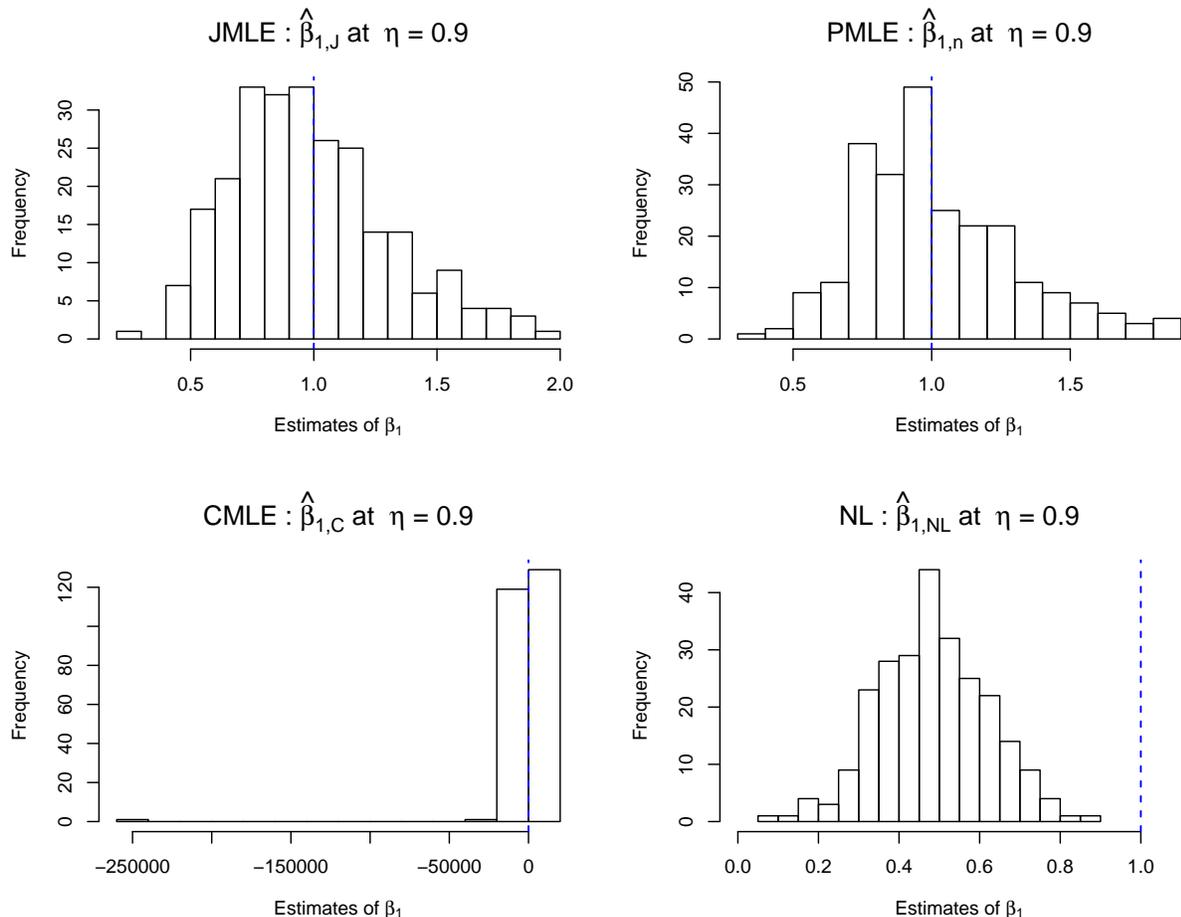


Figure 3: Frequency plots showing distributions of various estimators: JMLE $\hat{\beta}_{1,J}$, PMLE $\hat{\beta}_{1,n}$, CMLE $\hat{\beta}_{1,C}$ and NL $\hat{\beta}_{1,NL}$ at $\eta = 0.9$. Here, $n = 300$ and $n_1 = 60$. Blue vertical dotted line shows the true value $\beta_{1,0} = 1$. JMLE has slightly larger variability than PMLE, CMLE is the worst and NL has underestimation and higher variability.

data-sets. Tables 2, 3 and 4 show the results for $n = 300, 600$ and 1000 , respectively.

In Table 2, with $n = 300$, we notice that for $f_n = 0.2$ and 0.3 , the asymptotic and bootstrap methods have comparable empirical coverages, although the bootstrap intervals have slightly larger average length. The average length decreases as f_n increases. In case $f_n = 0.1$, the bootstrap based intervals have substantially larger average length, compared to the asymptotic intervals. This happens when the validation sample size n_1 is very small. The primary reason is, the appearance of a *bad* bootstrap validation sample. It leads to bad estimates of θ_0 and β_0 . A few such estimates influence the bootstrap percentiles, and hence lead to wider CI's obtained from (3.4).

Asymptotic CI's are also not exempt from such aberrations, if n_1 is very small. The lengths of asymptotic CI's for components of β_0 in model (c) with $n_1 = 30$, shown in Table 2 reveal this phenomenon, although to a lesser extent than that of the bootstrap CI's. Note, however, that the average lengths of asymptotic CI's are based on 250 simulated data-sets while the bootstrap based CI's are based on $B = 700$ bootstrap samples from each data-set. Therefore, if the original validation sample \mathcal{X}_{n_1} is bad, the chances of obtaining a worse bootstrap sample $\hat{\mathcal{X}}_{n_1}$ is higher. This magnifies the problem for the bootstrap case. In worst cases, obtaining even 50 or more extremely large estimates of a single component of $\hat{\beta}_n$ out of a total of 700 iterations can have a huge influence on the corresponding bootstrap 2.5% and 97.5% percentiles. However, this issue disappears if n_1 is slightly increased. Comparing the same figures for $n = 600$ and $n = 1000$ in Tables 3 and 4, we see a dramatic improvement over the earlier case even when $f_n = 0.1$.

Our simulations suggest that the phenomenon of unstable parameter estimates and large widths of bootstrap CI's is dependent on the validation sample size n_1 and not on the actual validation sampling fraction f_n . We find that $n_1 \approx 50$ or more to be a safe choice for obtaining reliable bootstrap based estimates, irrespective of the total sample size n . It could be achieved either by increasing n for a fixed f_n , or by increasing f_n for a fixed n . In fact, for $n = 300$ and $f = 0.2$ and $f = 0.3$, except for a few cases, the widths of bootstrap based CI's are only in between 5 - 20 % larger their asymptotic counterparts, and the situation improves further when n increases.

The good performance of the asymptotic CI's is not surprising, since the asymptotic covariance matrix has been estimated by using the exact expression in (2.13), and each term has been painstakingly computed. Also, the expectation with respect to the empirical distribution of the covariates is a good approximation to the true distribution $Q(\mathbf{x})$ for $n = 300$ or more. But, exact computation of the asymptotic covariance matrix may not be easy to implement in practice as is evident from the complicated expressions of the matrices involved in (2.14), and specially if p is large. Hence, for moderately large n_1 and n , clearly the bootstrap is a preferable method for inference on β_0 from the point of view of implementation.

5.3 Model with a large number of covariates

Here we consider a model with $p = 9$ covariates, including an intercept term with a mix of continuous, discrete and categorical covariates. The covariates $\{X_2, \dots, X_9\}$ are independent, with $X_2 \sim \text{Log-normal}(0, 1) - e^{1/2}$, $X_3 \sim \text{Bernoulli}(1/3) - 1/3$, $X_4 \sim \text{Uniform}(0, 1) - 1/2$, $(X_5, X_6) \sim N_2(0, 0, 1, 1, 0.6)$, $X_7 \sim \text{Poisson}(3) - 3$, $X_8 \sim \chi_2^2 - 2$ and $X_9 \sim (3/5) \cdot N(-1, 1) + (2/5) \cdot N(4, 2) - 1$. We choose, $\beta_0 = (-1, 0.7, 1.5, -0.6, 1, -0.75, -2, -1.5, 1)'$ with $n = 1000$. Table 5 presents the empirical coverages and average lengths. The observed patterns are similar to that of Tables 3 and 4.

Table 2: Empirical coverage probabilities and average lengths (in parenthesis below) for 95% confidence intervals for $\beta_{j,0}$, $j = 1, 2, 3, 4$, in models (a), (b) and (c), with sample size $n = 300$.

Model	Coefficient	$n_1 = 30$		$n_1 = 60$		$n_1 = 90$	
		Asymp.	Boot.	Asymp.	Boot.	Asymp.	Boot.
(a)	$\beta_{1,0}$	0.976 (2.235)	0.98 (15.529)	0.984 (1.389)	0.96 (1.427)	0.976 (1.06)	0.968 (1.037)
	$\beta_{2,0}$	0.972 (1.396)	0.956 (33.836)	0.944 (0.945)	0.936 (1.07)	0.96 (0.795)	0.944 (0.871)
	$\beta_{3,0}$	0.98 (2.953)	0.98 (44.891)	0.976 (2.021)	0.98 (2.381)	0.964 (1.708)	0.98 (1.811)
	$\beta_{4,0}$	0.996 (3.607)	0.996 (13.536)	0.984 (2.901)	0.996 (3.19)	0.976 (2.605)	0.988 (2.783)
(b)	$\beta_{1,0}$	0.964 (1.537)	0.992 (4.855)	0.968 (1.055)	0.992 (1.062)	0.964 (0.87)	0.98 (0.856)
	$\beta_{2,0}$	0.96 (1.025)	0.96 (10.015)	0.944 (0.759)	0.94 (0.854)	0.944 (0.672)	0.936 (0.739)
	$\beta_{3,0}$	0.984 (2.147)	0.984 (13.547)	0.976 (1.614)	0.988 (1.749)	0.948 (1.444)	0.972 (1.516)
	$\beta_{4,0}$	0.988 (2.734)	0.992 (10.466)	0.984 (2.389)	0.992 (2.563)	0.98 (2.244)	0.992 (2.357)
(c)	$\beta_{1,0}$	0.992 (74.772)	1 (561.488)	0.992 (2.359)	0.984 (4.62)	0.988 (1.61)	0.972 (1.557)
	$\beta_{2,0}$	0.964 (47.69)	0.908 (202.541)	0.968 (1.072)	0.952 (3.328)	0.972 (0.824)	0.948 (1.11)
	$\beta_{3,0}$	0.972 (218.372)	0.968 (629.814)	0.976 (2.69)	0.98 (4.968)	0.976 (2.084)	0.976 (2.582)
	$\beta_{4,0}$	0.996 (77.727)	0.996 (363.567)	0.968 (3.49)	0.992 (8.212)	0.96 (3.05)	0.988 (3.977)

5.4 Real data analysis

In this section we consider a data set obtained from a household literacy survey conducted across four Indian states, and apply our methodology for obtaining confidence intervals for the parameters of interest. For a detailed discussion about the survey we refer to [Kothari and Bandyopadhyay \(2011\)](#). One of the goals of the survey was to compare the literacy rate obtained from the data collected through indirect responses, which are subject to errors of misclassification, with the literacy rate obtained from the data collected through direct responses, which can be considered as gold-standard. Data on 7409 individuals within the age group 15-45 were collected. Literacy was judged directly by evaluating each individual through written/oral tests,

Table 3: Empirical coverage probabilities and average lengths (in parenthesis below) for 95% confidence intervals for $\beta_{j,0}$, $j = 1, 2, 3, 4$, in models (a), (b) and (c), with sample size $n = 600$.

Model	Coefficient	$n_1 = 60$		$n_1 = 120$		$n_1 = 180$	
		Asymp.	Boot.	Asymp.	Boot.	Asymp.	Boot.
(a)	$\beta_{1,0}$	0.964 (1.525)	0.956 (1.502)	0.952 (0.944)	0.956 (0.902)	0.956 (0.731)	0.952 (0.7)
	$\beta_{2,0}$	0.956 (0.906)	0.952 (0.935)	0.94 (0.643)	0.956 (0.67)	0.948 (0.549)	0.948 (0.567)
	$\beta_{3,0}$	0.992 (1.929)	0.984 (1.978)	0.964 (1.364)	0.964 (1.353)	0.96 (1.174)	0.968 (1.187)
	$\beta_{4,0}$	0.992 (2.335)	0.992 (2.506)	0.98 (1.958)	0.996 (2.035)	0.976 (1.789)	0.984 (1.841)
(b)	$\beta_{1,0}$	0.936 (1.008)	0.98 (1.231)	0.956 (0.72)	0.972 (0.702)	0.96 (0.608)	0.96 (0.586)
	$\beta_{2,0}$	0.948 (0.627)	0.948 (1.183)	0.932 (0.508)	0.936 (0.537)	0.92 (0.467)	0.924 (0.485)
	$\beta_{3,0}$	0.98 (1.349)	0.984 (2.098)	0.956 (1.09)	0.972 (1.113)	0.94 (1.003)	0.96 (1.025)
	$\beta_{4,0}$	0.968 (1.79)	0.988 (2.795)	0.968 (1.635)	0.976 (1.685)	0.972 (1.566)	0.98 (1.61)
(c)	$\beta_{1,0}$	0.996 (4.063)	0.976 (148.595)	0.996 (1.63)	0.96 (1.125)	0.988 (1.115)	0.92 (0.778)
	$\beta_{2,0}$	0.972 (2.557)	0.96 (104.463)	0.968 (0.719)	0.952 (0.705)	0.98 (0.564)	0.96 (0.542)
	$\beta_{3,0}$	0.98 (3.387)	0.944 (62.804)	0.984 (1.819)	0.98 (1.675)	0.992 (1.428)	0.98 (1.329)
	$\beta_{4,0}$	0.976 (3.79)	0.992 (24.522)	0.968 (2.284)	0.972 (2.505)	0.948 (2.058)	0.952 (2.125)

and the responses were, literate ($Y = 1$) or illiterate ($Y = 0$). The indirect responses were obtained from the head of the family, who reported on the literacy status (\tilde{Y}) of each member of the family, with $\tilde{Y} = 1$ if the member was considered literate and $\tilde{Y} = 0$, otherwise. The latter method is used in Indian census for collecting literacy data. Since the chances of misclassifying a literate person ($Y = 1$) as an illiterate by the head of the family ($\tilde{Y} = 0$) is zero, this type of misclassification is ignored. For each individual, his/her age X was also recorded and we treat it as a covariate.

Complete information on the triplet, (Y, \tilde{Y}, X) was available for 7409 individuals and based on this, we found $\theta_{1,0} = \mathbf{P}(\tilde{Y} = 1 | Y = 0) = 0.177$. For our data-analysis, we selected a random sample of size $n = 740$

Table 4: Empirical coverage probabilities and average lengths (in parenthesis below) for 95% confidence intervals for $\beta_{j,0}$, $j = 1, 2, 3, 4$, in models (a), (b) and (c), with sample size $n = 1000$.

Model	Coefficient	$n_1 = 100$		$n_1 = 200$		$n_1 = 300$	
		Asymp.	Boot.	Asymp.	Boot.	Asymp.	Boot.
(a)	$\beta_{1,0}$	0.964 (1.103)	0.956 (1.084)	0.956 (0.724)	0.964 (0.683)	0.976 (0.563)	0.96 (0.532)
	$\beta_{2,0}$	0.956 (0.641)	0.964 (0.632)	0.952 (0.484)	0.932 (0.484)	0.944 (0.418)	0.932 (0.42)
	$\beta_{3,0}$	0.984 (1.385)	0.972 (1.293)	0.952 (1.041)	0.952 (1.009)	0.948 (0.901)	0.944 (0.893)
	$\beta_{4,0}$	0.98 (1.681)	0.98 (1.74)	0.984 (1.488)	0.98 (1.519)	0.98 (1.372)	0.988 (1.402)
(b)	$\beta_{1,0}$	0.948 (0.768)	0.98 (0.774)	0.968 (0.554)	0.972 (0.538)	0.964 (0.465)	0.964 (0.446)
	$\beta_{2,0}$	0.968 (0.463)	0.984 (0.488)	0.952 (0.387)	0.96 (0.402)	0.932 (0.355)	0.948 (0.365)
	$\beta_{3,0}$	0.976 (1.005)	0.968 (1.008)	0.972 (0.832)	0.98 (0.834)	0.964 (0.765)	0.968 (0.774)
	$\beta_{4,0}$	0.964 (1.341)	0.972 (1.381)	0.964 (1.251)	0.964 (1.272)	0.964 (1.197)	0.968 (1.217)
(c)	$\beta_{1,0}$	0.988 (2.373)	0.944 (3.126)	0.988 (1.282)	0.928 (0.769)	1 (0.869)	0.932 (0.594)
	$\beta_{2,0}$	0.976 (0.921)	0.928 (2.42)	0.98 (0.549)	0.956 (0.475)	0.976 (0.427)	0.968 (0.401)
	$\beta_{3,0}$	0.976 (2.371)	0.924 (4.048)	0.988 (1.408)	0.964 (1.146)	0.98 (1.095)	0.972 (0.995)
	$\beta_{4,0}$	0.964 (2.143)	0.968 (5.082)	0.964 (1.746)	0.96 (1.767)	0.96 (1.578)	0.952 (1.605)

from the complete set (about 10% of the available number of observations), out of which the first $n_1 = 148$ observations were treated as the validation sample. This leads to a validation sample fraction $f \approx 0.2$. After including an intercept term in the logistic regression model, we have

$$\mathbf{P}_0(Y = 1 | X = x) = \frac{1}{1 + \exp\{-\beta_{1,0} - \beta_{2,0} \cdot x\}}, \quad x \in \mathbb{R}.$$

Based on complete information over 7409 individuals, the *true* parameter values were found to be, $(\beta_{1,0}, \beta_{2,0}) = (1.66, -0.0623)$. However, usually such detailed information would be unavailable. Using the validation sample, we found $\hat{\theta}_{1,n} = 0.137$. The pseudo-likelihood based estimates using the sample of size $n = 740$ are, $\hat{\beta}_{1,n} = 2.37$ and $\hat{\beta}_{2,n} = -0.0761$, and the corresponding bootstrap percentile 95% confidence intervals are,

Table 5: Empirical coverage probabilities and average lengths (in parenthesis below) for 95% confidence intervals for $\beta_{j,0}$, $j = 1, \dots, 9$, with sample size $n = 1000$.

		Covariates ^a and associated coefficients								
Method	n_1	Intercept	Log-Normal	Bernoulli	Uniform	BVN.1 ^b	BVN.2 ^b	Poisson	χ_2^2	Mixture ^c
		$\beta_{1,0}$	$\beta_{2,0}$	$\beta_{3,0}$	$\beta_{4,0}$	$\beta_{5,0}$	$\beta_{6,0}$	$\beta_{7,0}$	$\beta_{8,0}$	$\beta_{9,0}$
Asymp.	100	0.98 (3.407)	0.948 (1.718)	0.952 (3.528)	0.972 (5.14)	0.948 (2.514)	0.952 (2.424)	0.956 (4.431)	0.96 (3.38)	0.96 (2.167)
	200	0.988 (1.456)	0.956 (0.664)	0.976 (1.87)	0.948 (2.682)	0.96 (1.164)	0.972 (1.084)	0.948 (1.443)	0.956 (1.153)	0.968 (0.734)
	300	0.984 (1.095)	0.948 (0.559)	0.96 (1.644)	0.948 (2.437)	0.96 (1.011)	0.952 (0.959)	0.96 (1.129)	0.944 (0.919)	0.956 (0.578)
Boot.	100	0.976 (7.641)	0.932 (5.391)	0.968 (10.58)	0.992 (9.795)	0.944 (7.143)	0.956 (6.394)	0.944 (15.206)	0.948 (11.569)	0.944 (7.536)
	200	0.968 (1.448)	0.96 (0.888)	0.988 (2.335)	0.988 (3.316)	0.972 (1.467)	0.976 (1.362)	0.936 (1.949)	0.952 (1.547)	0.968 (0.989)
	300	0.972 (1.099)	0.956 (0.678)	0.984 (1.902)	0.988 (2.805)	0.972 (1.181)	0.972 (1.114)	0.916 (1.379)	0.928 (1.114)	0.952 (0.699)

^a All covariates are centered.

^b BVN.1 and BVN.2 denotes covariates X_5 and X_6 respectively, where $(X_5, X_6) \sim N_2(0, 0, 1, 1, 0.65)$.

^c Corresponds to $X_9 \sim (3/5) \cdot N(-1, 1) + (2/5) \cdot N(4, 2) - 1$.

(1.752, 2.961) and $(-0.09573, -0.05529)$ respectively. The CI for the intercept coefficient does not contain $\beta_{1,0}$, but the CI for the slope parameter $\beta_{2,0}$ contains the true value. Also, the CI for $\beta_{2,0}$ indicates that, plausibly with increase in age, the chances of a person being literate decreases, which is true in general for the Indian populace.

6 Concluding remarks

In this article, we propose a pseudo likelihood based approach to the estimation of parameters of a regression model that assumes a logistic link and incorporates misclassification probabilities of binary responses as parameters. Under a minimal set of assumptions, the consistency and asymptotic normality of the resulting estimators are proved. To the best of our knowledge, this is possibly the only paper in the existing literature, that provides a rigorous asymptotic theory for likelihood based inference for such models. Also, the asymptotic theory developed here is comprehensive enough to deal with differential classification errors.

For drawing inference on the model parameters we propose a bootstrap method supported by its distributional consistency. The method avoids direct estimation of asymptotic variance of the estimators. This helps in making the implementation of the methodology easier for practitioners in real-life situations. The bootstrap distributional consistency result is new to the literature of pseudo-likelihood based estimation. No result is available to handle bootstrap consistency in this set up. Some of the techniques used in proving the bootstrap consistency result are novel, and could possibly be used to prove bootstrap consistency in similar problems.

The extensive numerical studies presented here clearly show the superiority of the proposed pseudo-likelihood based estimation procedure over the other commonly used likelihood based methods for the estimation of the model parameters. In this context, it should be mentioned that no asymptotic theory is available for the joint maximum likelihood estimators (JMLE), and developing such a theory, at this point, seems to be a formidable problem.

Finally, it should be noted that the proposed methodology is applicable for binary regression with a logistic link only. It is an open problem to develop a similar theory for binary regression with arbitrary, but known link functions.

7 Proofs of main results

Throughout the proofs many convergence in probability statements involving the p -dimensional functions $h_{1,\beta}$ or $h_{2,\beta,\theta}$ (cf. (2.3)) have been proved by showing that those convergence in probability statements hold for each component of $h_{1,\beta}$ or $h_{2,\beta,\theta}$. This is possible since p is fixed. Without loss of generality, we have studied the first component of these p dimensional functions and for simplicity, we have used the same notation to denote the first components of $h_{1,\beta}$ or $h_{2,\beta,\theta}$. On some occasions, the same notational convention has been used for the p -dimensional functions, $Z_n(\beta)$, $Z_{n,i}(\beta)$, $Z(\beta)$ and $Z_i(\beta)$, (cf. (2.6) and (2.8)) for $i = 1, 2$. Whenever this has been done while proving a result, we have noted that in the proof. The cases where these functions are treated as p -dimensional functions will be obvious from the context of the concerned statements. In case of the with-intercept model (2.9), under the modified assumption (A3[†]), the same results can be proved by studying the second or other components of $h_{1,\beta}$ or $h_{2,\beta,\theta}$, by following exactly same arguments as in the without intercept case, as shown in the proofs below. Hence, separate proofs for the with-intercept case are not shown.

7.1 Proofs of Theorems 2.2 and 3.2

Proof of Theorem 2.2(i). Consider the first components of $h_{1,\beta}$, $h_{2,\beta,\theta}$, $Z_n(\beta)$, $Z_{n,i}(\beta)$, $Z(\beta)$ and $Z_i(\beta)$, for $i = 1, 2$, and denote them by the same symbols. Using assumptions (A3) and (A4) along with Lemma 8.7

and (8.1) we can write,

$$\begin{aligned}
& \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} |Z_n(\boldsymbol{\beta}) - Z(\boldsymbol{\beta})| \\
& \leq f_n \cdot \sup_{\boldsymbol{\beta}} |Z_{n,1}(\boldsymbol{\beta}) - Z_1(\boldsymbol{\beta})| + (1 - f_n) \cdot \sup_{\boldsymbol{\beta}} |Z_{n,2}(\boldsymbol{\beta}) - Z_2(\boldsymbol{\beta})| + |f_n - f| \cdot \sup_{\boldsymbol{\beta}} |\mathbf{P}_0 h_{1,\boldsymbol{\beta}}| \\
& \quad + |f_n - f| \cdot \sup_{\boldsymbol{\beta}} |\mathbf{P}_0 h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0}| \\
& \leq o_{\mathbf{P}_0^*}(1) + o_{\mathbf{P}_0^*}(1) + |f_n - f| \cdot 2\mathbf{E}|X_1| + |f_n - f| \cdot M_0 \cdot \mathbf{E}|X_1| = o_{\mathbf{P}_0^*}(1) + o(1) = o_{\mathbf{P}_0^*}(1).
\end{aligned}$$

Extending this argument to all p components we obtain, $\sup_{\boldsymbol{\beta}} \|Z_n(\boldsymbol{\beta}) - Z(\boldsymbol{\beta})\| = o_{\mathbf{P}_0^*}(1)$. From (2.7) we have, $Z_n(\hat{\boldsymbol{\beta}}_n) = o_{\mathbf{P}_0^*}(1)$. Combining these facts with Lemma 8.3 and using Theorem 2.10 of Kosorok (2008) leads to the proof. \square

Proof of Theorem 2.2(ii). We verify the conditions of Theorem 2.11 of Kosorok (2008) to complete the proof. Note that, from Lemma 8.3 we have $Z(\boldsymbol{\beta}_0) = 0$, from (2.7) it follows that $\sqrt{n}Z_n(\hat{\boldsymbol{\beta}}_n) = o_{\mathbf{P}_0^*}(1)$ and from Theorem 2.2(i) it follows that $\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\| = o_{\mathbf{P}_0^*}(1)$. From Lemma 8.10 we obtain the stochastic equicontinuity condition, $\|\sqrt{n}(Z_n - Z)(\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0)\| = o_{\mathbf{P}_0^*}(1 + \sqrt{n}\|\hat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\|)$. Also, as per assumption (A5), $\dot{Z}(\boldsymbol{\beta}_0)$ (cf. (8.41)) is nonsingular. It remains to study the limit distribution of $\sqrt{n}(Z_n(\boldsymbol{\beta}_0) - Z(\boldsymbol{\beta}_0))$. Using the definition of $h_{1,\boldsymbol{\beta}}$ (cf. (2.3)) and \mathbf{W}_i 's (cf. (7.6)), we define the $(p+3)$ dimensional i.i.d. random vectors \mathbf{T}_i 's as,

$$\mathbf{T}_i = \begin{pmatrix} \mathbf{T}_i^{(1)} \\ \mathbf{T}_i^{(2)} \end{pmatrix}, \quad \text{with} \quad \left. \begin{array}{l} \mathbf{T}_i^{(1)} = h_{1,\boldsymbol{\beta}_0}(Y_i, \mathbf{X}_i) - \mathbf{E}_0 h_{1,\boldsymbol{\beta}_0}(Y_1, \mathbf{X}_1), \\ \mathbf{T}_i^{(2)} = \mathbf{W}_i - \mathbf{E}_0 \mathbf{W}_1, \end{array} \right\} \quad i = 1, \dots, n_1. \quad (7.1)$$

Note that, $\mathbf{E}_0 \mathbf{T}_i = \mathbf{0}$ and $\mathbf{Var}_0(\mathbf{T}_i) = \boldsymbol{\Sigma}$, which is computed in (8.37). Using assumption (A3) and the CLT for i.i.d. random vectors, $\sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1} \xrightarrow{d} \mathbf{T} = (\mathbf{T}^{(1)}, \mathbf{T}^{(2)})' \sim N_{p+3}(\mathbf{0}, \boldsymbol{\Sigma})$. Based on the definition of the \mathbf{T}_i 's we can write,

$$\begin{aligned}
& \sqrt{n} \cdot (Z_n(\boldsymbol{\beta}_0) - Z(\boldsymbol{\beta}_0)) \\
& = \sqrt{f_n} \cdot \mathbf{G}_{n_1} h_{1,\boldsymbol{\beta}_0} + \sqrt{1-f_n} \cdot \mathbf{G}_{n_2} (h_{2,\boldsymbol{\beta}_0,\hat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}) + \sqrt{1-f_n} \cdot \mathbf{G}_{n_2} h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \\
& \quad + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot (\mathbf{P}_0 h_{2,\boldsymbol{\beta}_0,\hat{\boldsymbol{\theta}}_n} - \mathbf{P}_0 h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}) + o(1) \\
& \equiv A_{1,n} + A_{2,n} + A_{3,n} + A_{4,n} + o(1) \\
& = \sqrt{f_n} \cdot \sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1}^{(1)} + A_{2,n} + A_{3,n} + A_{4,n} + o(1).
\end{aligned} \quad (7.2)$$

Note that $A_{3,n}$ is independent of $A_{1,n}$ and $A_{4,n}$. Consider the function $\bar{f} : \Theta \mapsto \mathbf{P}_0 h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}}$, defined in Lemma 8.8(ii) and its total derivative map \mathbf{A}_0 evaluated at $\boldsymbol{\theta} = \boldsymbol{\theta}_0$, shown in (8.40). Note that \mathbf{A}_0 is a $p \times 2$ matrix and recall the definitions of $\boldsymbol{\phi}$ (cf. (7.8)) and \mathbf{B}_0 (cf. (2.10)). We can use Lemma 2.1 and the Delta method

(cf. Theorem 3.1 of [van der Vaart \(1998\)](#)) repeatedly to carry out the following simplification:

$$\begin{aligned}
A_{4,n} &= \sqrt{1-f_n} \cdot \sqrt{n_2} \left[\mathbf{P}_0 h_{2,\beta_0, \hat{\theta}_n} - \mathbf{P}_0 h_{2,\beta_0, \theta_0} \right] = \sqrt{1-f_n} \cdot \sqrt{n_2} \left[\bar{f}(\hat{\theta}_n) - \bar{f}(\theta_0) \right] \\
&= \sqrt{1-f_n} \cdot \sqrt{n_2} \left[\mathbf{A}_0(\hat{\theta}_n - \theta_0) + o_{\mathbf{P}_0^*}(\|\hat{\theta}_n - \theta_0\|) \right] \\
&= \mathbf{A}_0 \cdot \frac{1-f_n}{\sqrt{f_n}} \cdot \sqrt{n_1} (\phi(\bar{\mathbf{W}}_{n_1}) - \phi(\mathbf{E}_0 \mathbf{W}_1)) + o_{\mathbf{P}_0^*}(1) \\
&= \frac{1-f_n}{\sqrt{f_n}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1}^{(2)} + o_{\mathbf{P}_0^*}(1), \quad (\text{here, } \mathbf{A}_0 \mathbf{B}_0 \text{ is a } p \times 3 \text{ matrix}).
\end{aligned}$$

Let \mathbf{I}_p denote the p dimensional identity matrix. Following the above steps we have,

$$\begin{aligned}
A_{1,n} + A_{4,n} &= \left(\sqrt{f_n} \cdot \mathbf{I}_p \quad \frac{1-f_n}{\sqrt{f_n}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \right) \begin{pmatrix} \sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1}^{(1)} \\ \sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1}^{(2)} \end{pmatrix} + o_{\mathbf{P}_0^*}(1) \\
&\xrightarrow{d} \left(\sqrt{f} \cdot \mathbf{I}_p \quad \frac{1-f}{\sqrt{f}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \right) \begin{pmatrix} \mathbf{T}^{(1)} \\ \mathbf{T}^{(2)} \end{pmatrix}.
\end{aligned}$$

Using assumption (A3) and the CLT for i.i.d. random vectors,

$$\begin{aligned}
A_{3,n} &= \sqrt{1-f_n} \cdot \mathbb{G}_{n_2} h_{2,\beta_0, \theta_0} \\
&= \sqrt{1-f_n} \cdot \sqrt{n_2} \left(\frac{1}{n_2} \sum_{i=n_1+1}^n h_{2,\beta_0, \theta_0}(\tilde{Y}_i, \mathbf{X}_i) - \mathbf{E}_0 h_{2,\beta_0, \theta_0}(\tilde{Y}_i, \mathbf{X}_i) \right) \xrightarrow{d} \sqrt{1-f} \cdot N_p(\mathbf{0}, \mathbf{\Gamma}),
\end{aligned}$$

with $\mathbf{\Gamma}$ being defined as in (8.38). Further, using Lemma 8.6,

$$\|A_{2,n}\| \leq \sqrt{1-f_n} \cdot \sup_{\beta} \|\mathbb{G}_{n_2}(h_{2,\beta, \hat{\theta}_n} - h_{2,\beta, \theta_0})\| = o_{\mathbf{P}_0^*}(1).$$

Since \mathbf{T} and the $N_p(\mathbf{0}, \mathbf{\Gamma})$ random vectors can be considered as independent, it follows that

$$\sqrt{n}(Z_n(\beta_0) - Z(\beta_0)) \xrightarrow{d} \mathbf{U} \equiv \sqrt{f} \cdot \mathbf{T}^{(1)} + \frac{1-f}{\sqrt{f}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \cdot \mathbf{T}^{(2)} + \sqrt{1-f} \cdot N_p(\mathbf{0}, \mathbf{\Gamma}), \quad (7.3)$$

where, $\mathbf{U} \sim N_p(\mathbf{0}, \mathbf{\Sigma}_0)$ with $\mathbf{\Sigma}_0$ given by (2.14). This completes the verification of required conditions in Theorem 2.11 of [Kosorok \(2008\)](#), and it implies that,

$$\sqrt{n}(\hat{\beta}_n - \beta_0) \xrightarrow{d} - \left(\dot{Z}(\beta_0) \right)^{-1} \mathbf{U} \stackrel{d}{=} N_p \left(\mathbf{0}, \left[\dot{Z}(\beta_0) \right]^{-1} \mathbf{\Sigma}_0 \left(\left[\dot{Z}(\beta_0) \right]^{-1} \right)' \right),$$

as stated in (2.13). □

Proof of Theorem 3.2(i). Let us consider the first components of $h_{1,\beta}$, $h_{2,\beta,\theta}$, $Z_n(\beta)$, $Z_{n,i}(\beta)$, $Z(\beta)$ and $Z_i(\beta)$, for $i = 1, 2$, and denote them by the same symbols. Define the classes of functions, $\mathcal{G}_1 = \{h_{1,\beta} : \beta \in \mathbb{R}^p\}$

and $\mathcal{G}_2 = \{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$. Following the arguments used in the proof of Theorem 2.2, we have

$$\begin{aligned}
& \sup_{\beta} |\widehat{Z}_n(\beta) - Z(\beta)| \\
& \leq \sup_{\beta} |f_n \cdot \widehat{\mathbb{P}}_{n_1} h_{1,\beta} + (1 - f_n) \cdot \widehat{\mathbb{P}}_{n_2} h_{2,\beta,\widehat{\theta}_n} - f_n \cdot \mathbf{P}_0 h_{1,\beta} - (1 - f_n) \cdot \mathbf{P}_0 h_{2,\beta,\theta_0}| \\
& \quad + |f_n - f| \cdot \sup_{\beta} |\mathbf{P}_0 h_{1,\beta}| + |f_n - f| \cdot \sup_{\beta} |\mathbf{P}_0 h_{2,\beta,\theta_0}| \\
& \leq f_n \cdot \sup_{\beta} |(\widehat{\mathbb{P}}_{n_1} - \mathbf{P}_0) h_{1,\beta}| + (1 - f_n) \cdot \sup_{\beta} \left| \widehat{\mathbb{P}}_{n_2} \left(h_{2,\beta,\widehat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| \\
& \quad + (1 - f_n) \cdot \sup_{\beta} |(\widehat{\mathbb{P}}_{n_2} - \mathbf{P}_0) h_{2,\beta,\theta_0}| + o(1) \\
& \equiv E_{1,n} + E_{2,n} + E_{3,n} + o(1). \tag{7.4}
\end{aligned}$$

Consider $E_{1,n}$ in (7.4). We can express, $\widehat{\mathbb{P}}_{n_1} - \mathbf{P}_0 = n_1^{-1} \sum_{i=1}^{n_1} M_{n_1,i} \cdot (\delta_{(Y_i, \mathbf{X}_i)} - \mathbf{P}_0)$. Note that, the (scaled) multinomial weights $\{n_1^{-1} M_{n_1,i} : 1 \leq i \leq n_1\}$ are non-negative and exchangeable random variables and independent of \mathcal{X}_{n_1} . Further, $n_1^{-1} \sum_{i=1}^{n_1} M_{n_1,i} = 1$ and $\max\{|n_1^{-1} M_{n_1,i}| : 1 \leq i \leq n_1\} = o_{\mathbf{P}_M}(1)$, due to the following argument: for each $i \in \{1, \dots, n_1\}$, $M_{n_1,i} \sim \text{Bin}(n_1, n_1^{-1}) \stackrel{d}{=} \sum_{i=1}^{n_1} U_{n_1,i}$, where, $\{U_{n_1,i} : 1 \leq i \leq n_1\}$ are i.i.d. Bernoulli (n_1^{-1}) . Fix any $\epsilon \in (0, 1]$. Then, there exists some $n_1 \in \mathbb{N}$, such that $n_1^{-1} < \min\{\epsilon, 1/2\}$, which implies, $0 < \epsilon - \mathbf{E}(U_{n_1,1}) = \epsilon - n_1^{-1} < 1 - n_1^{-1}$. We can apply Hoeffding's inequality (Proposition A.6.1 of [van der Vaart and Wellner \(1996\)](#)), to obtain

$$\begin{aligned}
\mathbf{P}_M \left(\max_{1 \leq i \leq n_1} \frac{|M_{n_1,i}|}{n_1} \geq \epsilon \right) & \leq n_1 \cdot \mathbf{P}_M (M_{n_1,1} \geq n_1 \epsilon) \\
& = n_1 \cdot \mathbf{P} (\bar{U}_{n_1} - \mathbf{E}(U_{n_1,1}) \geq \epsilon - \mathbf{E}(U_{n_1,1})) \\
& \leq n_1 \cdot \exp \left\{ - \left(n_1 \epsilon^2 + \frac{1}{n_1} - 2\epsilon \right) \cdot \frac{n_1}{n_1 - 2} \cdot \log(n_1 - 1) \right\} \\
& \approx n_1 \cdot \exp \left\{ (-n_1 \epsilon^2 + 2\epsilon) \cdot \log n_1 \right\}, \quad (\text{for large enough } n_1) \\
& = n_1^{1+2\epsilon-n_1 \epsilon^2} \rightarrow 0, \quad \text{as } n_1 \rightarrow \infty.
\end{aligned}$$

Since, \mathcal{G}_1 is a Glivenko Cantelli (GC) class (from Lemma 8.7(i)), using Lemma 3.6.16 of [van der Vaart and Wellner \(1996\)](#) it follows that $E_{1,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Following the same arguments we can show that $E_{3,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, since \mathcal{G}_2 is a GC class (from Lemma 8.7(ii)) and the weights $\{n_2^{-1} M_{n_2,i} : 1 \leq i \leq n_2\}$ satisfy the required assumptions.

Now consider $E_{2,n}$. From Lemma 3.1(ii) and Lemma 8.11 we know that, $\|\widehat{\theta}_n - \theta_0\| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. This implies (using Lemma 8.18), there exists a sequence $\epsilon_n \downarrow 0$, such that $\mathbf{P}_M^*(\|\widehat{\theta}_n - \theta_0\| > \epsilon_n) = o_{\mathbf{P}_0^*}(1)$. Consider the events, $A_n = [\|\widehat{\theta}_n - \theta_0\| \leq \epsilon_n]$, $n \geq 1$. Since $\epsilon_n \downarrow 0$, there exists some $n_0 \in \mathbb{N}$, such that, $\{\theta : \|\theta - \theta_0\| \leq \epsilon_n\} \subset \Theta$ for all $n \geq n_0$. On the set A_n , using (8.6) we obtain,

$$\sup_{\beta} \left| \widehat{\mathbb{P}}_{n_2} (h_{2,\beta,\widehat{\theta}_n} - h_{2,\beta,\theta_0}) \right| \leq K_0 \cdot \epsilon_n \cdot \left(\frac{1}{n_2} \sum_{i=n_1+1}^n |\widehat{X}_{i,1}| - \mathbf{E}|X_1| \right) + K_0 \cdot \epsilon_n \cdot \mathbf{E}|X_1|.$$

Fix any $\eta > 0$. Since $\epsilon \downarrow 0$, assumption (A3) implies that for large enough n we must have, $2K_0 \cdot \epsilon_n \cdot \mathbf{E}|X_1| < \eta$. Also, $(n_2^{-1} \sum_i |\hat{X}_{1,i}| - \mathbf{E}|X_1|) = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* probability, under assumption (A3). Hence,

$$\begin{aligned} & \mathbf{P}_M^* \left(\sup_{\beta} |\hat{\mathbb{P}}_{n_2}(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0})| > \eta \right) \\ & \leq \mathbf{P}_M^* \left(K_0 \cdot \epsilon_n \cdot \left(\frac{1}{n_2} \sum_{i=n_1+1}^n |\hat{X}_{i,1}| - \mathbf{E}|X_1| \right) + K_0 \cdot \epsilon_n \cdot \mathbf{E}|X_1| > \eta, \|\hat{\theta}_n - \theta_0\| \leq \epsilon_n \right) + o_{\mathbf{P}_0^*}(1) \\ & \leq \mathbf{P}_M^* \left(n_2^{-1} \sum_{i=n_1+1}^n |\hat{X}_{i,1}| - \mathbf{E}|X_1| > \frac{\eta}{2\epsilon_n \cdot K_0} \right) + \mathbf{P}_M^* \left(K_0 \cdot \epsilon_n \cdot \mathbf{E}|X_1| > \frac{\eta}{2} \right) + o_{\mathbf{P}_0^*}(1) \\ & = o_{\mathbf{P}_0^*}(1) + 0 = o_{\mathbf{P}_0^*}(1). \end{aligned}$$

Extending this argument to all components we get, $\sup\{\|Z_n(\beta) - Z(\beta)\| : \beta \in \mathbb{R}^p\} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. The proof follows by using Theorem 13.1 of Kosorok (2008) along with Lemma 8.3. \square

Proof of Theorem 3.2(ii). Following the definition of \mathbf{T}_i 's in (7.1), define their bootstrap versions and denote them by $\{\hat{\mathbf{T}}_i : 1 \leq i \leq n_1\}$, where $\hat{\mathbf{T}}_i = (\hat{\mathbf{T}}_i^{(1)}, \hat{\mathbf{T}}_i^{(2)})'$. Following same arguments as in the proof of Theorem 2.2(ii), using Lemma 3.1(ii) and Lemma 8.11, we can write

$$\sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \mathbf{P}_0 \left(h_{2,\beta_0,\hat{\theta}_n} - h_{2,\beta_0,\theta_0} \right) = \frac{1-f_n}{\sqrt{f_n}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \cdot \sqrt{n_1} \cdot \bar{\mathbf{T}}_{n_1}^{(2)} + o_{\mathbf{P}_M^*}(1).$$

This implies,

$$\sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \mathbf{P}_0 \left(h_{2,\beta_0,\hat{\theta}_n} - h_{2,\beta_0,\hat{\theta}_n} \right) = \frac{1-f_n}{\sqrt{f_n}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \cdot \sqrt{n_1} \cdot \left(\bar{\mathbf{T}}_{n_1}^{(2)} - \bar{\mathbf{T}}_{n_1}^{(2)} \right) + o_{\mathbf{P}_M^*}(1). \quad (7.5)$$

Following Lemma (8.15) and Lemma 8.16 we can write,

$$\begin{aligned} \sqrt{n} \left(Z(\hat{\beta}_n) - Z(\beta_0) \right) &= \sqrt{n} \cdot \dot{Z}(\beta_0) \left(\hat{\beta}_n - \beta_0 \right) + o_{\mathbf{P}_M^*} \left(\sqrt{n} \cdot \|\hat{\beta}_n - \beta_0\| \right) \\ &= \sqrt{n} \cdot \dot{Z}(\beta_0) \left(\hat{\beta}_n - \beta_0 \right) + o_{\mathbf{P}_M^*}(1), \end{aligned}$$

and the proof of Theorem 2.2(ii) essentially shows that,

$$\begin{aligned} \sqrt{n} \left(Z(\hat{\beta}_n) - Z(\beta_0) \right) &= \sqrt{n} \cdot \dot{Z}(\beta_0) \left(\hat{\beta}_n - \beta_0 \right) + o_{\mathbf{P}_0^*} \left(\sqrt{n} \cdot \|\hat{\beta}_n - \beta_0\| \right) \\ &= \sqrt{n} \cdot \dot{Z}(\beta_0) \left(\hat{\beta}_n - \beta_0 \right) + o_{\mathbf{P}_0^*}(1), \end{aligned}$$

because, $\sqrt{n} \cdot \|\hat{\beta}_n - \beta_0\| = O_{\mathbf{P}_0^*}(1)$. Hence subtracting the above two equations, using (8.26), using the

second equation in (7.5) and using the limiting distribution of the bootstrap random vectors $\widehat{\mathbf{T}}_i$'s, we obtain

$$\begin{aligned}
& \sqrt{n} \cdot \left(Z(\widehat{\boldsymbol{\beta}}_n) - Z(\widehat{\boldsymbol{\beta}}_n) \right) = \dot{Z}(\boldsymbol{\beta}_0) \cdot \sqrt{n} \left(\widehat{\boldsymbol{\beta}}_n - \widehat{\boldsymbol{\beta}}_n \right) + o_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_0^*}(1) \\
& = -\sqrt{n} \cdot \left(\widehat{Z}_n(\boldsymbol{\beta}_0) - Z_n(\boldsymbol{\beta}_0) \right) + o_{\mathbf{P}_M^*}(1) \\
& = -\left[\sqrt{f_n} \cdot \widehat{\mathbb{G}}_{n_1} h_{1,\boldsymbol{\beta}_0} + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \mathbf{P}_0 \left(h_{2,\boldsymbol{\beta}_0,\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) + \sqrt{1-f_n} \cdot \widehat{\mathbb{G}}_{n_2} h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right] + o_{\mathbf{P}_M^*}(1) \\
& = -\left(\sqrt{f_n} \cdot \mathbf{I}_p \quad \frac{1-f_n}{\sqrt{f_n}} \cdot (\mathbf{A}_0 \mathbf{B}_0) \right) \begin{pmatrix} \sqrt{n_1} \cdot \left(\widehat{\mathbf{T}}_{n_1}^{(1)} - \bar{\mathbf{T}}_{n_1}^{(1)} \right) \\ \sqrt{n_1} \cdot \left(\widehat{\mathbf{T}}_{n_1}^{(2)} - \bar{\mathbf{T}}_{n_1}^{(2)} \right) \end{pmatrix} \\
& \quad - \sqrt{1-f_n} \cdot \frac{1}{\sqrt{n_2}} \sum_{i=n_1+1}^n \left\{ h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \left(\widehat{Y}_i, \widehat{\mathbf{X}}_i \right) - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \left(\widetilde{Y}_i, \mathbf{X}_i \right) \right\} + o_{\mathbf{P}_M^*}(1) \\
& \xrightarrow{d} -\mathbf{U}, \quad \text{conditionally in outer probability } \mathbf{P}_0^*,
\end{aligned}$$

where, \mathbf{U} is defined in (7.3). This follows because, $\sqrt{n_1} \left(\widehat{\mathbf{T}}_{n_1} - \bar{\mathbf{T}}_{n_1} \right)$ converges to the same limiting distribution as $\sqrt{n_1} \left(\bar{\mathbf{T}}_{n_1} - \mathbf{E}_0 \mathbf{T}_1 \right)$, conditionally in probability \mathbf{P}_0 (cf. Theorem 23.4 of van der Vaart (1998)). Also, for the same reason $\widehat{\mathbb{G}}_{n_2} h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}$ has the same limiting distribution as the term $A_{3,n}$ in the proof of Theorem 2.2(ii). Hence, $\sqrt{n} \left(\widehat{\boldsymbol{\beta}}_n - \widehat{\boldsymbol{\beta}}_n \right) \xrightarrow{d} -\left(\dot{Z}(\boldsymbol{\beta}_0) \right)^{-1} \mathbf{U}$, conditionally in outer probability \mathbf{P}_0^* . \square

Proof of Corollary 3.3. The proof follows by using Lemma 21.2 of van der Vaart (1998) and Theorem 3.2. We omit the details. \square

Proof of Corollary 3.4. The proof proceeds in the same manner as that of Corollary 3.3, by noting that: (i) the gradient of $\pi(\boldsymbol{\beta})$ is non-zero at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ and (ii) $\sqrt{n}(\pi(\widehat{\boldsymbol{\beta}}_n) - \pi(\boldsymbol{\beta}_0))$ and $\sqrt{n}(\pi(\widehat{\boldsymbol{\beta}}_n) - \pi(\widehat{\boldsymbol{\beta}}_n))$ converge to the same limiting normal distribution. \square

7.2 Proofs of Lemmas 2.1 and 3.1

Proof of Lemma 2.1. For each i.i.d. triplet $\{(Y_i, \widetilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq n_1\}$, define the random variables,

$$\left. \begin{aligned}
V_{1,i} &= \mathbf{1}(\widetilde{Y}_i = 0, Y_i = 0), & V_{2,i} &= \mathbf{1}(\widetilde{Y}_i = 0, Y_i = 1), \\
V_{3,i} &= \mathbf{1}(\widetilde{Y}_i = 1, Y_i = 0), & V_{4,i} &= \mathbf{1}(\widetilde{Y}_i = 1, Y_i = 1) \quad \text{and} \\
\mathbf{W}_i &= (V_{1,i}, V_{2,i}, V_{3,i})',
\end{aligned} \right\} \text{ for all } i = 1, \dots, n_1. \quad (7.6)$$

Write, $\pi_{k,0} = \mathbf{E}_0 V_{k,1}$, for $k = 1, 2, 3, 4$. It is easily seen that,

$$(\pi_{1,0}, \pi_{2,0}, \pi_{3,0}, \pi_{4,0})' = ((1 - \theta_{1,0})(1 - a_0), \theta_{2,0} a_0, \theta_{1,0}(1 - a_0), (1 - \theta_{2,0}) a_0)', \quad (7.7)$$

where, a_0 has been defined in (2.10). Since $a_0 \in (0, 1)$ and $\theta_{k,0}$ are bounded away from 0 and 1, this implies $\pi_{k,0} \in (0, 1)$ for all k , and satisfy the relations, $\pi_{1,0} + \pi_{3,0} = (1 - a_0)$ and $\pi_{2,0} + \pi_{4,0} = a_0$. Note

that, $\{\mathbf{W}_i : 1 \leq i \leq n_1\}$ are i.i.d. random vectors with $\mathbf{E}_0(\mathbf{W}_i) = (1 - a_0 - \pi_{3,0}, \pi_{2,0}, \pi_{3,0})'$ and it can be shown that $\mathbf{Var}_0(\mathbf{W}_i) = \boldsymbol{\Sigma}_{2,2}$ (cf. (2.11)). Using the multivariate CLT for i.i.d. random vectors we obtain, $\sqrt{n_1}(\bar{\mathbf{W}}_{n_1} - \mathbf{E}_0 \mathbf{W}_1) \xrightarrow{d} N_3(\mathbf{0}, \boldsymbol{\Sigma}_{2,2})$. Following the definition of $\hat{\boldsymbol{\theta}}_n$ in (2.2), we can express

$$\hat{\theta}_{1,n} = \frac{\bar{V}_{3,n_1} + (2n_1)^{-1}}{\bar{V}_{1,n_1} + \bar{V}_{3,n_1} + n_1^{-1}} \quad \text{and} \quad \hat{\theta}_{2,n} = \frac{\bar{V}_{2,n_1} + (2n_1)^{-1}}{1 - \bar{V}_{1,n_1} - \bar{V}_{3,n_1} + n_1^{-1}}.$$

Define the mapping

$$\boldsymbol{\phi} = (\phi_1, \phi_2) : \mathbb{R}^3 \mapsto \mathbb{R}^2, \quad \text{where} \quad \left. \begin{aligned} \phi_1(x, y, z) &= z/(x+z), \\ \phi_2(x, y, z) &= y/(1-x-z) \end{aligned} \right\}. \quad (7.8)$$

Based on the map $\boldsymbol{\phi}$ described above, define the random quantities

$$\tilde{\boldsymbol{\theta}}_n = (\tilde{\theta}_{1,n}, \tilde{\theta}_{2,n})', \quad \text{with} \quad \tilde{\theta}_{k,n} = \phi_k(\bar{V}_{1,n_1}, \bar{V}_{2,n_1}, \bar{V}_{3,n_1}), \quad \text{for } k = 1, 2, \text{ and } n \geq 1. \quad (7.9)$$

Then, it is possible to express

$$\sqrt{n_1}(\hat{\theta}_{1,n} - \tilde{\theta}_{1,n}) = \frac{1}{2\sqrt{n_1}} \cdot \frac{\bar{V}_{1,n_1} - \bar{V}_{3,n_1}}{(\bar{V}_{1,n_1} + \bar{V}_{3,n_1} + n_1^{-1}) \cdot (\bar{V}_{1,n_1} + \bar{V}_{3,n_1})}.$$

Using the WLLN for i.i.d. means and since $a_0 \in (0, 1)$, we can claim

$$\frac{\bar{V}_{1,n_1} - \bar{V}_{3,n_1}}{(\bar{V}_{1,n_1} + \bar{V}_{3,n_1} + n_1^{-1}) \cdot (\bar{V}_{1,n_1} + \bar{V}_{3,n_1})} \xrightarrow{\mathbf{P}_0} \frac{1 - a_0 - 2\pi_{3,0}}{(1 - a_0)^2} \in (0, \infty),$$

and as a result, $\sqrt{n_1}(\hat{\theta}_{1,n} - \tilde{\theta}_{1,n}) = o_{\mathbf{P}_0}(n_1^{-1/2})$. Similarly, we can show that, $\sqrt{n_1}(\hat{\theta}_{2,n} - \tilde{\theta}_{2,n}) = o_{\mathbf{P}_0}(n_1^{-1/2})$. Hence, in order to study the asymptotic distribution of $\sqrt{n_1}(\hat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$, it is enough to study the asymptotic distribution of $\sqrt{n_1}(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$, where $\tilde{\boldsymbol{\theta}}_n = (\tilde{\theta}_{1,n}, \tilde{\theta}_{2,n})'$. From (7.7) it follows that $\theta_{k,0} = \phi_k(\mathbf{E}V_{1,1}, \mathbf{E}V_{2,1}, \mathbf{E}V_{3,1})$, for $k = 1, 2$. The matrix of partial derivatives of $\boldsymbol{\phi}$ will be

$$\mathbf{B}(x, y, z) \equiv \begin{pmatrix} \frac{\partial}{\partial x} \phi_1 & \frac{\partial}{\partial y} \phi_1 & \frac{\partial}{\partial z} \phi_1 \\ \frac{\partial}{\partial x} \phi_2 & \frac{\partial}{\partial y} \phi_2 & \frac{\partial}{\partial z} \phi_2 \end{pmatrix} = \begin{pmatrix} -\frac{z}{(x+z)^2} & 0 & \frac{x}{(x+z)^2} \\ \frac{y}{(1-x-z)^2} & \frac{1}{1-x-z} & \frac{y}{(1-x-z)^2} \end{pmatrix}.$$

It follows that, $\mathbf{B}(\mathbf{E}_0 \mathbf{W}_1) = \mathbf{B}(1 - a_0 - \pi_{3,0}, \pi_{2,0}, \pi_{3,0}) = \mathbf{B}_0$, which is defined in (2.10). Following the stated assumptions, it is easy to check that each element of $\mathbf{B}(x, y, z)$ exists in a neighbourhood of $\mathbf{E}_0 \mathbf{W}_1$ and is continuous at $\mathbf{E}_0 \mathbf{W}_1$. The proof follows by using the Delta method (cf. Theorem 3.1 of van der Vaart (1998)) along with the asymptotic normality of $\bar{\mathbf{W}}_{n_1}$. \square

Proof of Lemma 3.1. Using the bootstrapped validation sample $\hat{\mathcal{X}}_{n_1}$, we define the bootstrap version of $V_{k,i}$'s and \mathbf{W}_i 's (cf. (7.6)) as,

$$\{\hat{V}_{k,i} : k = 1, 2, 3, 4\} \quad \text{and} \quad \hat{\mathbf{W}}_i = (\hat{V}_{1,i}, \hat{V}_{2,i}, \hat{V}_{3,i})', \quad \text{for all } i = 1, \dots, n_1.$$

Conditional on \mathcal{X}_{n_1} , the $\hat{\mathbf{W}}_i$'s are i.i.d. and can be considered as a with replacement sample from the finite set of random vectors $\{\mathbf{W}_i : i = 1, \dots, n_1\}$. Also, from the proof of Lemma 2.1, $\mathbf{E}_0 \mathbf{W}_i = \mathbf{0}$ and

$\text{Var}_0(\mathbf{W}_i) = \boldsymbol{\Sigma}_{2,2}$. Hence, using Theorem 23.4 of [van der Vaart \(1998\)](#) it follows that,

$$\frac{1}{\sqrt{n_1}} \sum_{i=1}^{n_1} \left(\widehat{\mathbf{W}}_i - \overline{\mathbf{W}}_{n_1} \right) \xrightarrow{d} N_3(\mathbf{0}, \boldsymbol{\Sigma}_{2,2}) \quad \text{conditionally almost surely } (\mathbf{P}_0), \quad (7.10)$$

where $\boldsymbol{\Sigma}_{2,2}$ is defined in (2.11). Now we handle two statements of Lemma 3.1 separately.

(i) This proof of this part follows from retracing the arguments used in the proof of Lemma 2.1, applied to the bootstrapped random vectors $\{\widehat{\mathbf{W}}_i : i = 1, \dots, n_1\}$, using (7.10) above, using the fact that the map ϕ in (7.8) is continuously differentiable in a neighbourhood of $\mathbf{E}_0 \mathbf{W}_1$ and by applying the continuous mapping theorem for bootstrapped random variables (cf. Theorem 23.5 of [van der Vaart \(1998\)](#)).

(ii) Following the proof of Lemma 2.1, it is enough to study the limiting distribution of $\sqrt{n_1}(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$, where $\tilde{\boldsymbol{\theta}}_n$ is the bootstrap version of $\tilde{\boldsymbol{\theta}}_n$, defined in (7.9). This follows because $\|\widehat{\mathbf{W}}_{n_1} - \mathbf{E}_0(\mathbf{W}_1)\| = o_{\mathbf{P}_0^*}(1)$. Write, $\iota = \sqrt{-1}$ and consider any fixed $\mathbf{t} (\neq \mathbf{0}) \in \mathbb{R}^2$. Define, $\sigma^2(\mathbf{t}) = \mathbf{t}' \mathbf{B}_0 \boldsymbol{\Sigma}_{2,2} \mathbf{B}_0' \mathbf{t}$ and the random variables,

$$\widehat{R}_n(\mathbf{t}) = \frac{\sqrt{n_1} \left(\mathbf{t}' \tilde{\boldsymbol{\theta}}_n - \mathbf{t}' \boldsymbol{\theta}_0 \right)}{\sigma(\mathbf{t})} \quad \text{and} \quad R_n(\mathbf{t}) = \frac{\sqrt{n_1} \left(\mathbf{t}' \tilde{\boldsymbol{\theta}}_n - \mathbf{t}' \boldsymbol{\theta}_0 \right)}{\sigma(\mathbf{t})}, \quad n \geq 1.$$

Write, $\chi(\cdot) =$ characteristic function of $N(0, 1)$. From Lemma 2.1 and (7.10), it follows that for any fixed $u \in \mathbb{R}$,

$$\left| \mathbf{E}_0 \left(e^{\iota \cdot u R_n(\mathbf{t})} \right) - \chi(u) \right| = o(1) \quad \text{and} \quad \left| \mathbf{E}_{\mathbf{M}} \left(e^{\iota \cdot u \widehat{R}_n(\mathbf{t})} \right) - \chi(u) \right| = o_{\mathbf{P}_0}(1). \quad (7.11)$$

Note that, $(\sigma(\mathbf{t}))^{-1} \cdot \sqrt{n_1} (\mathbf{t}' \tilde{\boldsymbol{\theta}}_n - \mathbf{t}' \boldsymbol{\theta}_0) = \widehat{R}_n(\mathbf{t}) + R_n(\mathbf{t})$. Since the characteristic function of any random variable is bounded, we can use the second assertion in (7.11) to claim that,

$$2 \geq \left| e^{\iota \cdot u R_n(\mathbf{t})} \cdot \left\{ \mathbf{E}_{\mathbf{M}} \left(e^{\iota \cdot u \widehat{R}_n(\mathbf{t})} \right) - \chi(u) \right\} \right| = o_{\mathbf{P}_0}(1). \quad (7.12)$$

Fix any $\epsilon > 0$ and define the sets,

$$A_n = \left[\left| e^{\iota \cdot u R_n(\mathbf{t})} \cdot \left\{ \mathbf{E}_{\mathbf{M}} \left(e^{\iota \cdot u \widehat{R}_n(\mathbf{t})} \right) - \chi(u) \right\} \right| > \epsilon \right], \quad n \geq 1.$$

Following (7.12), there exists some $n_0 = n_0(\epsilon) \in \mathbb{N}$, such that $\mathbf{P}_0(A_n) < \epsilon$, for all $n \geq n_0$. Thus, using Fubini's theorem for product measures we have

$$\begin{aligned} & \mathbf{E}_{\mathbf{r}} \left(\exp \left\{ \iota \cdot u \cdot (\sigma(\mathbf{t}))^{-1} \cdot \sqrt{n_1} \left(\mathbf{t}' \tilde{\boldsymbol{\theta}}_n - \mathbf{t}' \boldsymbol{\theta}_0 \right) \right\} \right) = \mathbf{E}_{\mathbf{r}} \left[e^{\iota \cdot u (R_n(\mathbf{t}) + \widehat{R}_n(\mathbf{t}))} \right] \\ & = \mathbf{E}_0 \left[\left(e^{\iota \cdot u R_n(\mathbf{t})} \right) \cdot \left\{ \mathbf{E}_{\mathbf{M}} \left(e^{\iota \cdot u \widehat{R}_n(\mathbf{t})} \right) - \chi(u) \right\} \cdot \{ \mathbf{1}(A_n^c) + \mathbf{1}(A_n) \} \right] + \chi(u) \cdot \mathbf{E}_0 \left(e^{\iota \cdot u R_n(\mathbf{t})} \right). \end{aligned} \quad (7.13)$$

Consider the first term on the right side of (7.13). We can write,

$$\mathbf{E}_0 \left[\left(e^{\iota \cdot u R_n(\mathbf{t})} \right) \cdot \left\{ \mathbf{E}_{\mathbf{M}} \left(e^{\iota \cdot u \widehat{R}_n(\mathbf{t})} \right) - \chi(u) \right\} \cdot \{ \mathbf{1}(A_n^c) + \mathbf{1}(A_n) \} \right] \leq \epsilon + 2 \cdot \mathbf{P}_0(A_n) \leq \epsilon + 2\epsilon = 3\epsilon,$$

whenever, $n \geq n_0$. From the first assertion in (7.11) it follows that, $\chi(u) \cdot \mathbf{E}_0 (e^{t \cdot u R_n(\mathbf{t})}) \rightarrow (\chi(u))^2$. Combining both parts, we can conclude that,

$$\mathbf{E}_{\mathbf{r}} \left(\exp \left\{ t \cdot u \cdot (\sigma(\mathbf{t}))^{-1} \cdot \sqrt{n_1} \left(\mathbf{t}' \tilde{\boldsymbol{\theta}}_n - \mathbf{t}' \boldsymbol{\theta}_0 \right) \right\} \right) \rightarrow e^{-u^2}, \quad \text{for any fixed } u \in \mathbb{R}.$$

The right side above corresponds to the characteristic function of $N(0, 2)$. By the uniqueness theorem for characteristic functions, the Cramer-Wold device and Slutsky's theorem we can claim that

$$\sqrt{n_1} \left(\tilde{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0 \right) \xrightarrow{d} N_2 \left(\mathbf{0}, 2\mathbf{B}_0 \boldsymbol{\Sigma}_{2,2} \mathbf{B}_0' \right),$$

where the convergence is understood in terms of the product probability \mathbf{Pr} . Combining this with $\sqrt{n_1} (\hat{\boldsymbol{\theta}}_n - \tilde{\boldsymbol{\theta}}_n) = o_{\mathbf{Pr}^*}(1)$, completes the proof. □

References

- Abrevaya, J. and Hausman, J. (1999). Semiparametric estimation with mismeasured dependent variables: An application to duration models for unemployment spells. *Annales d'conomie et de Statistique*, (55/56):243–275.
- Aerts, M. and Claeskens, G. (1999). Bootstrapping pseudolikelihood models for clustered binary data. *Annals of the Institute of Statistical Mathematics*, 51(3):515–530.
- Amemiya, T. (1985). *Advanced Econometrics*. Harvard University Press.
- Arcones, M. A. and Giné, E. (1992). On the bootstrap of M -estimators and other statistical functionals. In *Exploring the limits of bootstrap (East Lansing, MI, 1990)*, Wiley Ser. Probab. Math. Statist. Probab. Math. Statist., pages 13–47. Wiley, New York.
- Bhattacharya, R. N. and Ranga Rao, R. (1986). *Normal Approximation and Asymptotic Expansions*. Robert E. Krieger Publishing Co. Inc., Melbourne, FL. Reprint of the 1976 original.
- Bollinger, C. R. and David, M. H. (1997). Modeling discrete choice with response error: Food stamp participation. *Journal of the American Statistical Association*, 92(439):827–835.
- Buonaccorsi, J. P. (2010). *Measurement error: Models, methods, and applications*. Interdisciplinary Statistics. CRC Press, Boca Raton, FL.
- Carroll, R. J., Ruppert, D., Stefanski, L. A., and Crainiceanu, C. M. (2006). *Measurement error in nonlinear models*, volume 105 of *Monographs on Statistics and Applied Probability*. Chapman & Hall/CRC, Boca Raton, FL, second edition. A modern perspective.

- Carroll, R. J., Spiegelman, C. H., Lan, K. K. G., Bailey, K. T., and Abbott, R. D. (1984). On errors-in-variables for binary regression models. *Biometrika*, 71(1):19–25.
- Carroll, R. J. and Wand, M. P. (1991). Semiparametric estimation in logistic measurement error models. *J. Roy. Statist. Soc. Ser. B*, 53(3):573–585.
- Chatterjee, S. and Bose, A. (2005). Generalized bootstrap for estimating equations. *Ann. Statist.*, 33(1):414–436.
- Chen, Y. and Liang, K.-Y. (2010). On the asymptotic behaviour of the pseudolikelihood ratio test statistic with boundary problems. *Biometrika*, 97(3):603–620.
- Cheng, G. and Huang, J. Z. (2010). Bootstrap consistency for general semiparametric M-estimation. *Ann. Statist.*, 38(5):2884–2915.
- Copas, J. B. (1988). Binary regression models for contaminated data. *J. Roy. Statist. Soc. Ser. B*, 50(2):225–265. With discussion.
- Cox, D. R. and Snell, E. J. (1989). *Analysis of binary data*, volume 32 of *Monographs on Statistics and Applied Probability*. Chapman & Hall, London, second edition.
- Demidenko, E. (2004). *Mixed models*. Wiley Series in Probability and Statistics. Wiley-Interscience (John Wiley & Sons), Hoboken, NJ. Theory and applications.
- Dudley, R. M. (2014). *Uniform central limit theorems*, volume 142 of *Cambridge Studies in Advanced Mathematics*. Cambridge University Press, Cambridge.
- Duffy, S. W., Warwick, J., Williams, A. R. W., Keshavarz, H., Kaffashian, F., Rohan, T. E., Nili, F., and Sadeghi-Hassanabadi, A. (2004). A simple model for potential use with a misclassified binary outcome in epidemiology. *Journal of Epidemiology and Community Health*, 58(8):712–717.
- Edwards, J. K., Cole, S. R., Troester, M. A., and Richardson, D. B. (2013). Accounting for misclassified outcomes in binary regression models using multiple imputation with internal validation data. *American Journal of Epidemiology*, 177:904–912.
- Efron, B. (1979). Bootstrap methods: another look at the jackknife. *Ann. Statist.*, 7(1):1–26.
- Ekholm, A. and Palmgren, J. (1987). Correction for misclassification using doubly sampled data. *Journal of Official Statistics*, 3(4):419.
- Fahrmeir, L. and Kaufmann, H. (1985). Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models. *Ann. Statist.*, 13(1):342–368.

- Fuller, W. A. (2006). *Measurement error models*. Wiley Series in Probability and Statistics. Wiley-Interscience (John Wiley & Sons), Hoboken, NJ. Reprint of the 1987 original, Wiley-Interscience Paperback Series.
- Gart, J. J. and Zweifel, J. R. (1967). On the bias of various estimators of the logit and its variance with application of quantal bioassay. *Biometrika*, 54:181–187.
- Ghosh, A., Wright, F. A., and Zou, F. (2013). Unified analysis of secondary traits in case-control association studies. *J. Amer. Statist. Assoc.*, 108(502):566–576.
- Gilbert, P. B., Yu, X., and Rotnitzky, A. (2014). Optimal auxiliary-covariate-based two-phase sampling design for semiparametric efficient estimation of a mean or mean difference, with application to clinical trials. *Stat. Med.*, 33(6):901–917.
- Giné, E. and Zinn, J. (1990). Bootstrapping general empirical measures. *Ann. Probab.*, 18(2):851–869.
- Gong, G. and Samaniego, F. J. (1981). Pseudomaximum likelihood estimation: theory and applications. *Ann. Statist.*, 9(4):861–869.
- Gordis, L. (2009). *Epidemiology*. Philadelphia : Elsevier/Saunders, Fourth edition.
- Gouriéroux, C. and Monfort, A. (1981). Asymptotic properties of the maximum likelihood estimator in dichotomous logit models. *J. Econometrics*, 17(1):83–97.
- Guolo, A. (2011). Pseudo-likelihood inference for regression models with misclassified and mismeasured variables. *Statist. Sinica*, 21(4):1639–1663.
- Haldane, J. B. S. (1956). The estimation and significance of the logarithm of a ratio of frequencies. *Annals of Human Genetics*, 20(4):309–311.
- Hausman, J. (2001). Mismeasured variables in econometric analysis: Problems from the right and problems from the left. *The Journal of Economic Perspectives*, 15(4):57–67.
- Hausman, J. A., Abrevaya, J., and Scott-Morton, F. M. (1998). Misclassification of the dependent variable in a discrete-response setting. *Journal of Econometrics*, 87(2):239–269.
- Hilbe, J. M. (2009). *Logistic regression models*. Chapman & Hall/CRC Texts in Statistical Science Series. CRC Press, Boca Raton, FL.
- Hosmer, D. W. and Lemeshow, S. (2004). *Applied logistic regression*. John Wiley & Sons.
- Jewell, N. P. (2003). *Statistics for Epidemiology*. Chapman & Hall/CRC Texts in Statistical Science. Taylor & Francis.

- Kasahara, H. and Shimotsu, K. (2008). Pseudo-likelihood estimation and bootstrap inference for structural discrete markov decision models. *Journal of Econometrics*, 146(1):92 – 106.
- Kosorok, M. R. (2008). *Introduction to empirical processes and semiparametric inference*. Springer Series in Statistics. Springer, New York.
- Kothari, B. and Bandyopadhyay, T. (2011). Can india’s “literate” read? *International Review of Education*, 56(5):705–728.
- Lahiri, S. N. (2003). *Resampling methods for Dependent data*. Springer Series in Statistics. Springer-Verlag, New York.
- Lyles, R. H. and Kupper, L. L. (2013). Approximate and pseudo-likelihood analysis for logistic regression using external validation data to model log exposure. *J. Agric. Biol. Environ. Stat.*, 18(1):22–38.
- Lyles, R. H., Tang, L., Superak, H. M., King, C. C., Celentano, D. D., Lo, Y., and Sobel, J. D. (2011). Validation Data-based Adjustments for Outcome Misclassification in Logistic Regression: An Illustration. *Epidemiology*, 22(4):589–597.
- Magder, L. S. and Hughes, J. P. (1997). Logistic regression when the outcome is measured with uncertainty. *American Journal of Epidemiology*, 146(2):195–203.
- Meyer, B. and Mittag, N. (2016). Misclassification in binary choice models. Technical report. Preprint available at [http://harris.uchicago.edu/sites/default/files/Misreporting in BCM-March-8-2016.pdf](http://harris.uchicago.edu/sites/default/files/Misreporting%20in%20BCM-March-8-2016.pdf).
- Neuhaus, J. M. (1999). Bias and efficiency loss due to misclassified responses in binary regression. *Biometrika*, 86(4):843–855.
- Parzen, M., Lipsitz, S., Ibrahim, J., and Klar, N. (2002). An estimate of the odds ratio that always exists. *J. Comput. Graph. Statist.*, 11(2):420–436.
- Præstgaard, J. and Wellner, J. A. (1993). Exchangeably weighted bootstraps of the general empirical process. *Ann. Probab.*, 21(4):2053–2086.
- Roy, S., Banerjee, T., and Maiti, T. (2005). Measurement error model for misclassified binary responses. *Statistics in medicine*, 24(2):269–283.
- Savoca, E. (2011). Accounting for misclassification bias in binary outcome measures of illness: The case of post-traumatic stress disorder in male veterans. *Sociological Methodology*, 41(1):49–76.
- Sposto, R., Preston, D. L., Shimizu, Y., and Mabuchi, K. (1992). The effect of diagnostic misclassification on non-cancer and cancer mortality dose response in a-bomb survivors. *Biometrics*, 48(2):605–617.

- Stefanski, L. A. and Carroll, R. J. (1985). Covariate measurement error in logistic regression. *Ann. Statist.*, 13(4):1335–1351.
- van der Vaart, A. W. (1998). *Asymptotic statistics*, volume 3 of *Cambridge Series in Statistical and Probabilistic Mathematics*. Cambridge University Press, Cambridge.
- van der Vaart, A. W. (2002). *Semiparametric statistics*, volume 1781 of *Lecture Notes in Mathematics*. Springer-Verlag, Berlin. Lectures from the 29th Summer School on Probability Theory held in Saint-Flour, July 8–24, 1999, Edited by Pierre Bernard.
- van der Vaart, A. W. and Wellner, J. A. (1996). *Weak convergence and empirical processes*. Springer Series in Statistics. Springer-Verlag, New York. With applications to statistics.
- van der Vaart, A. W. and Wellner, J. A. (2007). Empirical processes indexed by estimated functions. In *Asymptotics: particles, processes and inverse problems*, volume 55 of *IMS Lecture Notes Monogr. Ser.*, pages 234–252. Inst. Math. Statist., Beachwood, OH.
- Varadhan, R. and Gilbert, P. (2009). BB: An R package for solving a large system of nonlinear equations and for optimizing a high-dimensional nonlinear objective function. *Journal of Statistical Software*, 32(1):1–26.
- Wang, C. Y. and Wang, S. (1997). Semiparametric methods in logistic regression with measurement error. *Statist. Sinica*, 7(4):1103–1120.
- Wang, Y.-G. and Zhao, Y. (2007). A modified pseudolikelihood approach for analysis of longitudinal data. *Biometrics*, 63(3):681–689.
- Wellner, J. A. and Zhan, Y. (1996). Bootstrapping Z-estimators. Technical report, Department of Statistics, University of Washington. Technical Report 308.

8 Appendix

We follow the same notational convention, as stated in the beginning of Section 7.

8.1 Auxiliary lemmas required for proving Theorem 2.2

We begin with a simple upper bound which will be used throughout the proofs.

Lemma 8.1. *Under assumption (A2), the following inequality holds:*

$$4 \leq \frac{1}{h_{3,\beta,\theta}(\mathbf{x}) \cdot \{1 - h_{3,\beta,\theta}(\mathbf{x})\}} < M_0 \equiv \frac{1}{\max_{i=1,2} \delta_i(1 - \delta_i)} < \infty, \quad \text{for all } \mathbf{x}, \beta \in \mathbb{R}^p \text{ and } \theta \in \Theta. \quad (8.1)$$

Proof of Lemma 8.1. From the expression of $h_{3,\beta,\theta}(x)$ (cf. (2.3)), it can be seen as a convex combination of θ_1 and $(1 - \theta_2)$, for any \mathbf{x}, β and θ . Hence, $h_{3,\beta,\theta}(\mathbf{x}) \in (\theta_1, 1 - \theta_2)$ or $(1 - \theta_2, \theta_1)$, depending on which probability is larger. Note that, $(1 - \theta_2) \neq \theta_1$, due to assumption (A2)(ii). In case $\theta_1 < (1 - \theta_2)$, using assumption (A2)(i) we must have, $\delta_1 < \theta_1 < (1 - \theta_2) < (1 - \delta_1)$. And if, $(1 - \theta_2) < \theta_1$, then $(1 - \delta_2) < (1 - \theta_2) < \theta_1 < \delta_2$. Combining both cases we have,

$$\frac{1}{4} \geq h_{3,\beta,\theta}(\mathbf{x}) \cdot \{1 - h_{3,\beta,\theta}(\mathbf{x})\} > \max_{i=1,2} \delta_i(1 - \delta_i) > 0, \quad \text{for all } \mathbf{x}, \beta \text{ and } \theta \in \Theta,$$

from which (8.1) follows. □

The following result shown in Lemma 8.2 shows the uniqueness and well-separability of β_0 as a root of the equation $Z_1(\beta) = \mathbf{0}$, under both without and with intercept models in (1.1) and (2.9) separately to illustrate some fine issues related to identifiability of the logistic regression model. For the rest of the article, it is enough to consider (1.1) for all calculations. In case (2.9) is considered, we will simply substitute $X_1 = 1$ with probability 1 and assume that the conditions in assumption (A3[†]) hold. Lemma 8.3 shows the uniqueness and well-separability of β_0 as a root of the equation $Z(\beta) = \mathbf{0}$, using Lemma 8.2.

Lemma 8.2. *The result is stated separately for the without and with intercept models.*

- (i) *Consider the model (1.1), without an intercept term. Suppose, assumptions (A1) and (A3) holds. Then, $Z_1(\beta)$ (cf. (2.8)) has a unique zero at $\beta = \beta_0$ and the unique zero at β_0 is well-separated in the following sense,*

$$\inf_{\beta: \|\beta - \beta_0\| \geq \delta} \|Z_1(\beta)\| > 0, \quad \text{for all } \delta > 0.$$

- (ii) *Consider the model (2.9), with an intercept term. Suppose, assumptions (A1) and (A3[†]) holds. Then, the statement about $Z_1(\beta)$ in part (i) above holds in this case.*

Proof of Lemma 8.2. We begin by proving the without intercept case.

- (i) In this case, the first component X_1 of the p -dimensional covariate vector \mathbf{X} is assumed to be a non-degenerate component. When the underlying regression coefficient is $\boldsymbol{\beta}$, then (Y, \mathbf{X}) has joint density

$$p_{\boldsymbol{\beta}}(y, \mathbf{x}) = \{\psi(\mathbf{x}'\boldsymbol{\beta})\}^y \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\}^{1-y}, \quad y \in \{0, 1\}, \quad \mathbf{x} \in \mathbb{R}^p,$$

with respect to the product measure $\mu = \nu \times Q$, with ν being the counting measure on $\{0, 1\}$ and Q denoting the marginal of \mathbf{X} . Also, the family of densities $\{p_{\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is identifiable. If not, then for some $\boldsymbol{\beta}_1 \neq \boldsymbol{\beta}_2 \in \mathbb{R}^p$, we will have $p_{\boldsymbol{\beta}_1} = p_{\boldsymbol{\beta}_2}$, a.e. μ . Using $y = 0$ and 1 , we obtain

$$\begin{aligned} Q(\psi(\mathbf{X}'\boldsymbol{\beta}_1) = \psi(\mathbf{X}'\boldsymbol{\beta}_2)) &= 1 \Rightarrow Q((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2)'\mathbf{X} = 0) = 1, \\ &\Rightarrow Q(\mathbf{X} \in \text{some } (p-1) \text{ dimensional subspace of } \mathbb{R}^p) = 1, \end{aligned}$$

However, under assumption (A3) this is impossible and hence, $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$, which implies identifiability. For each $\boldsymbol{\beta}$, define the expected log-likelihood as $M_1(\boldsymbol{\beta}) = \mathbf{P}_0 \log p_{\boldsymbol{\beta}}(Y, \mathbf{X})$, where \mathbf{P}_0 is the distribution corresponding to the density $p_{\boldsymbol{\beta}_0}$. Using assumption (A3) and the DCT (dominated convergence theorem), it is easy to check that the gradient of $M_1(\boldsymbol{\beta})$, denoted by $\nabla M_1(\boldsymbol{\beta})$, exists at all $\boldsymbol{\beta}$ and $\nabla M_1(\boldsymbol{\beta}) = Z_1(\boldsymbol{\beta})$ (cf. (2.8)). Extending the same argument it follows that all second order partial derivatives of $M_1(\boldsymbol{\beta})$ will exist at every $\boldsymbol{\beta}$ with,

$$\frac{\partial^2 M_1(\boldsymbol{\beta})}{\partial \beta_i \partial \beta_j} = - \int x_i x_j \psi(\mathbf{x}'\boldsymbol{\beta}) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\} dQ(\mathbf{x}), \quad \text{for all } i, j.$$

Consider any $\mathbf{u} \neq \mathbf{0} \in \mathbb{R}^p$. Then,

$$\begin{aligned} &\sum_{i=1}^p \sum_{j=1}^p u_i u_j \cdot \frac{\partial^2 M_1(\boldsymbol{\beta})}{\partial \beta_i \partial \beta_j} \\ &= - \int \psi(\mathbf{x}'\boldsymbol{\beta}) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\} \cdot (\mathbf{u}'\mathbf{x})^2 dQ(\mathbf{x}) \\ &= - \int_{\mathbf{u}'\mathbf{x}=0} \psi(\mathbf{x}'\boldsymbol{\beta}) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\} \cdot (\mathbf{u}'\mathbf{x})^2 dQ(\mathbf{x}) - \int_{\mathbf{u}'\mathbf{x} \neq 0} \psi(\mathbf{x}'\boldsymbol{\beta}) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\} \cdot (\mathbf{u}'\mathbf{x})^2 dQ(\mathbf{x}) < 0, \end{aligned}$$

since, $Q(\{\mathbf{x} : \mathbf{u}'\mathbf{x} = 0\}) < 1$ using assumption (A3) and $\psi(\mathbf{x}'\boldsymbol{\beta}) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\} < 0$, for all $\mathbf{x}, \boldsymbol{\beta} \in \mathbb{R}^p$. As a result, the Hessian of $M_1(\boldsymbol{\beta})$ is strictly negative definite at each $\boldsymbol{\beta} \in \mathbb{R}^p$, which implies $\boldsymbol{\beta} \mapsto M_1(\boldsymbol{\beta})$ is strictly concave. A strictly concave function has a unique point of maxima and in this case, due to the identifiability of the family $\{p_{\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ it follows from Lemma 5.35 of [van der Vaart \(1998\)](#) that the unique point of maxima of $M_1(\boldsymbol{\beta})$ is at the true parameter $\boldsymbol{\beta}_0$. Obviously, $\nabla M_1(\boldsymbol{\beta}_0) = Z_1(\boldsymbol{\beta}_0) = \mathbf{0}$. If possible, assume that there exists some $\boldsymbol{\beta}_1 \neq \boldsymbol{\beta}_0$, such that $\nabla M_1(\boldsymbol{\beta}_1) = Z_1(\boldsymbol{\beta}_1) = \mathbf{0}$. But due to the strict concavity of $M_1(\boldsymbol{\beta})$, this implies $\boldsymbol{\beta}_1$ will be a unique maxima of M_1 , which is not possible as per the above argument. Hence $\boldsymbol{\beta}_0$ is the unique zero of $Z_1(\boldsymbol{\beta})$ which proves the first part of the statement. To prove the second part of the statement we use Lemma 8.17, using the fact that $M_1 : \mathbb{R}^p \mapsto \mathbb{R}$ is strictly concave, with a unique maxima at $\boldsymbol{\beta}_0$. This completes the proof.

- (ii) In case an intercept term is included in the model, $\mathbf{X} = (1, X_2, \dots, X_p)'$ and

$$\mathbf{P}_0(Y = 1 | (X_2, \dots, X_p) = (x_2, \dots, x_p)) = \psi(\beta_{1,0} + \sum_{j=2}^p \beta_{j,0} x_j),$$

where, $\beta_{1,0}$ is the intercept term and $\boldsymbol{\beta}_0 = (\beta_{1,0}, \beta_{2,0}, \dots, \beta_{p,0})'$. Note that we can re-write,

$$\beta_{1,0} + \sum_{j=2}^p \beta_{j,0} \cdot x_j = \beta_{1,0} + \sum_{j=2}^p \beta_{j,0} \cdot \mathbf{E}X_j + \sum_{j=2}^p \beta_{j,0} \cdot (x_j - \mathbf{E}X_j) = \tilde{\beta}_{1,0} + \sum_{j=2}^p \beta_{j,0} \cdot (x_j - \mathbf{E}X_j),$$

where, $\tilde{\beta}_{1,0} = \beta_{1,0} + \sum_{j=2}^p \beta_{j,0} \cdot \mathbf{E}X_j$. Any of these two representations of the same model can be used and due to assumption (A3[†]), both these representations are equivalent. Similar to the no-intercept case, consider densities corresponding to two models $p_{\boldsymbol{\beta}^{(1)}}$ and $p_{\boldsymbol{\beta}^{(2)}}$, which satisfy $p_{\boldsymbol{\beta}^{(1)}} = p_{\boldsymbol{\beta}^{(2)}}$, a.e. μ , where, $\boldsymbol{\beta}^{(i)} = (\beta_1^{(i)}, \dots, \beta_p^{(i)})'$, $i = 1, 2$. Following earlier arguments, we can claim that

$$Q \left(\sum_{j=2}^p (\beta_j^{(1)} - \beta_j^{(2)}) X_j = - (\beta_1^{(1)} - \beta_1^{(2)}) \right) = 1,$$

where, Q is the joint distribution of $(X_2, \dots, X_p)'$. If possible, assume that, $(\beta_2^{(1)}, \dots, \beta_p^{(1)})' \neq (\beta_2^{(2)}, \dots, \beta_p^{(2)})'$. Then, using similar arguments as earlier applied to the vector $(X_2, \dots, X_p)'$, we can conclude that $(X_2, \dots, X_p)'$ is supported on a $(p-2)$ dimensional hyperplane. This hyperplane does not (does) pass through the origin if $(\beta_1^{(1)} - \beta_1^{(2)})$ is not equal (equal) to zero. This implies linear dependence among the $(p-1)$ components of $(X_2, \dots, X_p)'$, which contradicts assumption (A3[†]). Hence the only possibility is, $(\beta_2^{(1)}, \dots, \beta_p^{(1)})' = (\beta_2^{(2)}, \dots, \beta_p^{(2)})'$, which immediately implies $\beta_1^{(1)} = \beta_1^{(2)}$. Thus, $\boldsymbol{\beta}^{(1)} = \boldsymbol{\beta}^{(2)}$ and the class of models with intercept is identifiable. The remaining part of the proof follows the same arguments as in the without intercept case. □

Lemma 8.3. *Suppose, assumptions (A1)-(A4) hold. Then, $Z(\boldsymbol{\beta})$ (cf. (2.8)) has an unique zero at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ and*

$$0 = \|Z(\boldsymbol{\beta}_0)\| < \inf_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\| \geq \delta} \|Z(\boldsymbol{\beta})\|, \quad \text{for all } \delta > 0.$$

Proof of Lemma 8.3. The proof relies on Lemma 8.2. Using (2.8), we have

$$Z(\boldsymbol{\beta}) = \mathbf{E} \left[\mathbf{X} \cdot \{\psi(\mathbf{X}'\boldsymbol{\beta}_0) - \psi(\mathbf{X}'\boldsymbol{\beta})\} \cdot \left\{ f + (1-f) \cdot (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot \frac{\psi(\mathbf{X}'\boldsymbol{\beta})\{1 - \psi(\mathbf{X}'\boldsymbol{\beta})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0}(\mathbf{X})\{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0}(\mathbf{X})\}} \right\} \right].$$

Obviously, $Z(\boldsymbol{\beta}_0) = \mathbf{0}$. Since $\psi(u) \cdot \{1 - \psi(u)\} \in (0, 1/4]$, for all $u \in \mathbb{R}$, using (8.1) we obtain

$$f < f + (1-f) \cdot (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot \frac{\psi(\mathbf{x}'\boldsymbol{\beta})\{1 - \psi(\mathbf{x}'\boldsymbol{\beta})\}}{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0}(\mathbf{x})\{1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0}(\mathbf{x})\}} \leq f + \frac{(1-f) \cdot (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot M_0}{4},$$

for all \mathbf{x} and $\boldsymbol{\beta}$. Write, $C_1 = f$ and $C_2 = f + (1-f) \cdot (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot M_0/4$. Note that, $0 < C_1 < C_2$, due to assumption (A4). As a result we have a componentwise inequality, for all $\boldsymbol{\beta}$

$$C_1 \cdot Z_1(\boldsymbol{\beta}) = C_1 \cdot \mathbf{E} [\mathbf{X} \cdot \{\psi(\mathbf{X}'\boldsymbol{\beta}_0) - \psi(\mathbf{X}'\boldsymbol{\beta})\}] < Z(\boldsymbol{\beta}) < C_2 \cdot \mathbf{E} [\mathbf{X} \cdot \{\psi(\mathbf{X}'\boldsymbol{\beta}_0) - \psi(\mathbf{X}'\boldsymbol{\beta})\}] = C_2 \cdot Z_1(\boldsymbol{\beta}),$$

From Lemma 8.2, we know that $Z_1(\boldsymbol{\beta})$ has an unique zero at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$. Hence, for all other $\boldsymbol{\beta}_1 \neq \boldsymbol{\beta}_0$, there will be a component $j \in \{1, \dots, p\}$, such that $Z_{1,j}(\boldsymbol{\beta}_1) \neq 0$ (and either it will be negative or positive). Using the above componentwise inequality, we must have

$$C_1 \cdot Z_{1,j}(\boldsymbol{\beta}_1) < Z_j(\boldsymbol{\beta}_1) < C_2 \cdot Z_{1,j}(\boldsymbol{\beta}_1).$$

Since $C_1, C_2 > 0$, so $Z_j(\boldsymbol{\beta}_1) \neq 0$ and has the same sign as $Z_{1,j}(\boldsymbol{\beta}_1)$. As a result, we have showed that $Z(\boldsymbol{\beta})$ must have an unique zero at $\boldsymbol{\beta} = \boldsymbol{\beta}_0$. Now, using the componentwise upper and lower bounds on $Z(\boldsymbol{\beta})$, we have

$$C_1 \cdot \|Z_1(\boldsymbol{\beta})\| < \|Z(\boldsymbol{\beta})\| < C_2 \cdot \|Z_1(\boldsymbol{\beta})\|.$$

Now using part (i) of Lemma 8.2, it is easy to see that

$$\begin{aligned} 0 = \|Z(\boldsymbol{\beta}_0)\| = \|Z_1(\boldsymbol{\beta}_0)\| &< \inf_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\| \geq \delta} \|Z_1(\boldsymbol{\beta})\| < \frac{1}{C_1} \cdot \inf_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\| \geq \delta} \|Z(\boldsymbol{\beta})\|, \\ \Rightarrow 0 = \|Z(\boldsymbol{\beta}_0)\| &< \inf_{\boldsymbol{\beta}: \|\boldsymbol{\beta} - \boldsymbol{\beta}_0\| \geq \delta} \|Z(\boldsymbol{\beta})\|, \quad \text{for all } \delta > 0. \end{aligned}$$

This completes the proof. \square

Throughout the proofs, we will be dealing with the estimating functions $h_{1,\boldsymbol{\beta}}$ and $h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0}$ (cf. (2.3)). Lemma 8.4 states that these classes are \mathbf{P}_0 -Donsker classes.

Lemma 8.4. *Suppose, assumptions (A1) - (A3) hold. Then,*

- (i) $\{h_{1,\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is \mathbf{P}_0 -Donsker.
- (ii) $\{h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is \mathbf{P}_0 -Donsker.

Proof of Lemma 8.4. We prove each part separately.

- (i) Consider the class $\{g_{\boldsymbol{\beta}}(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}$. By referring to the standard orthonormal basis $\{\mathbf{e}_j : 1 \leq j \leq p\}$ of \mathbb{R}^p , we can write

$$g_{\boldsymbol{\beta}}(\mathbf{x}) = \sum_{i=1}^p \beta_i x_i = \sum_{i=1}^p \beta_i \cdot (\mathbf{x}'\mathbf{e}_i),$$

for all $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ and $\mathbf{x} = (x_1, \dots, x_p)' \in \mathbb{R}^p$. Hence, $\{g_{\boldsymbol{\beta}}(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is generated by linear combinations of the finite set $\{g_{\mathbf{e}_j}(\mathbf{x}) : 1 \leq j \leq p\} \subset \{g_{\boldsymbol{\beta}}(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}$. This implies, $\{g_{\boldsymbol{\beta}}(\mathbf{x}) = \mathbf{x}'\boldsymbol{\beta} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is a Vapnik-Chervonenkis (VC) class of functions, with VC dimension $\leq (p+2)$. Also, the logistic link function, $\psi : \mathbb{R} \mapsto [0, 1]$ defined as $\psi(u) = (1 + \exp\{-u\})^{-1}$, is monotone on \mathbb{R} . This leads to the representation, $\{\psi_{\boldsymbol{\beta}}(\mathbf{x}) = \psi(\mathbf{x}'\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\} = \{\psi(g_{\boldsymbol{\beta}}(\mathbf{x})) : \boldsymbol{\beta} \in \mathbb{R}^p\}$. Hence, $\{\psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$ will be a VC class. Finally,

$$\{h_{1,\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{R}^p\} = \{\mathbf{x} \cdot y - \mathbf{x} \cdot \psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\} = \mathbf{x} \cdot y - \mathbf{x} \cdot \{\psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\}.$$

Since, $(y, \mathbf{x}) \mapsto \mathbf{x} \cdot y$ and $(y, \mathbf{x}) \mapsto \mathbf{x}$, are two fixed functions defined on the sample space, it follows that $\{h_{1,\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is also a VC class. Now, consider the countable subclass, $\{h_{1,\boldsymbol{\beta}} : \boldsymbol{\beta} \in \mathbb{Q}^p\}$,

where \mathbb{Q} denotes the set of rationals. Since \mathbb{Q}^p is dense in \mathbb{R}^p , hence for any $\beta \in \mathbb{R}^p$, there will exist a sequence $\{\beta_m : m \geq 1\} \in \mathbb{Q}^p$, such that $h_{1,\beta}(y, \mathbf{x}) = \lim_{m \rightarrow \infty} h_{1,\beta_m}(y, \mathbf{x})$, for all (y, \mathbf{x}) . Thus, $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ will be a pointwise measurable class. Finally, note that due to assumption (A3), the j th component of $h_{1,\beta}(y, \mathbf{x})$ is bounded by the envelope function $2|x_j| \in L_2(\mathbf{P}_0)$, for all $j \in \{1, \dots, p\}$. Thus, $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ has the envelope function $2 \cdot \max\{\mathbf{E}|X_j| : 1 \leq j \leq p\} \cdot \mathbf{1}_{p \times 1} \in L_2(\mathbf{P}_0)$. Combining all steps we conclude that $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ is a \mathbf{P}_0 -Donsker class.

(ii) In order to show that the class of p -dimensional functions $\{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$ is \mathbf{P}_0 -Donsker, it is enough to show that classes corresponding to each component of $h_{2,\beta,\theta_0}(\tilde{y}, \mathbf{x})$ will be Donsker. Without loss of generality, consider the first component of $h_{2,\beta,\theta_0}(\tilde{y}, \mathbf{x})$ and define the class of functions,

$$\mathcal{G}_1 = \left\{ g_\beta(\tilde{y}, \mathbf{x}) = \frac{(1 - \theta_{1,0} - \theta_{2,0}) \cdot \psi(\mathbf{x}'\beta) \cdot \{1 - \psi(\mathbf{x}'\beta)\} \cdot \{\tilde{y} - h_{3,\beta,\theta_0}(\mathbf{x})\}}{h_{3,\beta,\theta_0}(\mathbf{x}) \cdot \{1 - h_{3,\beta,\theta_0}(\mathbf{x})\}} : \beta \in \mathbb{R}^p \right\}.$$

From the proof of part (i) above, we know $\{\psi(\mathbf{x}'\beta) : \beta \in \mathbb{R}^p\}$ is a VC class and the constants $\theta_{1,0}$ and $(1 - \theta_{1,0} - \theta_{2,0})$ can be considered as fixed functions on the sample space. Hence, $\{h_{3,\beta,\theta_0}(\mathbf{x}) : \beta \in \mathbb{R}^p\}$ will be a VC class. Since $u \mapsto 1/u$ is monotone for positive u , the class $\{h_{3,\beta,\theta_0}^{-1}(\mathbf{x}) : \beta \in \mathbb{R}^p\}$ will be a VC class. Following the proof of Lemma 8.1, without loss of generality we assume that, $h_{3,\beta,\theta_0}(\mathbf{x}) \in (\delta_1, 1 - \delta_1)$ (similar argument applies if the range is $(1 - \delta_2, \delta_2)$). Hence the constant δ_1^{-1} is an envelope function for this class. Similar arguments show that the class $\{(1 - h_{3,\beta,\theta_0}(\tilde{y}, \mathbf{x}))^{-1} : \beta \in \mathbb{R}^p\}$ will also be VC with envelope function $(1 - \delta_1)^{-1}$. Also, the classes $\{(1 - \theta_{1,0} - \theta_{2,0}) \cdot \psi(\mathbf{x}'\beta) : \beta \in \mathbb{R}^p\}$, $\{1 - \psi(\mathbf{x}'\beta) : \beta \in \mathbb{R}^p\}$ and $\{\tilde{y} - h_{3,\beta,\theta_0}(\mathbf{x}) : \beta \in \mathbb{R}^p\}$ will be VC, following usual VC preservation properties, and will have envelope functions 2, 1 and 2 respectively. It is known that any VC class has a bounded uniform entropy integral (BUEI) with respect to any envelope function for that class. Also products of individual BUEI classes will be BUEI with envelope function equal to the product of envelope functions of individual classes (cf. Theorem 9.15 of Kosorok (2008)). Hence, the class \mathcal{G}_1 which is a subset of the product of these above mentioned classes will also be BUEI with envelope function $G_1(\tilde{y}, \mathbf{x}) = 4(\delta_1 \cdot (1 - \delta_1))^{-1}$. Also define the trivial class $\mathcal{G}_2 = \{x_1\}$, consisting of only one function $(\tilde{y}, \mathbf{x}) \mapsto x_1$. For any $g = g_\beta \in \mathcal{G}_1$ and $g_2 \in \mathcal{G}_2$, define the map $\phi : \mathcal{G}_1 \times \mathcal{G}_2 \mapsto \mathbb{R}$ as, $\phi(g_\beta, g_2) = g_\beta \cdot g_2$. Due to the trivial nature of \mathcal{G}_2 , we have

$$\begin{aligned} |\phi(g_{\beta_1}, g_2) - \phi(g_{\beta_2}, g_2)|^2 &= |x_1|^2 \cdot |g_{\beta_1}(\tilde{y}, \mathbf{x}) - g_{\beta_2}(\tilde{y}, \mathbf{x})|^2 + 0 \cdot |g_2(\tilde{y}, \mathbf{x}) - g_2(\tilde{y}, \mathbf{x})|^2 \\ &= L_1^2(\tilde{y}, \mathbf{x}) \cdot |g_{\beta_1}(\tilde{y}, \mathbf{x}) - g_{\beta_2}(\tilde{y}, \mathbf{x})|^{2 \cdot \alpha_1} + L_2^2(\tilde{y}, \mathbf{x}) \cdot |g_2(\tilde{y}, \mathbf{x}) - g_2(\tilde{y}, \mathbf{x})|^{2 \cdot \alpha_2}, \end{aligned}$$

where, $L_1(\tilde{y}, \mathbf{x}) = |x_1|$, $L_2(\tilde{y}, \mathbf{x}) = 0$, $\alpha_1 = 1$ and $\alpha_2 = 1$. This satisfies equation (2.10.19) of van der Vaart and Wellner (1996) and shows that the map ϕ is Lipschitz of orders $(1, 1)$. Since $L_2(\tilde{y}, \mathbf{x}) = 0$, we can neglect the second Lipschitz order. Since the class \mathcal{G}_1 is BUEI with envelope function G_1 , hence the uniform entropy integral

$$\int_0^\delta \sup_R \sqrt{\log N \left(\epsilon \cdot \|G_1\|_{R,2}, \mathcal{G}_1, L_2(R) \right)} d\epsilon < \infty, \quad \text{for any } \delta > 0,$$

where the supremum is over probability measures R on the sample space with $\|G_1\|_{R,2} > 0$ and $N(\epsilon, \mathcal{G}_1, L_2(R))$ are the associated covering numbers. The corresponding covering number of the trivial class $\mathcal{G}_2 = \{x_1\}$ will be equal to 1, for any $\epsilon > 0$ and as a result the uniform entropy integral for \mathcal{G}_2 will be finite for all $\delta > 0$. Also, $\mathbf{P}_0|X_1|^2 \cdot G_1^2(\tilde{Y}, \mathbf{X}) < \infty$, under assumption (A3) and the class $\phi(\mathcal{G}_1, \mathcal{G}_2)$ can be shown as pointwise measurable by considering the subclass of \mathcal{G}_1 indexed by all $\beta \in \mathbb{Q}^p$. Now, using Theorem 2.10.20 of [van der Vaart and Wellner \(1996\)](#), we can claim that

$$\phi(\mathcal{G}_1, \mathcal{G}_2) = \{\text{first component of } h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$$

will be \mathbf{P}_0 -Donsker. Using the same argument for remaining components of h_{2,β,θ_0} we can complete the proof. □

Lemma 8.5. *Suppose, assumptions (A1)-(A4) hold. Then, for any fixed $\mathbf{u}_0 \in \mathbb{R}^2$,*

$$\sup_{\beta \in \mathbb{R}^p} \left\| \mathbb{G}_{n_2} \left(h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0} - h_{2,\beta,\theta_0} \right) \right\| = o_{\mathbf{P}_0^*}(1).$$

Proof of Lemma 8.5. It will be enough to show the convergence of each component of $h_{2,\beta,\theta}$ separately and without loss of generality we consider the first component, which is denoted by the same symbol. Define the classes of functions:

$$\mathcal{G}_n = \left\{ h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0} - h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p \right\}, \quad n \geq 1. \quad (8.2)$$

We use Lemma 2.2 of [van der Vaart and Wellner \(2007\)](#), which provides sufficient conditions to ensure that, $\sup_{g \in \mathcal{G}_n} |\mathbb{G}_{n_2} g| = o_{\mathbf{P}_0^*}(1)$. These sufficient conditions are essentially the following:

(C.1) $\sup_{g \in \mathcal{G}_n} \mathbf{P}_0 g^2 \rightarrow 0$ as $n \rightarrow \infty$.

(C.2) The class \mathcal{G}_n is pointwise measurable with envelope function G_n and satisfies $J(\delta_n, \mathcal{G}_n, L_2) \rightarrow 0$ as $\delta_n \downarrow 0$, where,

$$J(\delta, \mathcal{G}_n, L_2) = \int_0^\delta \sup_R \sqrt{\log N \left(\epsilon \|G_n\|_{R,2}, \mathcal{G}_n, L_2(R) \right)} d\epsilon.$$

(C.3) The envelope functions $\{G_n : n \geq 1\}$, used in part (C.2) above, satisfy the Lindeberg condition:

$$\mathbf{P}_0 G_n^2 = O(1) \quad \text{and} \quad \mathbf{P}_0 G_n^2 \mathbf{1}(G_n > \epsilon \sqrt{n}) \rightarrow 0 \quad \text{for all } \epsilon > 0.$$

These three conditions are verified below.

Verification of (C.1): Define, $D_k(h_{2,\beta,\theta})$, $k = 1, 2$, as the partial derivative of $h_{2,\beta,\theta}$ with respect to θ_k , $k = 1, 2$. We can show that, for all $\theta \in \Theta$ and any $\beta \in \mathbb{R}^p$,

$$\begin{aligned} & D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x}) \\ &= -x_1 \cdot \frac{\psi(\mathbf{x}'\beta)\{1 - \psi(\mathbf{x}'\beta)\} \cdot \{\tilde{y} - h_{3,\beta,\theta}(\mathbf{x})\}}{h_{3,\beta,\theta}(\mathbf{x})\{1 - h_{3,\beta,\theta}(\mathbf{x})\}} \cdot \left[1 + \frac{(1 - \theta_1 - \theta_2) \cdot \{1 - \psi(\mathbf{x}'\beta)\} \cdot \{\tilde{y} - h_{3,\beta,\theta}(\mathbf{x})\}}{h_{3,\beta,\theta}(\mathbf{x})\{1 - h_{3,\beta,\theta}(\mathbf{x})\}} \right], \end{aligned} \quad (8.3)$$

and

$$\begin{aligned}
& D_2(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x}) \\
&= -x_1 \cdot \frac{\psi(\mathbf{x}'\beta)\{1-\psi(\mathbf{x}'\beta)\} \cdot \{\tilde{y}-h_{3,\beta,\theta}(\mathbf{x})\}}{h_{3,\beta,\theta}(\mathbf{x})\{1-h_{3,\beta,\theta}(\mathbf{x})\}} \cdot \left[1 - \frac{(1-\theta_1-\theta_2) \cdot \{1-\psi(\mathbf{x}'\beta)\} \cdot \{\tilde{y}-h_{3,\beta,\theta}(\mathbf{x})\}}{h_{3,\beta,\theta}(\mathbf{x})\{1-h_{3,\beta,\theta}(\mathbf{x})\}}\right].
\end{aligned} \tag{8.4}$$

Combining (8.1) with (8.3) and (8.4), it follows that for $k = 1, 2$,

$$\sup_{\beta} \sup_{\theta \in \Theta} |D_k(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x})| \leq |x_1| \cdot \frac{M_0}{4} \cdot 2 \cdot [1 + 2 \cdot 1 \cdot 2 \cdot M_0] \leq K_0|x_1|, \quad \text{for all } (\tilde{y}, \mathbf{x}), \tag{8.5}$$

where, $K_0 = 5M_0^2/2$. Consider a non-random sequence θ_n satisfying $\|\theta_n - \theta_0\| \rightarrow 0$. Since $\theta \in \Theta$ and Θ is open (cf. assumption (A2)), there exists some $n_0 \in \mathbb{N}$, such that $\theta_n \in \Theta$ for all $n \geq n_0$. For any fixed β , there exists sequences $\{a_n\}, \{b_n\} \in (0, 1)$ (obtained from using the Mean-value theorem around θ_0), such that

$$\begin{aligned}
& |h_{2,\beta,\theta_n}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta_0}(\tilde{y}, \mathbf{x})| \\
&= |\theta_{1,n} - \theta_{1,0}| \cdot |D_1(h_{2,\beta,a_n \cdot \theta_n + (1-a_n) \cdot \theta_0})(\tilde{y}, \mathbf{x})| + |\theta_{2,n} - \theta_{2,0}| \cdot |D_2(h_{2,\beta,b_n \cdot \theta_n + (1-b_n) \cdot \theta_0})(\tilde{y}, \mathbf{x})| \\
&\leq \sum_{k=1}^2 |\theta_{k,n} - \theta_{k,0}| \cdot K_0|x_1|,
\end{aligned} \tag{8.6}$$

following (8.5). For any $\mathbf{u}_0 \in \mathbb{R}^2$, $\theta_n = \theta_0 + n_2^{-1/2}\mathbf{u}_0 \rightarrow \theta_0$. Using (8.6), for large enough n and for all (\tilde{y}, \mathbf{x}) , we have

$$\sup_{\beta \in \mathbb{R}^p} |h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta_0}(\tilde{y}, \mathbf{x})| \leq n_2^{-1/2} \cdot K_0\|\mathbf{u}_0\| \cdot |x_1| \equiv \tilde{G}_n(\tilde{y}, \mathbf{x}).$$

Then, \tilde{G}_n can be considered as a *particular* envelope function for the class \mathcal{G}_n . Due to assumption (A3),

$$\sup_{g \in \mathcal{G}_n} \mathbf{P}_0 g^2 \leq \mathbf{P}_0 \tilde{G}_n^2 = \frac{K_0^2 \cdot \|\mathbf{u}_0\|^2}{n_2} \cdot \mathbf{E}|X_1|^2 \rightarrow 0, \quad \text{as } n \rightarrow \infty,$$

which verifies (C.1).

Before verifying (C.2), the following remarks are essential. A class of functions \mathcal{F} has a bounded uniform entropy integral (BUEI) with respect to a envelope function F (cf. Section 9.1.2 of [Kosorok \(2008\)](#) for details), if $J(1, \mathcal{F}, L_2) < \infty$. Any VC class of functions \mathcal{F} will be BUEI with respect to any choice of envelope function for \mathcal{F} . For a sequence of VC classes of functions $\{\mathcal{F}_n : n \geq 1\}$ the main goal will be to study the VC dimensions $V(\mathcal{F}_n)$ and show that it remains constant whenever n is large enough. This will ensure that for large enough n , $J(1, \mathcal{F}_n, L_2)$ remains finite and bounded by a constant and in that case, $J(\delta_n, \mathcal{F}_n, L_2) \rightarrow 0$ for any $\delta_n \downarrow 0$.

Verification of (C.2): For any $\mathbf{u}_0 \in \mathbb{R}^2$ define the classes of functions:

$$\mathcal{G}_{n,1} = \left\{ h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0} : \beta \in \mathbb{R}^p \right\} \quad \text{and} \quad \mathcal{G}_{n,2} = \{ h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p \}, \quad n \geq 1.$$

Then, $\mathcal{G}_n \subset \mathcal{G}_{n,1} - \mathcal{G}_{n,2}$. We will study the properties of the class $\mathcal{G}_{n,1}$ by considering its individual components as separate function classes.

Consider the classes of functions, $\mathcal{F}_{n,1} = \{x_1 \cdot \psi(\mathbf{x}'\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$, $n \geq 1$. Although this class is not dependent on n , but we will use this notation for ease of explanation. Following arguments given in the proof of Lemma 8.4(i), each $\mathcal{F}_{n,1}$ will be a VC class with VC dimension independent of n . Hence, $\mathcal{F}_{n,1}$ will be BUEI with respect to the envelope function $F_{n,1}(\mathbf{x}) = |x_1|$. This leads to $J(1, \mathcal{F}_{n,1}, L_2) \leq C_1$, for some $C_1 \in (0, \infty)$, for all n . For same reasons, the classes $\mathcal{F}_{n,2} = \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$, $n \geq 1$, will be BUEI with envelope function $F_{n,2}(\mathbf{x}) = 2$, for all n . Hence, $J(1, \mathcal{F}_{n,2}, L_2) \leq C_2$, for some $C_2 \in (0, \infty)$, for all n .

Now, consider the classes $\mathcal{F}_{n,3} = \{\tilde{y} - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0+n_2^{-1/2}\mathbf{u}_0}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$, $n \geq 1$. Define the class of functions,

$$\mathcal{G}_{n,3} \equiv \left\{ \left\{ 1 - \theta_{1,0} - \theta_{2,0} - n_2^{-1/2}(u_{1,0} + u_{2,0}) \right\} \cdot \psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p \right\}, \quad n \geq 1.$$

Due to assumption (A2)(ii) and since $\mathbf{u}_0 \in \mathbb{R}^2$ is fixed, there exists some $n_0 = n_0(\mathbf{u}_0) \in \mathbb{N}$, such that

$$\text{sign}(1 - \theta_{1,0} - \theta_{2,0} - n_2^{-1/2}(u_{1,0} + u_{2,0})) = \text{sign}(1 - \theta_{1,0} - \theta_{2,0}), \quad \text{for all } n \geq n_0.$$

Using the arguments given in the proof of Lemma 9.9 (vi) of Kosorok (2008), it can be seen that the VC dimension of $\{(1 - \theta_{1,0} - \theta_{2,0}) \cdot \psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\}$, will be same as the VC dimension of $\mathcal{G}_{n,3}$ for all $n \geq n_0$. The proof of Lemma 9.9 (v) of Kosorok (2008) shows that the VC dimension of

$$\begin{aligned} & \{h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0+n_2^{-1/2}\mathbf{u}_0}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p\} \\ & = \left\{ (\theta_{1,0} + n_2^{-1/2}u_{1,0}) + \left\{ 1 - \theta_{1,0} - \theta_{2,0} - n_2^{-1/2}(u_{1,0} + u_{2,0}) \right\} \cdot \psi_{\boldsymbol{\beta}}(\mathbf{x}) : \boldsymbol{\beta} \in \mathbb{R}^p \right\}, \quad n \geq n_0, \end{aligned}$$

will be same as the VC dimension of $\mathcal{G}_{n,3}$ for $n \geq n_0$. This shows that the VC dimension of $\mathcal{F}_{n,3}$ remains constant for all $n \geq n_0$ and it will be BUEI with respect to the envelope function $F_{n,3}(\tilde{y}, \mathbf{x}) = 2$. As a result, $J(1, \mathcal{F}_{n,3}, L_2) \leq C_3$, for some $C_3 \in (0, \infty)$, for all $n \geq n_0$.

The above statement about $\mathcal{F}_{n,3}$ and the arguments used in the proof of Lemma 8.4(ii) imply that, $\mathcal{F}_{n,4} \equiv \left\{ \left\{ h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0+n_2^{-1/2}\mathbf{u}_0}(\mathbf{x}) \right\}^{-1} : \boldsymbol{\beta} \in \mathbb{R}^p \right\}$ will be a VC class with a constant VC dimension for all $n \geq n_0$, and BUEI with envelope function $F_{n,4}(\mathbf{x}) = 1/\delta_1$ (or $(1 - \delta_2)^{-1}$). So, $J(1, \mathcal{F}_{n,4}, L_2) \leq C_4$, for some constant $C_4 \in (0, \infty)$, for all $n \geq n_0$. Using the same reasoning, it follows that $\mathcal{F}_{n,5} = \left\{ \left\{ 1 - h_{3,\boldsymbol{\beta},\boldsymbol{\theta}_0+n_2^{-1/2}\mathbf{u}_0}(\mathbf{x}) \right\}^{-1} : \boldsymbol{\beta} \in \mathbb{R}^p \right\}$ will be BUEI with envelope function $F_{n,5}(\mathbf{x}) = 1/(1 - \delta_1)$ (or δ_2^{-1}) and $J(1, \mathcal{F}_{n,5}, L_2) \leq C_5$, for some constant $C_5 \in (0, \infty)$, for all $n \geq n_0$.

Since products of individual BUEI classes are BUEI with envelope function equal to the product of individual envelope functions (cf. Theorem 9.15 of Kosorok (2008)), it follows that $\times_{i=1}^5 \mathcal{F}_{n,i}$ will be BUEI with envelope functions,

$$G_{n,1}(\tilde{y}, \mathbf{x}) \equiv \frac{4|x_1|}{\delta_1(1 - \delta_1)}, \quad \text{for all } n \geq n_0. \quad (8.7)$$

Since $\mathcal{G}_{n,1} \subset \times_{i=1}^5 \mathcal{F}_{n,i}$, the covering numbers satisfy, $N(\epsilon, \mathcal{G}_{n,1}, L_2(R)) \leq N(\epsilon, \times_{i=1}^5 \mathcal{F}_{n,i}, L_2(R))$, for all $\epsilon > 0$ and as a result, $\mathcal{G}_{n,1}$ will be BUEI with envelope function $G_{n,1}$, for all $n \geq n_0$. Hence,

$$J(1, \mathcal{G}_{n,1}, L_2) = \int_0^1 \sup_R \sqrt{\log N\left(\epsilon \|G_{n,1}\|_{R,2}, \mathcal{G}_{n,1}, L_2(R)\right)} d\epsilon \leq \sum_{i=1}^5 C_i \in (0, \infty), \quad \text{for all } n \geq n_0.$$

A similar arguments shows that the above statement holds if $\mathcal{G}_{n,1}$ is replaced by $\mathcal{G}_{n,2}$. As a result the class, $\mathcal{G}_n \subset \mathcal{G}_{n,1} - \mathcal{G}_{n,2}$, will be BUEI with envelope function $2 \cdot G_{n,1}$, using Lemma 9.14 (part (iii)) of [Kosorok \(2008\)](#). Hence, $J(1, \mathcal{G}_n, L_2) \leq 2 \cdot \sum_{i=1}^5 C_i < \infty$, for all $n \geq n_0$. Using DCT, $J(\delta_n, \mathcal{G}_n, L_2) \rightarrow 0$, for any $\delta_n \downarrow 0$. The class \mathcal{G}_n will be pointwise measurable by considering the countable subclass indexed over all $\beta \in \mathbb{Q}^p$. This verifies part (C.2), with envelope function $G_n = 2 \cdot G_{n,1}$.

Verification of (C.3): The envelope function G_n defined above is of the form $G_n(\tilde{y}, \mathbf{x}) = K_2 \cdot |x_1|$, for some $K_2 \in (0, \infty)$. Using assumption (A3) we obtain,

$$\mathbf{P}_0 G_n^2 = K_2^2 \cdot \mathbf{P}_0 |X_1|^2 = O(1) \quad \text{and} \quad \mathbf{P}_0 G_n^2 \cdot \mathbf{1}(G_n > \epsilon \sqrt{n}) = K_2^2 \cdot \mathbf{P}_0 X_1^2 \cdot \mathbf{1}\left(|X_1| > \frac{\epsilon \sqrt{n}}{K_2}\right) \rightarrow 0,$$

for all $\epsilon > 0$.

Hence all the three sufficient conditions required by Lemma 2.2 of [van der Vaart and Wellner \(2007\)](#) are verified. This completes the proof for the first component of $h_{2,\beta,\theta}$ and similar arguments can be used for the other components. Combining all parts we complete the proof. \square

Lemma 8.6. *Suppose, assumptions (A1)-(A4) hold. Then,*

$$\sup_{\beta \in \mathbb{R}^p} \|\mathbb{G}_{n_2}(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0})\| = o_{\mathbf{P}_0^*}(1).$$

Proof of Lemma 8.6. It will be enough to show that the above convergence result is true for the first component of $h_{2,\beta,\theta}$. We will denote the first component of $h_{2,\beta,\theta}$ by the same symbol. Note that, due to Lemma 2.1 and assumption (A4), $\sqrt{n_2} \cdot (\hat{\theta}_n - \theta_0)$ converges in distribution to a tight limiting (normal) random variable, which takes values in the σ -compact set \mathbb{R}^2 . Consider any fixed $\mathbf{u}_0 \in \mathbb{R}^2$, $\beta \in \mathbb{R}^p$ and a $\delta > 0$. Following Theorem 2.3 of [van der Vaart and Wellner \(2007\)](#), define the class of functions,

$$\mathcal{F}_n(\beta, \mathbf{u}_0, \delta) = \left\{ h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0}(\tilde{y}, \mathbf{x}) : \mathbf{u} \in \mathbb{R}^2, \|\mathbf{u} - \mathbf{u}_0\| < \delta \right\}. \quad (8.8)$$

Let, $\mathbf{u} = (u_1, u_2)'$, $\mathbf{u}_0 = (u_{1,0}, u_{2,0})' \in \mathbb{R}^2$. Using (8.5) and (8.6) we can write,

$$\sup_{\|\mathbf{u} - \mathbf{u}_0\| < \delta} |h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0}(\tilde{y}, \mathbf{x})| \leq \frac{K_1 \delta}{\sqrt{n_2}} \cdot |x_1| \equiv F_n(\beta, \mathbf{u}_0, \delta)(\tilde{y}, \mathbf{x}), \quad (8.9)$$

for some constant $K_1 \in (0, \infty)$. The function $F_n(\beta, \mathbf{u}_0, \delta)$ can be considered as a *particular* envelope function for the class $\mathcal{F}_n(\beta, \mathbf{u}_0, \delta)$. Following (8.9) and using assumption (A3), we can apply the CLT for i.i.d. random variables to obtain,

$$\sup_{\beta \in \mathbb{R}^p} |\mathbb{G}_{n_2} F_n(\beta, \mathbf{u}_0, \delta)| = \frac{K_1 \delta}{\sqrt{n_2}} \cdot \left| \frac{1}{\sqrt{n_2}} \sum_{i=n_1+1}^{n_1+n_2} (|X_{1,i}| - \mathbf{E}|X_1|) \right| = \frac{K_1 \delta}{\sqrt{n_2}} \cdot O_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_0^*}(1). \quad (8.10)$$

Finally, note that if $\delta_n \rightarrow 0$, then due to assumption (A3),

$$\sqrt{n_2} \cdot \mathbf{P}_0 F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta_n) = K_1 \delta_n \cdot \mathbf{E}|X_1| \rightarrow 0. \quad (8.11)$$

This is true uniformly for all $\boldsymbol{\beta} \in \mathbb{R}^p$ and for all $\mathbf{u}_0 \in A$, where A is any compact set in \mathbb{R}^2 . Thus, Conditions (ii) and (iii) of Theorem 2.3 of [van der Vaart and Wellner \(2007\)](#) are verified. Condition (i) of Theorem 2.3 of [van der Vaart and Wellner \(2007\)](#) is verified in Lemma 8.5. The remaining components can be handled similarly to complete the proof. \square

Lemma 8.7. *Recall the definitions of $Z_{n,1}$, $Z_{n,2}$, Z_1 and Z_2 provided in (2.6) and (2.8). Suppose, assumptions (A1)-(A4) hold. Then,*

$$(i) \sup\{\|Z_{n,1}(\boldsymbol{\beta}) - Z_1(\boldsymbol{\beta})\| : \boldsymbol{\beta} \in \mathbb{R}^p\} = o_{\mathbf{P}_0^*}(1).$$

$$(ii) \sup\{\|Z_{n,2}(\boldsymbol{\beta}) - Z_2(\boldsymbol{\beta})\| : \boldsymbol{\beta} \in \mathbb{R}^p\} = o_{\mathbf{P}_0^*}(1).$$

Proof of Lemma 8.7. We will split the proof into two parts.

(i) Note that,

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|Z_{n,1}(\boldsymbol{\beta}) - Z_1(\boldsymbol{\beta})\| = \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|(\mathbb{P}_{n_1} - \mathbf{P}_0)h_{1,\boldsymbol{\beta}}\|.$$

The proof now follows from Lemma 8.4(i) and by noting that any \mathbf{P}_0 -Donsker class will be a GC class.

(ii) Define the p -dimensional process, $\{\tilde{Z}_{n,2}(\boldsymbol{\beta}) = \mathbb{P}_{n_2} h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} : \boldsymbol{\beta} \in \mathbb{R}^p\}$. Note that,

$$\sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|\tilde{Z}_{n,2}(\boldsymbol{\beta}) - Z_2(\boldsymbol{\beta})\| = \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|(\mathbb{P}_{n_2} - \mathbf{P}_0)h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0}\| = o_{\mathbf{P}_0^*}(1),$$

using Lemma 8.4(ii). Also,

$$\begin{aligned} \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \|Z_{n,2}(\boldsymbol{\beta}) - \tilde{Z}_{n,2}(\boldsymbol{\beta})\| &= \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\| \mathbb{P}_{n_2} \left(h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} \right) \right\| \\ &\leq \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\| (\mathbb{P}_{n_2} - \mathbf{P}_0) \left(h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} \right) \right\| + \sup_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\| \mathbf{P}_0 \left(h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} \right) \right\| \\ &= \frac{1}{\sqrt{n_2}} \cdot o_{\mathbf{P}_0^*}(1) + \sup_{\boldsymbol{\beta}} \left\| \mathbf{P}_0 \left(h_{2,\boldsymbol{\beta},\hat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} \right) \right\|, \end{aligned} \quad (8.12)$$

which follows from Lemma 8.6. In order to handle the last term of (8.12), consider the first component of the p -dimensional function $h_{2,\boldsymbol{\beta},\boldsymbol{\theta}}$ and denote it by the same symbol. For any $\boldsymbol{\theta} \in \Theta$, define the function, $g(\boldsymbol{\theta}) \equiv \sup_{\boldsymbol{\beta}} |\mathbf{P}_0(h_{2,\boldsymbol{\beta},\boldsymbol{\theta}} - h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0})|$. Note that, $g(\boldsymbol{\theta}_0) = 0$. For any nonrandom sequence $\boldsymbol{\theta}_n = (\theta_{1,n}, \theta_{2,n})' \rightarrow \boldsymbol{\theta}_0$, using (8.6) it follows $g(\boldsymbol{\theta}_n) \rightarrow 0$. Hence, $g : \Theta \mapsto \mathbb{R}$ is continuous at $\boldsymbol{\theta} = \boldsymbol{\theta}_0$. Now using Lemma 2.1 and continuous mapping theorem, it follows that $g(\hat{\boldsymbol{\theta}}_n) = o_{\mathbf{P}_0^*}(1)$. Applying the same argument over other components shows that the last term of (8.12) is $o_{\mathbf{P}_0^*}(1)$. Combining this with other convergence statements given above completes the proof. \square

Lemma 8.8. *Suppose, assumptions (A1)-(A3) hold. Then, the following results are true.*

- (i) *Consider the first component of $h_{2,\beta,\theta}$ and denote it by the same symbol. Define the map, $f : \mathbb{R}^p \times \Theta \mapsto \mathbb{R}$ as, $f(\beta, \theta) = \mathbf{P}_0 h_{2,\beta,\theta}$. Let $g_i(\beta, \theta)$ denote the partial derivative of $f(\beta, \theta)$ with respect to θ_i , $i = 1, 2$. Then, $g_i(\beta, \theta)$ exists, is continuous everywhere and can be obtained by interchanging the integration and differentiation symbols:*

$$g_i(\beta, \theta) = \int \left(\frac{\partial}{\partial \theta_i} h_{2,\beta,\theta} \right) d\mathbf{P}_0, \quad i = 1, 2.$$

The result holds if any other component of $h_{2,\beta,\theta}$ is used in defining $f(\beta, \theta)$.

- (ii) *Now, use the original definition of $h_{2,\beta,\theta}$ as a p -dimensional function from (2.3). Define the map, $\bar{f} : \Theta \mapsto \mathbb{R}^p$ as, $\bar{f}(\theta) = \mathbf{P}_0 h_{2,\beta_0,\theta}$. For any $\theta \in \Theta$, the total derivative of $\bar{f}(\theta)$ is $\mathbf{A}_\theta : \Theta \mapsto \mathbb{R}^p$, a $p \times 2$ matrix, with the (i, j) -th element being $\mathbf{A}_\theta(i, j) =$ partial derivative of the i -th component of $\bar{f}(\theta)$ with respect to θ_j , $i = 1, \dots, p$ and $j = 1, 2$.*

Proof of Lemma 8.8. We will prove each part separately.

- (i) Consider any $\beta \in \mathbb{R}^p$ and any $\theta = (\theta_1, \theta_2) \in \Theta$. Consider a real valued sequence $\{h_n : n \geq 1\}$ satisfying $h_n \rightarrow 0$. Write, $\theta_n = (\theta_1 + h_n, \theta_2)$, $n \geq 1$. As because Θ is an open set the partial derivative of $h_{2,\beta,\theta}$, with respect to θ_1 exists at each $\theta \in \Theta$, with $[h_{2,\beta,\theta_n}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta}(\tilde{y}, \mathbf{x})] / h_n \rightarrow D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x})$, for all (\tilde{y}, \mathbf{x}) (cf. (8.3)). Following (8.5), $D_1(h_{2,\beta,\theta}) \in L_1(\mathbf{P}_0)$. Hence, we can use the DCT to claim that

$$\begin{aligned} g_1(\beta, \theta) &= \lim_{n \rightarrow \infty} \frac{f(\beta, \theta_n) - f(\beta, \theta)}{h_n} \\ &= \int \left(\lim_{n \rightarrow \infty} \frac{1}{h_n} [h_{2,\beta,\theta_n}(\tilde{y}, \mathbf{x}) - h_{2,\beta,\theta}(\tilde{y}, \mathbf{x})] \right) d\mathbf{P}_0(\tilde{y}, \mathbf{x}) \\ &= \int D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x}) d\mathbf{P}_0(\tilde{y}, \mathbf{x}) = \int \left(\frac{\partial}{\partial \theta_1} h_{2,\beta,\theta}(\tilde{y}, \mathbf{x}) \right) d\mathbf{P}_0(\tilde{y}, \mathbf{x}). \end{aligned}$$

This proves the existence of the partial derivatives $g_1(\beta, \theta)$ and also shows that they can be obtained by interchanging the differentiation and integration symbols. Now consider any sequence $\{(\beta_n, \theta_n) : n \geq 1\} \in \mathbb{R}^p \times \Theta$, such that $\|\beta_n - \beta\| \rightarrow 0$ and $\|\theta_n - \theta\| \rightarrow 0$, where, $(\beta, \theta) \in \mathbb{R}^p \times \Theta$. For any fixed choice of (\tilde{y}, \mathbf{x}) , the map $(\beta, \theta) \mapsto D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x})$ (cf. (8.3)) is continuous everywhere. Hence, for any (\tilde{y}, \mathbf{x}) , $D_1(h_{2,\beta_n,\theta_n})(\tilde{y}, \mathbf{x}) \rightarrow D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x})$. Combining this with (8.5) and applying the DCT along with assumption (A3) leads to,

$$\lim_{n \rightarrow \infty} g_1(\beta_n, \theta_n) = \lim_{n \rightarrow \infty} \int D_1(h_{2,\beta_n,\theta_n})(\tilde{y}, \mathbf{x}) d\mathbf{P}_0(\tilde{y}, \mathbf{x}) = \int D_1(h_{2,\beta,\theta})(\tilde{y}, \mathbf{x}) d\mathbf{P}_0(\tilde{y}, \mathbf{x}) = g_1(\beta, \theta).$$

This proves the continuity of $g_1(\beta, \theta)$ for all $(\beta, \theta) \in \mathbb{R}^p \times \Theta$. A similar argument works for $g_2(\beta, \theta)$.

- (ii) A sufficient condition for the existence of the total derivative map \mathbf{A}_θ requires the following: the partial derivatives exists for all θ and are continuous in θ (cf. Chapter 3 of van der Vaart (1998)). Following the proof of part (a) given above, this claim follows by keeping β fixed at β_0 and carrying the same arguments. The elements of \mathbf{A}_θ are given in (8.39).

□

Remark 8.9. Lemma 8.10 is a stochastic equicontinuity result about the process $\sqrt{n}(Z_n - Z)$, required for proving asymptotic normality of $\hat{\beta}_n$ and uses the consistency property of $\hat{\beta}_n$ (cf. Theorem 2.2(i)). However, the proof of consistency does not require Lemma 8.10 and is based on separate arguments, only requiring Lemma's 8.3 and 8.7.

Lemma 8.10. Suppose, assumptions (A1)-(A4). Then,

$$\|\sqrt{n}(Z_n - Z)(\hat{\beta}_n - \beta_0)\| = o_{\mathbf{P}_0^*}(1 + \sqrt{n}\|\hat{\beta}_n - \beta_0\|).$$

Proof of Lemma 8.10. It will be enough to prove, $\|\sqrt{n}(Z_n - Z)(\hat{\beta}_n - \beta_0)\| = o_{\mathbf{P}_0^*}(1)$. Using (2.8) and (2.6) we can write,

$$\begin{aligned} & \sqrt{n}(Z_n - Z)(\hat{\beta}_n - \beta_0) \\ &= \sqrt{f_n} \cdot \mathbb{G}_{n_1}(h_{1,\hat{\beta}_n} - h_{1,\beta_0}) + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot [Z_{n,2}(\hat{\beta}_n) - Z_2(\hat{\beta}_n) - Z_{n,2}(\beta_0) + Z_2(\beta_0)] \\ & \quad + \sqrt{n}(f_n - f)\mathbf{P}_0(h_{1,\hat{\beta}_n} - h_{1,\beta_0}) - \sqrt{n}(f_n - f) [Z_2(\hat{\beta}_n) - Z_2(\beta_0)] \\ & \equiv B_{1,n} + B_{2,n} + B_{3,n} + B_{4,n}. \end{aligned} \tag{8.13}$$

Note that the right side of (8.13) is p -dimensional. It is enough to show convergence of each component separately. Without loss of generality consider the first component of each term on the right side of (8.13). Also, denote the first component of $h_{1,\beta}$ and $h_{2,\beta,\theta}$ by the same symbols. Consider any non-random sequence $\{\beta_n : n \geq 1\}$ satisfying $\|\beta_n - \beta_0\| \rightarrow 0$. For each fixed (y, \mathbf{x}) , due to continuity of $\beta \mapsto h_{1,\beta}$, $|h_{1,\beta_n}(y, \mathbf{x}) - h_{1,\beta_0}(y, \mathbf{x})| \rightarrow 0$. Using assumption (A3) we have, $|h_{1,\beta_n}(y, \mathbf{x}) - h_{1,\beta_0}(y, \mathbf{x})|^2 \leq 8\|\mathbf{x}\|^2 \in L_2(\mathbf{P}_0)$, for all $n \geq 1$ and all (y, \mathbf{x}) . Define the map, $\beta \mapsto g(\beta) = \mathbf{P}_0|h_{1,\beta} - h_{1,\beta_0}|^2$. We can now apply the DCT to obtain,

$$\lim_{n \rightarrow \infty} g(\beta_n) = \lim_{n \rightarrow \infty} \int |h_{1,\beta_n} - h_{1,\beta_0}|^2 d\mathbf{P}_0 = \int \lim_{n \rightarrow \infty} |h_{1,\beta_n} - h_{1,\beta_0}|^2 d\mathbf{P}_0 = 0.$$

Since $g(\beta_0) = 0$, hence $g(\beta)$ is continuous at $\beta = \beta_0$. Using consistency of $\hat{\beta}_n$ (cf. Theorem 2.2(i)) and the continuous mapping theorem, we obtain $g(\hat{\beta}_n) = o_{\mathbf{P}_0^*}(1)$. Using the Donsker property shown in Lemma 8.4(i) and applying Lemma 19.24 of van der Vaart (1998) it follows that, $B_{1,n} = o_{\mathbf{P}_0^*}(1)$. Note that, $B_{2,n}$ can be expressed as

$$B_{2,n} = \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left[\mathbb{P}_{n_2} h_{2,\hat{\beta}_n, \hat{\theta}_n} - \mathbf{P}_0 h_{2,\hat{\beta}_n, \theta_0} - \mathbb{P}_{n_2} h_{2,\beta_0, \hat{\theta}_n} + \mathbf{P}_0 h_{2,\beta_0, \theta_0} \right]. \tag{8.14}$$

Replacing $\hat{\theta}_n$ in (8.14) by θ_0 , define

$$\begin{aligned} \tilde{B}_{2,n} &= \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left[\mathbb{P}_{n_2} h_{2,\hat{\beta}_n, \theta_0} - \mathbf{P}_0 h_{2,\hat{\beta}_n, \theta_0} - \mathbb{P}_{n_2} h_{2,\beta_0, \theta_0} + \mathbf{P}_0 h_{2,\beta_0, \theta_0} \right] \\ &= \sqrt{1-f_n} \cdot \mathbb{G}_{n_2} \left(h_{2,\hat{\beta}_n, \theta_0} - h_{2,\beta_0, \theta_0} \right), \end{aligned} \tag{8.15}$$

Then, it can be shown that

$$\begin{aligned}
|B_{2,n} - \tilde{B}_{2,n}| &\leq 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta} \left| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| \\
&\quad + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left| \mathbf{P}_0 \left(h_{2,\hat{\beta}_n,\hat{\theta}_n} - h_{2,\hat{\beta}_n,\theta_0} - h_{2,\beta_0,\hat{\theta}_n} + h_{2,\beta_0,\theta_0} \right) \right| \\
&\equiv C_{1,n} + C_{2,n}.
\end{aligned}$$

Due to Lemma 8.6, $C_{1,n} = o_{\mathbf{P}_0^*}(1)$. Note that, for any sequence of non-random constants $\{a_n : n \geq 1\} \in (0, 1)$, the sequence, $a_n \cdot \hat{\theta}_n + (1-a_n) \cdot \theta_0 \xrightarrow{\mathbf{P}_0^*} \theta_0$, using Lemma 2.1. Using Lemma 8.8(i) and Mean-Value theorem, we can write,

$$\begin{aligned}
C_{2,n} &= \sqrt{1-f_n} \\
&\quad \times \sqrt{n_2} (\hat{\theta}_n - \theta_0)' \begin{pmatrix} g_1(\hat{\beta}_n, a_{1,n}\hat{\theta}_n + (1-a_{1,n})\theta_0) - g_1(\beta_0, a_{2,n}\hat{\theta}_n + (1-a_{2,n})\theta_0) \\ g_2(\hat{\beta}_n, b_{1,n}\hat{\theta}_n + (1-b_{1,n})\theta_0) - g_2(\beta_0, b_{2,n}\hat{\theta}_n + (1-b_{2,n})\theta_0) \end{pmatrix}, \quad (8.16)
\end{aligned}$$

for some sequences $\{a_{j,n}\}, \{b_{j,n}\} \in (0, 1)$, for $j = 1, 2$, and $g_j, j = 1, 2$, are the partial derivative functions defined in Lemma 8.8(i). Using Lemma 2.1, Theorem 2.2(i), Lemma 8.8(i) and the continuous mapping theorem, we obtain,

$$\begin{aligned}
g_1(\hat{\beta}_n, a_{1,n}\hat{\theta}_n + (1-a_{1,n})\theta_0) &\xrightarrow{\mathbf{P}_0^*} g_1(\beta_0, \theta_0), \\
g_1(\beta_0, a_{2,n}\hat{\theta}_n + (1-a_{2,n})\theta_0) &\xrightarrow{\mathbf{P}_0^*} g_1(\beta_0, \theta_0).
\end{aligned}$$

This implies, $g_1(\hat{\beta}_n, a_{1,n}\hat{\theta}_n + (1-a_{1,n})\theta_0) - g_1(\beta_0, a_{2,n}\hat{\theta}_n + (1-a_{2,n})\theta_0) = o_{\mathbf{P}_0^*}(1)$. For same reasons, $g_2(\hat{\beta}_n, b_{1,n}\hat{\theta}_n + (1-b_{1,n})\theta_0) - g_2(\beta_0, b_{2,n}\hat{\theta}_n + (1-b_{2,n})\theta_0) = o_{\mathbf{P}_0^*}(1)$. Applying Lemma 2.1, we get $C_{2,n} = O_{\mathbf{P}_0^*}(1) \cdot o_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_0^*}(1)$, which implies $|B_{2,n} - \tilde{B}_{2,n}| = o_{\mathbf{P}_0^*}(1)$. Note that, $\{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$ is a Donsker class, following Lemma 8.4(ii). So we can apply the same arguments as in the case of $B_{1,n}$ to show that $\mathbf{P}_0(h_{2,\hat{\beta}_n,\theta_0} - h_{2,\beta_0,\theta_0})^2 = o_{\mathbf{P}_0^*}(1)$. Using Lemma 19.24 of van der Vaart (1998) it follows that $\tilde{B}_{2,n} = o_{\mathbf{P}_0^*}(1)$ and hence $B_{2,n} = o_{\mathbf{P}_0^*}(1)$.

Finally, due to assumption (A3) and the DCT, $\beta \mapsto \mathbf{P}_0 h_{1,\beta}$ is continuous everywhere. Applying the continuous mapping theorem, $\mathbf{P}_0 h_{1,\hat{\beta}_n} - \mathbf{P}_0 h_{1,\beta_0} = o_{\mathbf{P}_0^*}(1)$ and by assumption (A4), $\sqrt{n}|f_n - f| = o(1)$. Hence, $B_{3,n} = o_{\mathbf{P}_0^*}(1)$. Similarly, using Lemma 8.8(i) and the stated assumptions we can claim that, $\beta \mapsto \mathbf{P}_0 h_{2,\beta,\theta_0}$ will be continuous at all $\beta \in \mathbb{R}^p$. Similar arguments imply $B_{4,n} = o_{\mathbf{P}_0^*}(1)$. Combining the results above shows that the right side (8.13) is $o_{\mathbf{P}_0^*}(1)$. Applying this over all p components separately shows that $\|\sqrt{n}(Z_n - Z)(\hat{\beta}_n - \beta_0)\| = o_{\mathbf{P}_0^*}(1)$. This completes the proof. \square

8.2 Auxiliary lemmas required for proving Theorem 3.2

We begin with a result which describes the inter-relations among the different probability orders and will frequently be used in the proofs. Extensions to Lemma 8.11 can be found in Lemma 3 of Cheng and

Huang (2010). Further details on outer modes of convergence can be found in van der Vaart and Wellner (1996). Consider any class of measurable (vector valued) functions \mathcal{H} . For any measure \mathbf{R} , we define, $\|\mathbf{R}\|_{\mathcal{H}} \equiv \sup_{h \in \mathcal{H}} \|\mathbf{R}h\|$.

Lemma 8.11. *Consider a sequence of functions $\{\Delta_n : n \geq 1\}$.*

- (a) *Suppose, $\Delta_n = o_{\mathbf{Pr}^*}(1)$. Then, $\Delta_n = o_{\mathbf{P}_0^*}(1)$ in \mathbf{P}_0^* -probability.*
- (b) *Suppose, $\{\Delta_n : n \geq 1\}$ is a sequence of functions defined only the first product probability space $(\mathcal{X}^\infty, \mathcal{A}^\infty, \mathbf{P}_0^\infty)$ and $\Delta_n = o_{\mathbf{P}_0^*}(1)$. Then, $\Delta_n = o_{\mathbf{Pr}^*}(1)$ and $\Delta_n = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.*
- (c) *Suppose, $\Delta_n = O_{\mathbf{Pr}^*}(1)$. Then, $\Delta_n = O_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.*

Proof of Lemma 8.11. We consider the three parts separately.

- (a) Fix any $\epsilon, \eta > 0$. Then, using Markov's inequality and Fubini's theorem for product measures (cf. Lemma 1.2.6 of van der Vaart and Wellner (1996)),

$$\mathbf{P}_0^*(\mathbf{P}_M^*(|\Delta_n| > \epsilon) > \eta) \leq \frac{1}{\eta} \cdot \mathbf{E}_0^*[\mathbf{P}_M^*(|\Delta_n| > \epsilon)] \leq \frac{1}{\eta} \cdot \mathbf{Er}^*\mathbf{1}(|\Delta_n| > \epsilon) = \frac{1}{\eta} \cdot \mathbf{Pr}^*(|\Delta_n| > \epsilon) \rightarrow 0,$$

by definition of outer convergence in probability in \mathbf{Pr} . This completes the proof.

- (b) In this case,

$$\mathbf{Pr}^*(|\Delta_n| > \epsilon) = \mathbf{Pr}(|\Delta_n| > \epsilon)^* = \mathbf{P}_0(|\Delta_n| > \epsilon)^* = \mathbf{P}_0^*(|\Delta_n| > \epsilon) \rightarrow 0, \quad \text{as } n \rightarrow \infty.$$

Since Δ_n depends only the product space $(\mathcal{X}^\infty, \mathcal{A}^\infty, \mathbf{P}_0^\infty)$ so the outer majorant with respect to \mathbf{Pr} is same as the outer majorant with respect to \mathbf{P}_0^∞ (cf. page 10 of van der Vaart and Wellner (1996)). Since, $\Delta_n = o_{\mathbf{Pr}^*}(1)$, we can apply part (a) above to obtain, $\Delta_n = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. In case, $\Delta_n = O_{\mathbf{Pr}^*}(1)$, the a very similar argument implies that $\Delta_n = O_{\mathbf{Pr}^*}(1)$.

- (c) For any $\delta > 0$ and $M > 0$, we can write

$$\mathbf{P}_0^*(\mathbf{P}_M^*(|\Delta_n| > M) > \delta) \leq \frac{1}{\delta} \cdot \mathbf{E}_0^*\mathbf{E}_M^*\mathbf{1}(|\Delta_n| > M) \leq \frac{1}{\delta} \cdot \mathbf{Er}^*\mathbf{1}(|\Delta_n| > M) \leq \frac{1}{\delta} \cdot \delta\epsilon = \epsilon,$$

provided we choose $M = M(\epsilon, \delta)$ and $n \geq n_0(\epsilon, \delta)$, such that $\mathbf{Pr}^*(|\Delta_n| > M(\epsilon, \delta)) < \delta\epsilon$, for all $n \geq n_0(\epsilon, \delta)$. The latter condition will be true, because $\Delta_n = O_{\mathbf{Pr}^*}(1)$. This completes the proof. □

Remark 8.12. *Lemma's 8.13, 8.14, 8.15 and 8.16 are needed for studying the asymptotic distribution of the bootstrapped PMLE $\hat{\hat{\beta}}_n$ and uses consistency of $\hat{\hat{\beta}}_n$. The proof of bootstrap consistency of $\hat{\hat{\beta}}_n$ is not dependent on these lemmas.*

Lemma 8.13. *Suppose, assumptions (A1)-(A4) holds. Then,*

- (i) $\mathbb{G}_{n_1} \left(h_{1, \hat{\beta}_n} - h_{1, \beta_0} \right) = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.
- (ii) $\mathbb{G}_{n_2} \left(h_{2, \hat{\beta}_n, \theta_0} - h_{2, \beta_0, \theta_0} \right) = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.
- (iii) $\sup \left\{ \mathbb{G}_{n_2} \left(h_{2, \beta, \hat{\theta}_n} - h_{2, \beta, \theta_0} \right) : \beta \in \mathbb{R}^p \right\} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.

Proof of Lemma 8.13. We will prove each part separately.

- (i) Without loss of generality consider the first component of $h_{1, \beta}$ and denote it by the same symbol. For any $f, g \in L_2(\mathbf{P}_0)$, define the semi-metric $\rho_2(f, g) \equiv (\mathbf{Var}_{\mathbf{P}_0}(f - g))^{1/2}$. Note that, $h_{1, \beta} \in L_2(\mathbf{P}_0)$ for all β . We claim that there exists a sequence $\delta_n \downarrow 0$, such that

$$\mathbf{P}_M^* \left(\int |h_{1, \hat{\beta}_n} - h_{1, \beta_0}|^r d\mathbf{P}_0 > \delta_n \right) = o_{\mathbf{P}_0^*}(1), \quad \text{for } r = 1, 2. \quad (8.17)$$

The conclusion for $r = 2$ in (8.17) follows from consistency of $\hat{\beta}_n$ in Theorem 3.2(i), the arguments used in handling the term $B_{1, n}$ in (8.13) and finally using Lemma 8.18. The case of $r = 1$ follows from the case of $r = 2$ (or it can be proved directly). It now follows that there exists a sequence $\delta_n \downarrow 0$, such that

$$\mathbf{P}_M^* \left(\rho_2 \left(h_{1, \hat{\beta}_n}, h_{1, \beta_0} \right) > \delta_n \right) = o_{\mathbf{P}_0^*}(1). \quad (8.18)$$

Based on the above chosen sequence $\{\delta_n : n \geq 1\}$, define the classes of functions

$$\mathcal{H}_n = \{h_{1, \beta} - h_{1, \beta_0} : \rho_2(h_{1, \beta}, h_{1, \beta_0}) \leq \delta_n\}, \quad n \geq 1. \quad (8.19)$$

Following Lemma 8.4(i), it is known that $\{h_{1, \beta} : \beta \in \mathbb{R}^p\}$ is a Donsker class. Applying Corollary 2.3.12 of van der Vaart and Wellner (1996), it follows that, $\|\mathbb{G}_{n_1}\|_{\mathcal{H}_n} = o_{\mathbf{P}_0^*}(1)$, which in turn implies, $\|\mathbb{G}_{n_1}\|_{\mathcal{H}_n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, using Lemma 8.11. Fix any $\epsilon > 0$. Then,

$$\begin{aligned} & \mathbf{P}_M^* \left(\left| \mathbb{G}_{n_1} \left(h_{1, \hat{\beta}_n} - h_{1, \beta_0} \right) \right| > \epsilon \right) \\ & \leq \mathbf{P}_M^* \left(\left| \mathbb{G}_{n_1} \left(h_{1, \hat{\beta}_n} - h_{1, \beta_0} \right) \right| > \epsilon, \rho_2 \left(h_{1, \hat{\beta}_n}, h_{1, \beta_0} \right) \leq \delta_n \right) + \mathbf{P}_M^* \left(\rho_2 \left(h_{1, \hat{\beta}_n}, h_{1, \beta_0} \right) > \delta_n \right) \\ & \leq \mathbf{P}_M^* \left(\|\mathbb{G}_{n_1}\|_{\mathcal{H}_n} > \epsilon, \rho_2 \left(h_{1, \hat{\beta}_n}, h_{1, \beta_0} \right) \leq \delta_n \right) + o_{\mathbf{P}_0^*}(1) \\ & \leq \mathbf{P}_M^* \left(\|\mathbb{G}_{n_1}\|_{\mathcal{H}_n} > \epsilon \right) + o_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_0^*}(1). \end{aligned}$$

We can apply this argument to all other components of $h_{1, \beta}$ separately to complete the proof.

- (ii) The proof proceeds in exactly the same manner as part (i) above and using Lemma 8.4(ii), which states that $\{h_{2, \beta, \theta_0} : \beta \in \mathbb{R}^p\} \subset L_2(\mathbf{P}_0)$, is a Donsker class.
- (iii) Using Lemma 3.1(ii) and assumption (A4), we know $\sqrt{n_2} \cdot (\hat{\theta}_n - \theta_0)$ converges to a tight limiting distribution (with respect to \mathbf{Pr}), hence we can apply Lemma 8.6 (with all probability statements

understood in terms of \mathbf{Pr}). Thus, for any $\epsilon > 0$,

$$\mathbf{Pr}^* \left(\sup_{\beta} \left\| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right\| > \epsilon \right) = o(1).$$

Applying Lemma 8.11 completes the proof. □

Lemma 8.14. *Suppose, assumptions (A1)-(A4) holds. Then,*

- (i) $\widehat{\mathbb{G}}_{n_1} \left(h_{1,\hat{\beta}_n} - h_{1,\beta_0} \right) = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.
- (ii) $\sup \left\{ \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) : \beta \in \mathbb{R}^p \right\} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.
- (iii) $\sup \left\{ \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) : \beta \in \mathbb{R}^p \right\} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability.

Proof of Lemma 8.14. We will prove each part separately. As earlier, we will work with the first components of $h_{1,\beta}$ and $h_{2,\beta,\theta}$ and denote them by the same symbol.

- (i) Note that, $\widehat{\mathbb{G}}_{n_1}$ remains unchanged if $h_{1,\beta}$ is re-centered at $h_{1,\beta} - \mathbf{P}_0 h_{1,\beta}$. Hence, for the rest of the proof of part (i), we will work with the re-centered versions and assume that $\mathbf{P}_0 h_{1,\beta} = 0$ for all β . Following Markov's inequality, it will be enough to show $\mathbf{E}_M^* \left| \widehat{\mathbb{G}}_{n_1} \left(h_{1,\hat{\beta}_n} - h_{1,\beta_0} \right) \right| = o_{\mathbf{P}_0^*}(1)$. Note that, the class \mathcal{H}_n defined in (8.19) is not affected, since the ρ_2 semi-metric uses centered versions of $h_{1,\beta}$. Fix any $\epsilon > 0$. Following (8.18) we can write,

$$\begin{aligned} & \mathbf{P}_0^* \left(\mathbf{E}_M^* \left| \widehat{\mathbb{G}}_{n_1} \left(h_{1,\hat{\beta}_n} - h_{1,\beta_0} \right) \right| > \epsilon \right) \\ & \leq \mathbf{P}_0^* \left(\mathbf{E}_M^* \left| \widehat{\mathbb{G}}_{n_1} \left(h_{1,\hat{\beta}_n} - h_{1,\beta_0} \right) \right| > \epsilon, \rho_2 \left(h_{1,\hat{\beta}_n}, h_{1,\beta_0} \right) \leq \delta_n \right) + o_{\mathbf{P}_0^*}(1) \\ & \leq \mathbf{P}_0^* \left(\mathbf{E}_M^* \left\| \widehat{\mathbb{G}}_{n_1} \right\|_{\mathcal{H}_n} > \epsilon \right) + o_{\mathbf{P}_0^*}(1) \leq \frac{1}{\epsilon} \cdot \mathbf{E}_0^* \left(\mathbf{E}_M^* \left\| \widehat{\mathbb{G}}_{n_1} \right\|_{\mathcal{H}_n} \right) + o_{\mathbf{P}_0^*}(1). \end{aligned}$$

Since, $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ is a Donsker class, we can use the arguments given in the proof of Theorem 9.3 of Dudley (2014) (cf. pg 344) to conclude that

$$\limsup_{n \rightarrow \infty} \mathbf{E}_0^* \left(\mathbf{E}_M^* \left\| \widehat{\mathbb{G}}_{n_1} \right\|_{\mathcal{H}_n} \right) = 0.$$

- (ii) We will adapt the proof of Theorem 2.3 of van der Vaart and Wellner (2007) in this case. For any $\delta > 0$ and any compact set K , define $K^\delta = \{\mathbf{y} \in \mathbb{R}^2 : d(\mathbf{y}, K) < \delta\}$, where $d(\mathbf{y}, K) = \inf\{d(\mathbf{y}, \mathbf{x}) : \mathbf{x} \in K\}$ and d denotes the Euclidean distance in \mathbb{R}^2 . Since K is compact, for any $\delta > 0$ there exists finitely many points $\{\mathbf{u}_1, \dots, \mathbf{u}_p\} \in K$, with $p = p(\delta)$ such that $K \subset \bigcup_{i=1}^p B(\mathbf{u}_i, \delta/2)$, where $B(\mathbf{u}, \delta)$ denotes an open ball of radius δ around the point $\mathbf{u} \in \mathbb{R}^2$. Hence, $K^\delta \subset \bigcup_{i=1}^p B(\mathbf{u}_i, \delta)$. Note that,

$$\left\{ \sqrt{n_2} (\widehat{\theta}_n - \theta_0) \in K^{\delta/2} \right\} \Rightarrow \bigcup_{i=1}^p \left\{ \widehat{\theta}_n \in B(\theta_0 + n_2^{-1/2} \mathbf{u}_i, \delta) \right\}.$$

Fix any $\eta > 0$. Then, for any $\delta > 0$,

$$\begin{aligned}
& \mathbf{P}_M^* \left(\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\widehat{\boldsymbol{\theta}}_n} - h_{2,\beta,\boldsymbol{\theta}_0} \right) \right| > \eta \right) \\
& \leq \mathbf{P}_M^* \left(\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\widehat{\boldsymbol{\theta}}_n} - h_{2,\beta,\boldsymbol{\theta}_0} \right) \right| > \eta, \sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \in K^{\delta/2} \right) + \mathbf{P}_M^* \left(\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \notin K^{\delta/2} \right) \\
& \leq \mathbf{P}_M^* \left(\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\widehat{\boldsymbol{\theta}}_n} - h_{2,\beta,\boldsymbol{\theta}_0} \right) \right| > \eta, \sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \in K^{\delta/2} \right) \\
& \quad + \mathbf{P}_M^* \left(\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \notin K \right). \tag{8.20}
\end{aligned}$$

From Lemma 3.1(iii), it follows $\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) = O_{\mathbf{P}_M^*}(1)$. Fix any arbitrary $\epsilon_1, \epsilon_2 > 0$. Then, following Lemma 8.11, there exists a compact set $K_{\epsilon_1, \epsilon_2}$ and an integer $n_0 = n_0(\epsilon_1, \epsilon_2)$, such that

$$\mathbf{P}_0^* \left(\mathbf{P}_M^* \left(\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \notin K_{\epsilon_1, \epsilon_2} \right) > \epsilon_2 \right) \leq \epsilon_1, \quad \text{for all } n \geq n_0. \tag{8.21}$$

Use this choice of $K = K_{\epsilon_1, \epsilon_2}$ in the right side of (8.20). Following the proof of Theorem 2.3 of [van der Vaart and Wellner \(2007\)](#) and using the same choice of $F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta)$ (cf. (8.9)), the first term on the right side of (8.20) with $K = K_{\epsilon_1, \epsilon_2}$, can be written as

$$\begin{aligned}
& \mathbf{P}_M^* \left(\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\widehat{\boldsymbol{\theta}}_n} - h_{2,\beta,\boldsymbol{\theta}_0} \right) \right| > \eta, \sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0) \in K_{\epsilon_1, \epsilon_2}^{\delta/2} \right) \\
& \leq \mathbf{P}_M^* \left(\sup_{\beta} \max_{1 \leq i \leq p} \left| \widehat{\mathbb{G}}_{n_2} F_n(\boldsymbol{\beta}, \mathbf{u}_i, \delta) \right| > \eta/3 \right) + \mathbf{P}_M^* \left(2 \sup_{\beta} \sup_{\mathbf{u}_0 \in K_{\epsilon_1, \epsilon_2}} \sqrt{n_2} \cdot \mathbb{P}_{n_2} F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta) > \eta/3 \right) \\
& \quad + \mathbf{P}_M^* \left(\sup_{\beta} \max_{1 \leq i \leq p} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\boldsymbol{\theta}_0 + n_2^{-1/2} \mathbf{u}_i} - h_{2,\beta,\boldsymbol{\theta}_0} \right) \right| > \eta/3 \right) \\
& \equiv L_{1,n} + L_{2,n} + L_{3,n}. \tag{8.22}
\end{aligned}$$

Consider the first term of (8.22). Since $\delta > 0$ is a fixed quantity and ϵ_1, ϵ_2 are arbitrary but *fixed* positive constants, hence $p = p(\delta)$ is finite. Due to finiteness of p , showing $L_{1,n} = o_{\mathbf{P}_0^*}(1)$, is equivalent to showing that $\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta) \right| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability for any fixed $\mathbf{u}_0 \in \mathbb{R}^2$. Following (8.10), we can write

$$\begin{aligned}
\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta) \right| &= K_1 \cdot \delta \cdot \left| \left(\widehat{\mathbb{P}}_{n_2} - \mathbb{P}_{n_2} \right) |x_1| \right| \\
&= K_1 \cdot \delta \cdot \left| \frac{1}{n_2} \sum_{i=n_1+1}^n \left(|\widehat{X}_{1,i}| - \mathbf{E}|X_1| \right) - \frac{1}{n_2} \sum_{i=n_1+1}^n \left(|X_{1,i}| - \mathbf{E}|X_1| \right) \right|.
\end{aligned}$$

Under the stated assumptions, $n_2^{-1} \sum_{i=n_1+1}^n \left(|\widehat{X}_{1,i}| - \mathbf{E}|X_1| \right) = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Also by Lemma 8.11, $n_2^{-1} \sum_{i=n_1+1}^n \left(|X_{1,i}| - \mathbf{E}|X_1| \right) = o_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_M^*}(1) = o_{\mathbf{P}_R^*}(1)$ in \mathbf{P}_0^* -probability.

Now, consider the second term of (8.22). The definition of $F_n(\boldsymbol{\beta}, \mathbf{u}_0, \delta)$ in (8.9) shows that it is independent of $\boldsymbol{\beta}$ and \mathbf{u}_0 . Hence, for any $\delta_n \rightarrow 0$, using the WLLN, assumption (A3) and (8.11) we

have,

$$\sup_{\beta \in \mathbb{R}^p} \sup_{\mathbf{u} \in \mathbb{R}^2} \sqrt{n_2} \cdot \mathbb{P}_{n_2} F_n(\beta, \mathbf{u}_0, \delta_n) = K_1 \delta_n \cdot \left(n_2^{-1} \sum_{i=n_1+1}^n |X_{1,i}| \right) = o_{\mathbf{P}_0^*}(1).$$

Hence, $L_{2,n} = o_{\mathbf{P}_0^*}(1)$.

Finally, consider the third term of (8.22). Due to finiteness of p , it enough to show that for any $\mathbf{u}_0 \in \mathbb{R}^2$, $\|\widehat{\mathbb{G}}_{n_2}\|_{\mathcal{G}_n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, where, \mathcal{G}_n is defined in (8.2). Let $\{N_i : i \geq n_1 + 1\}$ be a sequence of i.i.d. Poisson(1) random variables, independent of $\{(\tilde{Y}_i, \mathbf{X}_i) : i \geq n_1 + 1\}$ and defined on a different probability space. Let \mathbf{E}^* denote the outer expectation with respect to the product probability measure corresponding to the product space $(\mathcal{X}^\infty, \mathcal{A}^\infty, \mathbf{P}_0^\infty) \times$ the probability space of N_i 's. Following the arguments given in Dudley (2014) (Lemma 9.12 and pg. 344), we can write

$$\begin{aligned} \mathbf{E}_0^* \mathbf{E}_M^* \left(\|\widehat{\mathbb{G}}_{n_2}\|_{\mathcal{G}_n} \right) &= n_2^{-1/2} \cdot \mathbf{E}_0^* \mathbf{E}_M^* \left(\left\| \sum_{i=n_1+1}^n \left(\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbb{P}_{n_2} \right) \right\|_{\mathcal{G}_n} \right) \\ &\leq n_2^{-1/2} \cdot \frac{e}{e-1} \cdot \mathbf{E}^* \left\| \sum_{i=n_1+1}^n (N_i - 1) \cdot \left(\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbb{P}_{n_2} \right) \right\|_{\mathcal{G}_n} \\ &\leq n_2^{-1/2} \cdot \frac{e}{e-1} \cdot \mathbf{E}^* \left\| \sum_{i=n_1+1}^n (N_i - 1) \cdot \delta_{(\tilde{Y}_i, \mathbf{X}_i)} \right\|_{\mathcal{G}_n} + \frac{e}{e-1} \cdot \mathbf{E}_0^* \|\mathbb{P}_{n_2}\|_{\mathcal{G}_n} \\ &= n_2^{-1/2} \cdot \frac{e}{e-1} \cdot \mathbf{E}^* \left\| \sum_{i=n_1+1}^n (N_i - 1) \cdot \delta_{(\tilde{Y}_i, \mathbf{X}_i)} \right\|_{\mathcal{G}_n} + \frac{e}{e-1} \cdot \mathbf{E}_0^* \|n_2^{-1/2} \cdot \mathbb{G}_{n_2} + \mathbf{P}_0\|_{\mathcal{G}_n} \\ &\leq n_2^{-1/2} \cdot \frac{e}{e-1} \cdot \mathbf{E}^* \left\| \sum_{i=n_1+1}^n (N_i - 1) \cdot \delta_{(\tilde{Y}_i, \mathbf{X}_i)} \right\|_{\mathcal{G}_n} + n_2^{-1/2} \cdot \frac{e}{e-1} \cdot \mathbf{E}_0^* \|\mathbb{G}_{n_2}\|_{\mathcal{G}_n} + \frac{e}{e-1} \cdot \|\mathbf{P}_0\|_{\mathcal{G}_n} \\ &\equiv P_{1,n} + P_{2,n} + P_{3,n}. \end{aligned} \tag{8.23}$$

We will handle each term of (8.23) separately.

Starting with the third term of (8.23), note that for large enough n_2 and fixed $\mathbf{u}_0 = (u_{1,0}, u_{2,0})'$, $\boldsymbol{\theta}_0 + n_2^{-1/2} \mathbf{u}_0 \in \Theta$. Following (8.1) and (8.5) we can write, for some constant $K_0 \in (0, \infty)$,

$$P_{3,n} \leq \sum_{k=1}^2 \frac{|u_{k,0}|}{\sqrt{n_2}} \cdot K_0 \cdot \mathbf{E}|X_1| \rightarrow 0,$$

due to assumption (A3).

Next, consider the second term of (8.23). Note that, as per the proof of Lemma 8.5, for any $\delta > 0$, there exists some $n_0 \in \mathbb{N}$, such that $\sup_{g \in \mathcal{G}_n} \mathbf{P}_0 g^2 \leq \delta^2$, for all $n \geq n_0$. Now we can apply the sufficient conditions stated in the proof Lemma 8.5 and arguments given in Theorem 6.16 of van der Vaart (2002) (pg. 405) to show that $\mathbf{E}_0^* \|\mathbb{G}_{n_2}\|_{\mathcal{G}_n} \rightarrow 0$. Hence, $P_{2,n} \rightarrow 0$.

Finally, consider the first term of (8.23). Define i.i.d. Radamacher random variables $\{\epsilon_i : i \geq 1\}$ on a different factor of a probability space, which are independent of $\{N_i : i \geq 1\}$ and $\{(\tilde{Y}_i, \mathbf{X}_i) : i \geq 1\}$ and

denote the outer expectation with respect to the product measure of $\{\epsilon_i : i \geq 1\}$ and $\{(\tilde{Y}_i, \mathbf{X}_i) : i \geq n_1 + 1\}$ as¹ \mathbf{E}^* . In order to handle $P_{1,n}$, for convenience of notation we re-index the indices, so that the index set is $\{1, \dots, n_2\}$ instead of $\{n_1 + 1, \dots, n_1 + n_2\}$. Firstly,

$$\begin{aligned} P_{1,n} &= n_2^{-1/2} \mathbf{E}^* \left\| \sum_{i=n_1+1}^n (N_i - 1) \cdot \delta_{(\tilde{Y}_i, \mathbf{X}_i)} \right\|_{\mathcal{G}_n} \\ &\leq n_2^{-1/2} \cdot \mathbf{E}^* \left\| \sum_{i=1}^{n_2} (N_i - 1) \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} + n_2^{-1/2} \cdot \mathbf{E} \left| \sum_{i=1}^{n_2} (N_i - 1) \right| \cdot \|\mathbf{P}_0\|_{\mathcal{G}_n} \\ &\leq n_2^{-1/2} \cdot \mathbf{E}^* \left\| \sum_{i=1}^{n_2} (N_i - 1) \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} + \sum_{k=1}^2 \frac{|u_{k,0}|}{\sqrt{n_2}} \cdot K_0 \cdot \mathbf{E}|X_1|. \end{aligned}$$

The second term in the right side converges to zero for any fixed \mathbf{u}_0 . Note that $\{N_i - 1 : i \geq 1\}$ are centered i.i.d. random variables, independent of $\{(\tilde{Y}_i, \mathbf{X}_i) : i \geq 1\}$. Define, $\|(N_1 - 1)\|_{2,1} = \int_0^\infty \sqrt{\mathbf{P}(|N_1 - 1| > t)} dt$. Also, $\{\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0 : 1 \leq i \leq n_2\}$ are i.i.d. stochastic processes with finite expectation over \mathcal{G}_n , hence

$$\begin{aligned} &\mathbf{E}_0^* \|\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0\|_{\mathcal{G}_n} \\ &= \mathbf{E}_0^* \sup_{\beta} \left| h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_0}(\tilde{Y}_i, \mathbf{X}_i) - h_{2,\beta,\theta_0}(\tilde{Y}_i, \mathbf{X}_i) \right| + \|\mathbf{P}_0\|_{\mathcal{G}_n} \leq \frac{2K_0\|\mathbf{u}_0\|}{\sqrt{n_2}} \cdot \mathbf{E}|X_1| < \infty, \end{aligned}$$

whenever n is large enough, for any fixed \mathbf{u}_0 (cf. (8.5)). Now, $\mathbf{E}|N_1 - 1|^r < \infty$, for all $r > 2$. Hence, $\|(N_1 - 1)\|_{2,1} < \infty$. Also, $\mathbf{E}|N_1 - 1|^2 < \infty$, implies, $\mathbf{E}(\max_{1 \leq i \leq n_2} |N_i - 1|/\sqrt{n_2}) \rightarrow 0$. Then we can apply the Multiplier inequality given in Lemma 2.9.1 of [van der Vaart and Wellner \(1996\)](#) to obtain,

$$\begin{aligned} &\mathbf{E}^* \left\| \frac{1}{\sqrt{n_2}} \sum_{i=1}^{n_2} (N_i - 1) \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} \\ &\leq 2(n_0 - 1) \cdot \mathbf{E}_0^* \|\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0\|_{\mathcal{G}_n} \cdot \mathbf{E} \left(\frac{\max_{1 \leq i \leq n_2} |N_i - 1|}{\sqrt{n_2}} \right) \\ &\quad + 2\sqrt{2} \|(N_1 - 1)\|_{2,1} \cdot \max_{n_0 \leq k \leq n_2} \mathbf{E}^* \left\| \frac{1}{\sqrt{k}} \sum_{i=n_0}^k \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} \\ &\leq \frac{4(n_0 - 1)K_0\|\mathbf{u}_0\|}{\sqrt{n_2}} \cdot \mathbf{E} \left(\frac{\max_{1 \leq i \leq n_2} |N_i - 1|}{\sqrt{n_2}} \right) \\ &\quad + 2\sqrt{2} \|(N_1 - 1)\|_{2,1} \cdot \max_{n_0 \leq k \leq n_2} \mathbf{E}^* \left\| \frac{1}{\sqrt{k}} \sum_{i=1}^k \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n}. \end{aligned} \tag{8.24}$$

Following earlier arguments, for any sequence $n_0 \rightarrow \infty$, satisfying, $n_0/\sqrt{n_2} \rightarrow 0$, the first term on the

¹We have used the same notation to denote outer expectation with respect to the probability measure corresponding to the product space of $\{(\tilde{Y}_i, \mathbf{X}_i) : i \geq n_1 + 1\}$ and $\{N_i : i \geq n_1 + 1\}$, but this will not create any confusion as the context will be clear each time.

right side of (8.24) converges to zero as $n_2 \rightarrow \infty$. And

$$\begin{aligned} & \max_{n_0 \leq k \leq n_2} \mathbf{E}^* \left\| \frac{1}{\sqrt{k}} \sum_{i=n_0}^k \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} \\ & \leq \max_{n_0 \leq k \leq n_2} \mathbf{E}^* \left\| \frac{1}{\sqrt{k}} \sum_{i=1}^k \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} + \mathbf{E}^* \left\| \frac{1}{\sqrt{n_0}} \sum_{i=1}^{n_0} \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n}. \end{aligned} \quad (8.25)$$

For any integer m , using Lemma 2.3.6 of [van der Vaart and Wellner \(1996\)](#), we can remove the Radamacher random variables, and obtain

$$\mathbf{E}^* \left\| \frac{1}{\sqrt{m}} \sum_{i=1}^m \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} \leq 2 \cdot \mathbf{E}_0^* \left\| \frac{1}{\sqrt{m}} \sum_{i=1}^m (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} = 2 \cdot \mathbf{E}_0^* \|\mathbb{G}_m\|_{\mathcal{G}_n},$$

where, \mathbb{G}_m is the analogue of \mathbb{G}_{n_2} , based on the i.i.d. sample $\{(\tilde{Y}_i, \mathbf{X}_i) : 1 \leq i \leq m\}$. Now, use $m = n_0$ and let $n_0 \rightarrow \infty$. It can be verified that the sufficient conditions for Lemma 8.5 continue to hold if we replace \mathbb{G}_{n_2} by \mathbb{G}_{n_0} in the statement of Lemma 8.5. Following the arguments used in handling $P_{2,n}$ and removing the Radamacher random variables in the above manner, we can claim that $\mathbf{E}_0^* \|\mathbb{G}_{n_0}\|_{\mathcal{G}_n} \rightarrow 0$. In order to handle the first term on the right side of (8.25), firstly note that

$$\max_{n_0 \leq k \leq n_2} \mathbf{E}^* \left\| \frac{1}{\sqrt{k}} \sum_{i=1}^k \epsilon_i \cdot (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n} \leq \max_{n_0 \leq k \leq n_2} 2 \cdot \mathbf{E}_0^* \left\| \frac{1}{\sqrt{k}} \sum_{i=1}^k (\delta_{(\tilde{Y}_i, \mathbf{X}_i)} - \mathbf{P}_0) \right\|_{\mathcal{G}_n}.$$

Also note that, for any real valued sequence $\{b_n : n \geq 1\}$ satisfying $|b_n| \rightarrow 0$, it follows that $\max\{|b_j| : m_1 \leq j \leq m_2\} \rightarrow 0$, provided $m_1 \rightarrow \infty$. So, it will be enough to show that $\mathbf{E}_0^* \|\mathbb{G}_{n_2}\|_{\mathcal{G}_n} \rightarrow 0$, which itself follows from the arguments used in handling $P_{2,n}$. Since $n_0 \rightarrow \infty$, the first term in the right side of (8.25) converges to zero. Thus the right side of (8.24) converges to zero. Hence, $P_{1,n} \rightarrow 0$.

Combining all these steps and going back to (8.23), we conclude that $\mathbf{E}_0^* \mathbf{E}_M^* \left(\|\widehat{\mathbb{G}}_{n_2}\|_{\mathcal{G}_n} \right) \rightarrow 0$, which implies for any $\eta > 0$, $\mathbf{P}_M^* (\|\widehat{\mathbb{G}}_{n_2}\|_{\mathcal{G}_n} > \eta) = o_{\mathbf{P}_0^*}(1)$. Thus,

$$\mathbf{P}_M^* \left(\sup_{\beta} \max_{1 \leq i \leq p} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\theta_0+n_2^{-1/2}\mathbf{u}_i} - h_{2,\beta,\theta_0} \right) \right| > \eta/3 \right) = o_{\mathbf{P}_0^*}(1).$$

By Lemma 8.18, there will exist some sequence $\delta_n \downarrow 0$, such that $L_{1,n}$ and $L_{3,n}$ (cf. (8.22)) are $o_{\mathbf{P}_0^*}(1)$, with the fixed δ replaced by δ_n . Choose this sequence δ_n in $L_{2,n}$. Hence the left side of (8.22) satisfies

$$\mathbf{P}_M^* \left(\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| > \eta, \sqrt{n_2}(\hat{\theta}_n - \theta_0) \in K_{\epsilon_1, \epsilon_2}^{\delta_n/2} \right) = o_{\mathbf{P}_0^*}(1).$$

Combining this with (8.20), (8.21) and since $\epsilon_1, \epsilon_2 > 0$, are arbitrary, we can conclude that

$$\sup_{\beta} \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| = o_{\mathbf{P}_M^*}(1), \quad \text{in } \mathbf{P}_0^*\text{-probability.}$$

(iii) In this case, from Theorem 2.1 $\sqrt{n_2}(\hat{\theta}_n - \theta_0) = O_{\mathbf{P}_0^*}(1)$ and hence $O_{\mathbf{P}_M^*}(1)$. The remaining argument can be carried out along the same lines as done in part (ii) above. We skip the details.

□

Lemma 8.15. *Suppose, assumptions (A1)-(A4) hold. Then,*

$$\sqrt{n} \left(Z(\hat{\beta}_n) - Z(\beta_n) \right) + \sqrt{n} \cdot \left(\hat{Z}_n(\beta_0) - Z_n(\beta_0) \right) = o_{\mathbf{P}_M^*}(1) \quad \text{in } \mathbf{P}_0^* \text{-probability.} \quad (8.26)$$

Proof of Lemma 8.15. Recall the definitions of $\hat{Z}_n(\beta)$, $Z_n(\beta)$ and $Z(\beta)$ from (3.2), (2.6) and (2.8) respectively. Then, we can write

$$\begin{aligned} & \sqrt{n} \left(Z(\hat{\beta}_n) - Z(\beta_n) \right) + \sqrt{n} \left(\hat{Z}_n(\beta_0) - Z_n(\beta_0) \right) \\ &= \left[\sqrt{n} \hat{Z}_n(\hat{\beta}_n) - \sqrt{n} Z_n(\hat{\beta}_n) \right] + \sqrt{n} (Z_n - Z) (\hat{\beta}_n - \beta_0) + \sqrt{n} (Z_n - Z) (\hat{\beta}_n - \beta_0) \\ & \quad - \sqrt{n} (\hat{Z}_n - Z_n) (\hat{\beta}_n - \beta_0) \\ &\equiv F_{1,n} + F_{2,n} + F_{3,n} - F_{4,n}. \end{aligned} \quad (8.27)$$

We will work with the first components of $h_{1,\beta}$ and $h_{2,\beta,\theta}$ and denote them by the same symbols. Consider the first components of each term of (8.27) separately. Since $\hat{\beta}_n$ is defined to be an exact zero of $\hat{Z}_n(\beta)$, hence $\sqrt{n} \cdot \hat{Z}_n(\hat{\beta}_n) = o_{\mathbf{P}_M^*}(1)$, trivially. Similarly, $\sqrt{n} \cdot Z_n(\hat{\beta}_n) = o_{\mathbf{P}_M^*}(1)$. From Lemma 8.11 it now follows that,

$$F_{1,n} = o_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_M^*}(1) = o_{\mathbf{P}_M^*}(1) = o_{\mathbf{P}_M^*}(1) \quad \text{in } \mathbf{P}_0^* \text{-probability.}$$

Also, from Lemma 8.10 and Lemma 8.11 we obtain,

$$F_{2,n} = \sqrt{n} \cdot (Z_n - Z) (\hat{\beta}_n - \beta_0) = o_{\mathbf{P}_0^*}(1) = o_{\mathbf{P}_M^*}(1) = o_{\mathbf{P}_M^*}(1) \quad \text{in } \mathbf{P}_0^* \text{-probability.}$$

In order to study $F_{3,n}$, we will re-trace the arguments of Lemma 8.10, with some important modifications. Following (8.13), we can write,

$$\begin{aligned} F_{3,n} &= \sqrt{n} \cdot (Z_n - Z) (\hat{\beta}_n - \beta_0) \\ &= \sqrt{f_n} \cdot \mathbb{G}_{n_1} (h_{1,\hat{\beta}_n} - h_{1,\beta_0}) + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot [Z_{n,2}(\hat{\beta}_n) - Z_2(\hat{\beta}_n) - Z_{n,2}(\beta_0) + Z_2(\beta_0)] \\ & \quad + \sqrt{n} \cdot (f_n - f) \cdot \mathbf{P}_0 (h_{1,\hat{\beta}_n} - h_{1,\beta_0}) - \sqrt{n} \cdot (f_n - f) \cdot [Z_2(\hat{\beta}_n) - Z_2(\beta_0)] \\ &\equiv H_{1,n} + H_{2,n} + H_{3,n} + H_{4,n}. \end{aligned} \quad (8.28)$$

We consider each term in the right side of (8.28) separately. From Lemma 8.13(i) it follows that $H_{1,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Next, following (8.14) and (8.15) write,

$$\begin{aligned} H_{2,n} &= \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left[\mathbb{P}_{n_2} h_{2,\hat{\beta}_n,\hat{\theta}_n} - \mathbf{P}_0 h_{2,\hat{\beta}_n,\theta_0} - \mathbb{P}_{n_2} h_{2,\beta_0,\hat{\theta}_n} + \mathbf{P}_0 h_{2,\beta_0,\theta_0} \right], \\ \tilde{H}_{2,n} &= \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left[\mathbb{P}_{n_2} h_{2,\hat{\beta}_n,\theta_0} - \mathbf{P}_0 h_{2,\hat{\beta}_n,\theta_0} - \mathbb{P}_{n_2} h_{2,\beta_0,\theta_0} + \mathbf{P}_0 h_{2,\beta_0,\theta_0} \right] \\ &= \sqrt{1-f_n} \cdot \mathbb{G}_{n_2} \left(h_{2,\hat{\beta}_n,\theta_0} - h_{2,\beta_0,\theta_0} \right), \end{aligned}$$

From Lemma 8.13(ii) it follows that $\tilde{H}_{2,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Also,

$$\begin{aligned} |H_{2,n} - \tilde{H}_{2,n}| &\leq 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta \in \mathbb{R}^p} \left| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| \\ &\quad + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left| \mathbf{P}_0 \left(h_{2,\hat{\beta}_n,\hat{\theta}_n} - h_{2,\hat{\beta}_n,\theta_0} - h_{2,\beta_0,\hat{\theta}_n} + h_{2,\beta_0,\theta_0} \right) \right|. \end{aligned}$$

From Lemma 8.6 and Lemma 8.11,

$$2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta \in \mathbb{R}^p} \left| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| = o_{\mathbf{P}_M^*}(1) \quad \text{in } \mathbf{P}_0^*\text{-probability.}$$

Also, from Lemma 2.1 and Lemma 8.11 it follows that, $\|\hat{\theta}_n - \theta_0\| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Combining this with Theorem 3.2(i) we obtain, $\|\hat{\beta}_n - \beta_0\| + \|\hat{\theta}_n - \theta_0\| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Following the approach used in handling the term $C_{2,n}$ in the proof of Lemma 8.10 and using the continuous mapping theorem repeatedly, we obtain

$$\sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left| \mathbf{P}_0 h_{2,\hat{\beta}_n,\hat{\theta}_n} - \mathbf{P}_0 h_{2,\hat{\beta}_n,\theta_0} - \left(\mathbf{P}_0 h_{2,\beta_0,\hat{\theta}_n} - \mathbf{P}_0 h_{2,\beta_0,\theta_0} \right) \right| = o_{\mathbf{P}_M^*}(1),$$

in \mathbf{P}_0^* -probability. Hence, $|H_{2,n} - \tilde{H}_{2,n}| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, which implies $H_{2,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Finally, mimicking the arguments used for the third and fourth terms in the proof of Lemma 8.10 and applying the continuous mapping theorem for bootstrapped random variables, we obtain $H_{j,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, for $j = 3, 4$. Applying the same arguments to all p components of $h_{1,\beta}$ and $h_{2,\beta,\theta}$, we can conclude $F_{3,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability (cf. (8.28)). Now, consider the term $F_{4,n}$ in (8.27). Then,

$$\begin{aligned} F_{4,n} &= \sqrt{n} \cdot \left[\hat{Z}_n(\hat{\beta}_n) - Z_n(\hat{\beta}_n) - \hat{Z}_n(\beta_0) + Z_n(\beta_0) \right] \\ &= \sqrt{f_n} \cdot \hat{\mathbb{G}}_{n_1} \left(h_{1,\hat{\beta}_n} - h_{1,\beta_0} \right) + \sqrt{n} \cdot (1-f_n) \cdot \left[\hat{\mathbb{P}}_{n_2} h_{2,\hat{\beta}_n,\hat{\theta}_n} - \mathbb{P}_{n_2} h_{2,\hat{\beta}_n,\hat{\theta}_n} - \hat{\mathbb{P}}_{n_2} h_{2,\beta_0,\hat{\theta}_n} + \mathbb{P}_{n_2} h_{2,\beta_0,\hat{\theta}_n} \right] \\ &\equiv I_{1,n} + I_{2,n}. \end{aligned} \tag{8.29}$$

From Lemma 8.14(i) it follows that $I_{1,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Define the new term,

$$\begin{aligned} \tilde{I}_{2,n} &= \sqrt{n} \cdot (1-f_n) \cdot \left[\hat{\mathbb{P}}_{n_2} h_{2,\hat{\beta}_n,\hat{\theta}_n} - \mathbb{P}_{n_2} h_{2,\hat{\beta}_n,\hat{\theta}_n} - \hat{\mathbb{P}}_{n_2} h_{2,\beta_0,\hat{\theta}_n} + \mathbb{P}_{n_2} h_{2,\beta_0,\hat{\theta}_n} \right] \\ &= \sqrt{1-f_n} \cdot \hat{\mathbb{G}}_{n_2} \left(h_{2,\hat{\beta}_n,\hat{\theta}_n} - h_{2,\beta_0,\hat{\theta}_n} \right). \end{aligned} \tag{8.30}$$

Then, after some algebra it can be shown that

$$\begin{aligned} |I_{2,n} - \tilde{I}_{2,n}| &\leq 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta} \left| \hat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| + 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta} \left| \hat{\mathbb{G}}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| \\ &\quad + 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta} \left| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| + 2 \cdot \sqrt{1-f_n} \cdot \sup_{\beta} \left| \mathbb{G}_{n_2} \left(h_{2,\beta,\hat{\theta}_n} - h_{2,\beta,\theta_0} \right) \right| \\ &\quad + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \left| \mathbf{P}_0 \left(h_{2,\hat{\beta}_n,\hat{\theta}_n} - h_{2,\hat{\beta}_n,\theta_0} - h_{2,\beta_0,\hat{\theta}_n} + h_{2,\beta_0,\theta_0} \right) \right| \\ &\equiv J_{1,n} + J_{2,n} + J_{3,n} + J_{4,n} + J_{5,n}. \end{aligned} \tag{8.31}$$

Consider each term of (8.31) separately. Using Lemma 8.14 it follows that, $J_{1,n} = o_{\mathbf{P}_M^*}(1)$ and $J_{2,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Next, note that $\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \boldsymbol{\theta}_0)$ converges to a tight limiting distribution under the product probability measure \mathbf{Pr} . Hence, the proof of Lemma 8.6 holds in this case, without any changes except that all probability statements are understood in terms of \mathbf{Pr} instead of \mathbf{P}_0 . We skip the details about verification of the conditions. Hence, $J_{3,n} = o_{\mathbf{P}_M^*}(1)$ which implies, $J_{3,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. From Lemma 8.6 and Lemma 8.11, it follows that $J_{4,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. We can use the arguments used to handle the term $C_{2,n}$ (cf. (8.16)) in Lemma 8.10 along with the fact that $\sqrt{n_2}(\widehat{\boldsymbol{\theta}}_n - \widehat{\boldsymbol{\theta}}_n) = O_{\mathbf{P}_M^*}(1)$ and $\|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. We skip the details. Hence, $J_{5,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Combining all parts in the right side of (8.31), we obtain $|I_{2,n} - \tilde{I}_{2,n}| = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Also from (8.30),

$$|\tilde{I}_{2,n}| = \sqrt{1-f_n} \cdot \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\widehat{\boldsymbol{\beta}}_n,\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) \right| \leq J_{2,n} + \sqrt{1-f_n} \cdot \left| \widehat{\mathbb{G}}_{n_2} \left(h_{2,\widehat{\boldsymbol{\beta}}_n,\boldsymbol{\theta}_0} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) \right|, \quad \text{cf. (8.31)}.$$

We already know that $J_{2,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. The second term can be shown to be $o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, by exactly following the argument given in the proof of Lemma 8.13(i) and by noting that $\{h_{2,\boldsymbol{\beta},\boldsymbol{\theta}_0} : \boldsymbol{\beta} \in \mathbb{R}^p\}$ is a Donsker class. Hence, $\tilde{I}_{2,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Combining all steps and from (8.29) we obtain, $F_{4,n} = o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Following the argument separately for all p components in the right side of (8.27) completes the proof. \square

Lemma 8.16. *Suppose, assumptions (A1)-(A5) hold. Then,*

$$\sqrt{n} \cdot \|\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0\| = O_{\mathbf{P}_M^*}(1), \quad \text{in } \mathbf{P}_0^*\text{-probability.}$$

Proof of Lemma 8.16. Note that $Z(\boldsymbol{\beta}_0) = \mathbf{0}$ and $\widehat{Z}_n(\widehat{\boldsymbol{\beta}}_n) = \mathbf{0}$. Recall the decomposition obtained in (8.27). Then, we can write

$$\begin{aligned} \sqrt{n}(Z(\widehat{\boldsymbol{\beta}}_n) - Z(\boldsymbol{\beta}_0)) &= \sqrt{n}Z(\widehat{\boldsymbol{\beta}}_n) = \sqrt{n}(Z - Z_n)(\widehat{\boldsymbol{\beta}}_n) + \sqrt{n}(Z_n - \widehat{Z}_n)(\widehat{\boldsymbol{\beta}}_n) + \sqrt{n}\widehat{Z}_n(\widehat{\boldsymbol{\beta}}_n) \\ &= \sqrt{n}(Z - Z_n)(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) + \sqrt{n}(Z_n - \widehat{Z}_n)(\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0) + \sqrt{n}(Z - Z_n)(\boldsymbol{\beta}_0) + \sqrt{n}(Z_n - \widehat{Z}_n)(\boldsymbol{\beta}_0) \\ &\quad + \sqrt{n}\widehat{Z}_n(\widehat{\boldsymbol{\beta}}_n) \\ &\equiv -F_{3,n} - F_{4,n} - Q_{1,n} - Q_{2,n} + Q_{3,n}, \end{aligned} \tag{8.32}$$

where, $F_{3,n}$ and $F_{4,n}$ have been defined earlier in (8.27). Consider each term in the right side of (8.32) separately. Following the proof of Lemma 8.15, it follows that $F_{3,n}$ and $F_{4,n}$ are $o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Also, $Q_{3,n} = \mathbf{0}$, by definition of $\widehat{\boldsymbol{\beta}}_n$ (cf. (3.3)). Finally, after some manipulations it can be shown that

$$\begin{aligned} &Q_{1,n} + Q_{2,n} \\ &= \sqrt{f_n} \cdot \widehat{\mathbb{G}}_{n_1} h_{1,\boldsymbol{\beta}_0} + \sqrt{1-f_n} \cdot \widehat{\mathbb{G}}_{n_2} h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} + \sqrt{1-f_n} \cdot \widehat{\mathbb{G}}_{n_2} \left(h_{2,\boldsymbol{\beta}_0,\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) \\ &\quad - \sqrt{1-f_n} \cdot \widehat{\mathbb{G}}_{n_2} \left(h_{2,\boldsymbol{\beta}_0,\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) + \sqrt{1-f_n} \cdot \widehat{\mathbb{G}}_{n_2} \left(h_{2,\boldsymbol{\beta}_0,\widehat{\boldsymbol{\theta}}_n} - h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0} \right) + \sqrt{f_n} \cdot \widehat{\mathbb{G}}_{n_1} h_{1,\boldsymbol{\beta}_0} \end{aligned}$$

$$\begin{aligned}
& + \sqrt{1-f_n} \cdot \mathbb{G}_{n_2} \left(h_{2,\beta_0,\hat{\theta}_n} - h_{2,\beta_0,\theta_0} \right) + \sqrt{1-f_n} \cdot \mathbb{G}_{n_2} h_{2,\beta_0,\theta_0} \\
& + \sqrt{1-f_n} \cdot \sqrt{n_2} \cdot \mathbf{P}_0 \left(h_{2,\beta_0,\hat{\theta}_n} - h_{2,\beta_0,\theta_0} \right) + \sqrt{n}(f_n - f) \cdot \mathbf{P}_0 \left(h_{1,\beta_0} - h_{2,\beta_0,\theta_0} \right). \tag{8.33}
\end{aligned}$$

Since $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ is a Donsker class (cf. Lemma 8.4(i)), hence $\widehat{\mathbb{G}}_{n_1} h_{1,\beta_0} = O_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, by Theorem 3.6.1 of [van der Vaart and Wellner \(1996\)](#). A similar argument shows that, $\widehat{\mathbb{G}}_{n_2} h_{2,\beta_0,\theta_0} = O_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, since $\{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$ is Donsker. Also, using the Delta method, Lemma 3.1 and the methods used in the proof of Theorem 2.2 we can claim that

$$\sqrt{n_2} \cdot \mathbf{P}_0 \left(h_{2,\beta_0,\hat{\theta}_n} - h_{2,\beta_0,\theta_0} \right) = O_{\mathbf{P}_M^*}(1), \quad \text{in } \mathbf{P}_0^*\text{-probability.}$$

Following the proof of Theorem 2.2, $\sqrt{n}(f_n - f) \cdot \mathbf{P}_0 \left(h_{1,\beta_0} - h_{2,\beta_0,\theta_0} \right) = o(1)$. Since $\{h_{1,\beta} : \beta \in \mathbb{R}^p\}$ and $\{h_{2,\beta,\theta_0} : \beta \in \mathbb{R}^p\}$ are Donsker, hence $\mathbb{G}_{n_1} h_{1,\beta_0}$ and $\mathbb{G}_{n_2} h_{2,\beta_0,\theta_0}$ are $O_{\mathbf{P}_R^*}(1)$. The remaining terms in the right side of (8.33) are $o_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability, using Lemma's 8.13 and 8.14. So, $Q_{1,n} + Q_{2,n} = O_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. Therefore, continuing from (8.32) and applying triangle inequality we obtain,

$$\begin{aligned}
& \left\| \sqrt{n} \left(Z(\widehat{\beta}_n) - Z(\beta_0) \right) \right\| - \|Q_{1,n} + Q_{2,n}\| \leq \|F_{3,n}\| + \|F_{4,n}\| + \|Q_{3,n}\| \\
& = o_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_R^*}(1) = o_{\mathbf{P}_M^*} \left(1 + \sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\| \right).
\end{aligned}$$

because a $o_{\mathbf{P}_M^*}(1)$ term is obviously a $o_{\mathbf{P}_M^*}(1 + \sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\|)$ term in \mathbf{P}_0^* -probability. Since $Z(\beta)$ is differentiable at β_0 and the derivative is nonsingular (cf. assumption (A5)), we can write for some $c > 0$,

$$\begin{aligned}
\sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\| \cdot \left(c + o_{\mathbf{P}_M^*}(1) \right) & \leq \sqrt{n} \left\| \left(Z(\widehat{\beta}_n) - Z(\beta_0) \right) \right\| \leq \|Q_{1,n} + Q_{2,n}\| + o_{\mathbf{P}_M^*} \left(1 + \sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\| \right) \\
& = O_{\mathbf{P}_M^*}(1) + o_{\mathbf{P}_M^*} \left(1 + \sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\| \right).
\end{aligned}$$

This implies $\sqrt{n} \cdot \|\widehat{\beta}_n - \beta_0\| = O_{\mathbf{P}_M^*}(1)$ in \mathbf{P}_0^* -probability. \square

8.3 Two useful mathematical results

The result in Lemma 8.17 describes an useful property of the gradient vector of a strictly concave function with domain \mathbb{R}^p . This result has been used in the proof of Lemma 8.2 to show the well-separated condition for the unique root of $Z_1(\beta)$. Lemma 8.18 is a simple result about sequences indexed by two indices and has been used in the proofs of Lemma 8.13, Lemma 8.14 and Theorem 3.2(i). The symbol $\|\cdot\|$ is used to denote the Euclidean norm in \mathbb{R}^p .

Lemma 8.17. *Suppose, $f : \mathbb{R}^p \mapsto \mathbb{R}$ is a strictly concave function which is differentiable at all points and has an unique maxima at some $\mathbf{x}_0 \in \mathbb{R}^p$. Denote the gradient of f at a point \mathbf{x} by $\nabla f(\mathbf{x})$. Then, the following*

strict inequalities hold:

$$\left. \begin{aligned} f(\mathbf{x}_0) &> \sup_{\mathbf{x}: \|\mathbf{x}-\mathbf{x}_0\| \geq \delta} f(\mathbf{x}), \\ 0 = \|\nabla f(\mathbf{x}_0)\| &< \inf_{\mathbf{x}: \|\mathbf{x}-\mathbf{x}_0\| \geq \delta} \|\nabla f(\mathbf{x})\|, \end{aligned} \right\}, \quad \text{for all } \delta > 0. \quad (8.34)$$

Proof of Lemma 8.17. We complete the proof in separate parts.

- (i) Consider any $\mathbf{x} \in \mathbb{R}^p$ (with $\mathbf{x} \neq \mathbf{x}_0$) and the line passing segment joining \mathbf{x}_0 and \mathbf{x} . Any point on the interior of this line segment can be represented as, $\mathbf{y}(\lambda) = \lambda\mathbf{x} + (1-\lambda)\mathbf{x}_0$, for some $\lambda \in (0, 1)$. Since \mathbf{x}_0 is the unique maxima and f is strictly concave on this line segment, we must have

$$f(\mathbf{x}_0) > f(\mathbf{y}(\lambda)) > f(\mathbf{x}), \quad \text{for all } \lambda \in (0, 1).$$

Now, for any $\delta > 0$, define the sets $G_\delta = \{\mathbf{x} : \|\mathbf{x} - \mathbf{x}_0\| = \delta\}$. We claim that,

$$\sup_{\mathbf{x} \in G_{\delta_1}} f(\mathbf{x}) > \sup_{\mathbf{x} \in G_{\delta_2}} f(\mathbf{x}), \quad \text{for any } 0 < \delta_1 < \delta_2. \quad (8.35)$$

If possible, assume that the claim is false and the strict inequality in (8.35) fails for some pair of values, (δ_1, δ_2) with $\delta_1 < \delta_2$. Note that G_{δ_i} are compact. Then, as f is continuous, there exists a point $\mathbf{x}_{\delta_i} \in G_{\delta_i}$, such that $f(\mathbf{x}_{\delta_i}) = \sup_{\mathbf{x} \in G_{\delta_i}} f(\mathbf{x})$, $i = 1, 2$. As per our assumption we have, $f(\mathbf{x}_{\delta_1}) \leq f(\mathbf{x}_{\delta_2})$. Join the points \mathbf{x}_0 and \mathbf{x}_{δ_2} by a line segment which intersects G_{δ_1} (which is the inner circle) at some interior point (on the line) $\mathbf{x}_1 \in G_{\delta_1}$. By strict concavity of f over this line segment, we have

$$f(\mathbf{x}_{\delta_2}) < f(\mathbf{x}_1) \leq f(\mathbf{x}_{\delta_1}),$$

which contradicts our assumption and as a result (8.35) holds. For any $\delta > 0$,

$$\sup_{\mathbf{x}: \|\mathbf{x}-\mathbf{x}_0\| \geq \delta} f(\mathbf{x}) = \sup \left\{ f(\mathbf{x}) : \mathbf{x} \in \bigcup_{\eta \geq \delta} G_\eta \right\} = \sup_{\mathbf{x} \in G_\delta} f(\mathbf{x}) = f(\mathbf{x}_\delta),$$

for some $\mathbf{x}_\delta \in G_\delta$. Now, draw a line segment through the points \mathbf{x}_0 and \mathbf{x}_δ . Since $\|\mathbf{x}_\delta - \mathbf{x}_0\| = \delta > 0$, we can use strict concavity to claim that, $f(\mathbf{x}_\delta) < f(\mathbf{x}_0)$. This completes the proof of the first statement.

- (ii) Note that $\nabla f(\mathbf{x}) = (D_1 f(\mathbf{x}), \dots, D_p f(\mathbf{x}))'$, where $D_j f(\mathbf{x})$ denotes the j -th partial derivative of f at the point \mathbf{x} . Also, $D_j f(\mathbf{x})$ is the directional derivative of f at the point \mathbf{x} in the direction of \mathbf{e}_j , where \mathbf{e}_j denotes the j -th unit vector in \mathbb{R}^p . Finally note that, for a strictly concave and differentiable function f , $\|\nabla f(\mathbf{x})\| = 0$, if and only if \mathbf{x} is a global maxima of f .

Consider the line segment through \mathbf{x}_0 in the direction \mathbf{e}_1 . Define the function, $g_{\mathbf{x}_0}(t) = f(\mathbf{x}_0 + t\mathbf{e}_1)$, $t \in \mathbb{R}$. It is easily seen that $t \mapsto g_{\mathbf{x}_0}(t)$ is strictly concave and also differentiable (since f is differentiable

everywhere). Fix any $\delta > 0$. Using the proof of part (i) above, we know there exists a point \mathbf{x}_δ on G_δ such that,

$$\begin{aligned}
f(\mathbf{x}_0) &> f(\mathbf{x}_\delta) = \sup_{\mathbf{x} \in G_\delta} f(\mathbf{x}) \geq f(\mathbf{x}_0 + \delta \mathbf{e}_1), \quad \text{since, } \mathbf{x}_0 + \delta \mathbf{e}_1 \in G_\delta, \\
\Rightarrow 0 &< c_\delta \equiv f(\mathbf{x}_0) - f(\mathbf{x}_\delta) \leq f(\mathbf{x}_0) - f(\mathbf{x}_0 + \delta \mathbf{e}_1) = g_{\mathbf{x}_0}(0) - g_{\mathbf{x}}(\delta), \\
\Rightarrow 0 &< c_\delta \leq (0 - \delta) \cdot g'_{\mathbf{x}_0}(\theta\delta) = D_1 f(\mathbf{x}_0 + \theta\delta \mathbf{e}_1), \quad \text{using the Mean Value Theorem, for some } \theta \in (0, 1), \\
\Rightarrow 0 &< \frac{|c_\delta|}{|\delta|} \leq |D_1 f(\mathbf{x}_0 + \theta\delta \mathbf{e}_1)| \leq \|\nabla f(\mathbf{x}_0 + \theta\delta \mathbf{e}_1)\|. \tag{8.36}
\end{aligned}$$

Note that, $\|\mathbf{x}_0 - (\mathbf{x}_0 + \theta\delta \mathbf{e}_1)\| = \theta\delta$. Hence,

$$0 < \frac{|c_\delta|}{|\delta|} \leq \inf_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}_0\| \geq \theta\delta} \|\nabla f(\mathbf{x})\| \leq \inf_{\mathbf{x}: \|\mathbf{x} - \mathbf{x}_0\| \geq \delta} \|\nabla f(\mathbf{x})\|.$$

This completes the proof. □

Lemma 8.18 (A result on sequences). *Suppose $\{a_{n,k} : n \geq 1, k \geq 1\}$ is a sequence of non-negative real numbers such that, for each fixed $k \geq 1$, $\lim_{n \rightarrow \infty} a_{n,k} = 0$. Then, there exists a sequence $\{k_n : n \geq 1\}$ such that $\lim_{n \rightarrow \infty} a_{n,k_n} = 0$.*

Proof of Lemma 8.18. The proof can be found in [Lahiri \(2003\)](#) (pg. 78). □

8.4 Detailed expressions for matrices used in Theorem 2.2 and Lemma 8.8

We provide detailed expressions of the matrices used in the statement of Theorem 2.2 and Lemma 8.8. Note that \mathbf{E}_0 , \mathbf{Var}_0 and \mathbf{cov}_0 denote expectation, variance and covariance under \mathbf{P}_0 . In case of the with intercept model (2.9), computation of the matrices shown below can be carried out by replacing x_1 with 1.

Define, $\mathbf{\Sigma} = \mathbf{Var}_0(\mathbf{T}_1) = (\sigma_{i,j} : 1 \leq i, j \leq p+3)$, where \mathbf{T}_1 is defined in (7.1). Note that, $\mathbf{\Sigma}_{2,2} = \mathbf{Var}_0(\mathbf{T}^{(2)})$ has been found in (2.11). The remaining components of $\mathbf{\Sigma}$ are, $\mathbf{\Sigma}_{1,1} = \mathbf{Var}_0(\mathbf{T}^{(1)})$ and $\mathbf{\Sigma}_{1,2} = \mathbf{\Sigma}'_{2,1} = \mathbf{cov}_0(\mathbf{T}^{(1)}, \mathbf{T}^{(2)})$. It is easy to obtain the following expressions,

$$\left. \begin{aligned}
\sigma_{i,j} &= \int x_i x_j \cdot \psi(\mathbf{x}'\boldsymbol{\beta}_0) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}_0)\} dQ(\mathbf{x}), \quad \text{for all } j = 1, \dots, p, \\
\sigma_{i,p+1} &= -(1 - \theta_{1,0}) \cdot \int x_i \psi(\mathbf{x}'\boldsymbol{\beta}_0) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}_0)\} dQ(\mathbf{x}), \\
\sigma_{i,p+2} &= \theta_{2,0} \cdot \int x_i \psi(\mathbf{x}'\boldsymbol{\beta}_0) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}_0)\} dQ(\mathbf{x}), \quad \text{and} \\
\sigma_{i,p+3} &= -\theta_{1,0} \int x_i \psi(\mathbf{x}'\boldsymbol{\beta}_0) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}_0)\} dQ(\mathbf{x}).
\end{aligned} \right\} \forall i = 1, \dots, p. \tag{8.37}$$

Next, write $\mathbf{\Gamma} = \mathbf{Var}_0(h_{2,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}(\tilde{Y}, \mathbf{X})) = (\Gamma_{i,j} : 1 \leq i, j \leq p)$ (cf. (2.3)). The expression for $\Gamma_{i,j}$ is,

$$\Gamma_{i,j} = (1 - \theta_{1,0} - \theta_{2,0})^2 \int x_i x_j \cdot \frac{[\psi(\mathbf{x}'\boldsymbol{\beta}_0) \{1 - \psi(\mathbf{x}'\boldsymbol{\beta}_0)\}]^2}{h_{3,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}(\mathbf{x}) \cdot (1 - h_{3,\boldsymbol{\beta}_0,\boldsymbol{\theta}_0}(\mathbf{x}))} dQ(\mathbf{x}), \quad \text{for all } 1 \leq i, j \leq p. \tag{8.38}$$

The $p \times 2$ matrix $\mathbf{A}_\theta = (\mathbf{A}_\theta(i, j) : 1 \leq i \leq p, 1 \leq j \leq 2)$ is the total derivative map of $\bar{f}(\theta) = \mathbf{P}_0 h_{2, \beta_0, \theta}$ (cf. Lemma 8.8(ii)) with respect to θ . Also, \mathbf{A}_0 is the value of \mathbf{A}_θ at $\theta = \theta_0$, and is used in Theorem 2.2(ii). To simplify the expression, define

$$u_\theta(\tilde{y}, \mathbf{x}) \equiv \frac{\{1 - \psi(\mathbf{x}'\beta_0)\} \cdot \{\tilde{y} - h_{3, \beta_0, \theta}(\mathbf{x})\}}{h_{3, \beta_0, \theta}(\mathbf{x}) \cdot \{1 - h_{3, \beta_0, \theta}(\mathbf{x})\}}, \quad \text{for all } (\tilde{y}, \mathbf{x}) \text{ and all } \theta \in \Theta.$$

Then, for all $i = 1, \dots, p$ and $j = 1, 2$, elements of \mathbf{A}_θ can be expressed as,

$$\mathbf{A}_\theta(i, j) = - \int x_i \cdot \psi(\mathbf{x}'\beta_0) \cdot u_\theta(\tilde{y}, \mathbf{x}) \cdot \left[1 + (-1)^{j+1} \cdot (1 - \theta_1 - \theta_2) \cdot u_\theta(\tilde{y}, \mathbf{x}) \right] d\mathbf{P}_0(\tilde{y}, \mathbf{x}). \quad (8.39)$$

In case $\theta = \theta_0$, $\mathbf{E}_0(\tilde{Y}|\mathbf{X} = \mathbf{x}) = h_{3, \beta_0, \theta_0}(\mathbf{x})$, which implies $\mathbf{E}_0(u_{\beta_0, \theta_0}(\tilde{Y}, \mathbf{X})|\mathbf{X} = \mathbf{x}) = 0$. Hence, $\mathbf{Var}_0(\tilde{Y}|\mathbf{X} = \mathbf{x}) = h_{3, \beta_0, \theta_0}(\mathbf{x}) \cdot \{1 - h_{3, \beta_0, \theta_0}(\mathbf{x})\}$. After simplification, the expression for elements of \mathbf{A}_0 will be,

$$\mathbf{A}_0(i, j) = (-1)^j \cdot (1 - \theta_{1,0} - \theta_{2,0}) \cdot \int x_i \cdot \frac{\psi(\mathbf{x}'\beta_0) \cdot \{1 - \psi(\mathbf{x}'\beta_0)\}^2}{h_{3, \beta_0, \theta_0}(\mathbf{x}) \{1 - h_{3, \beta_0, \theta_0}(\mathbf{x})\}} dQ(\mathbf{x}), \quad \text{for all } (i, j). \quad (8.40)$$

Write, $\dot{Z}(\beta_0) = (\dot{Z}(\beta_0)_{(i,j)} : 1 \leq i, j \leq p)$. It can be checked that,

$$\begin{aligned} \dot{Z}(\beta_0)_{(i,j)} &= -f \cdot \int x_i x_j \cdot \psi(\mathbf{x}'\beta_0) \{1 - \psi(\mathbf{x}'\beta_0)\} dQ(\mathbf{x}) \\ &\quad - (1 - f) \cdot (1 - \theta_{1,0} - \theta_{2,0})^2 \cdot \int x_i x_j \cdot \frac{(\psi(\mathbf{x}'\beta_0) \{1 - \psi(\mathbf{x}'\beta_0)\})^2}{h_{3, \beta_0, \theta_0}(\mathbf{x}) \{1 - h_{3, \beta_0, \theta_0}(\mathbf{x})\}} dQ(\mathbf{x}). \end{aligned} \quad (8.41)$$