

# Informational Content of Special Regressors in Heteroskedastic Binary Response Models<sup>1</sup>

Songnian Chen<sup>2</sup>, Shakeeb Khan<sup>3</sup> and Xun Tang<sup>4</sup>

This Version: April 9, 2013

We quantify the identifying power of special regressors in heteroskedastic binary regressions with median-independent or conditionally symmetric errors. We measure the identifying power using two criteria: the set of regressor values that help point identify coefficients in latent payoffs as in (Manski 1988); and the Fisher information of coefficients as in (Chamberlain 1986). We find for median-independent errors, requiring one of the regressors to be special does not add to the identifying power or the information for coefficients. Nonetheless it does help identify the error distribution and the average structural function. For conditionally symmetric errors, the presence of a special regressor improves the identifying power by the criterion in (Manski 1988), and the Fisher information for coefficients is strictly positive under mild conditions. We propose a new estimator for coefficients that converges at the parametric rate under symmetric errors and a special regressor, and report its decent performance in small samples through simulations.

*Key words: Binary regression, heteroskedasticity, identification, information, median independence, conditional symmetry*

---

<sup>1</sup>We are grateful to Arthur Lewbel, Jim Powell, Frank Schorfheide, Brendan Kline, Haiqing Xu and seminar participants at UT Austin and Yale for comments. We thank Bingzhi (Ben) Zhao for capable research assistance. The usual disclaimer applies.

<sup>2</sup>Economics Department, Hong Kong University of Science and Technology. Email: snchen@ust.hk

<sup>3</sup>Economics Department, Duke University. Email: shakeebk@duke.edu.

<sup>4</sup>Economics Department, University of Pennsylvania. Email: xuntang@sas.upenn.edu.

# 1 Introduction

In a binary response model, a special regressor is one that is additively separable from all other components in the latent payoffs and that satisfies an exclusion restriction (i.e. being independent from the error conditional on all other regressors). It is worth noting that we do not include any “large support” requirement in the definition of special regressors. In this paper, we examine how a special regressor contributes to the identification and the Fisher information of coefficients in semiparametric binary regressions with heteroskedastic errors. We focus on the role of special regressors in two models where errors are median independence or conditional symmetry respectively. These models are of particular interest, because identification of coefficients in them does not require the “large support” condition (i.e. the support of special regressors includes that of the error), a condition typically used in identification-at-infinity arguments. To our knowledge, this paper marks the first effort to address this question in the literature.

Special regressor arise in various social-economic contexts. (Lewbel 2000) used a special regressor to recover coefficients in semiparametric binary regressions where heteroskedastic errors are mean-independent from regressors. He showed coefficients for all regressors along with the error distribution are identified up to scale, provided the support of special regressor is large enough. (Lewbel 2000) then proposed a two-step inverse-density-weighted estimator. Since then, arguments based on special regressors have been used to identify structural micro-econometric models in a variety of contexts. These include multinomial-choice demand models with heterogeneous consumers (Berry and Haile 2010); static games of incomplete information with player-specific regressors excluded from interaction effects (Lewbel and Tang 2012); and matching games with unobserved heterogeneity (Fox and Yang 2012).

Using a special regressor to identify coefficients in binary regressions with heteroskedastic errors typically requires additional conditions on the support of the special regressor. For instance, in the case with mean-independent errors, identification of linear coefficients requires the support of special regressors to be at least as large as that of errors. (Khan and Tamer 2010) argued that point identification of coefficients under mean independent errors is lost whenever the support of special regressor is bounded.<sup>5</sup> They also showed that when support of special regressor is unbounded, Fisher information for coefficients becomes zero when the second moment of regressors is finite.

The econometrics literature on semiparametric binary regressions has largely been silent about how to use special regressors in combination of alternative stochastic restrictions on errors that require less stringent conditions on the support of special regressors. Perhaps the only two exceptions are (Magnac and Maurin 2007) and (Chen 2005). (Magnac and Maurin 2007) introduced a new restriction on the tail behavior of latent

---

<sup>5</sup>They showed in a stylized model that there is no informative partial identification result for the intercept in this case.

utility distribution outside the support of special regressors. They established identification and positive Fisher information for coefficients under such restrictions. Nonetheless the tail condition they use is not directly linked to more conventional stochastic restrictions on heteroskedastic errors, such as median independence or conditional symmetry. (We show in Appendix B that the tail conditions in (Magnac and Maurin 2007) and the conditional symmetry considered in our paper are non-nested.) (Chen 2005) estimated a binary regression with special regressors and conditional symmetric errors, and proved it is root-N consistent and asymptotically normal. He does not quantify the incremental contribution of special regressors to the identifying power and Fisher information of the model. To our knowledge, no previous work has measured the incremental identifying power or Fisher information for coefficients in binary regressions due to special regressors under assumptions studied in this paper.

We fill in this gap in the literature by deriving several new results. First, we quantify the change in identifying power of the model due to the presence of special regressors under median independent or conditionally symmetric errors. This is done following the approach used in (Manski 1988), which amounts to comparing the size of the set of states where the propensity scores can be used for distinguishing true coefficients from other elements in the parameter space. For the model with median independent errors, we find that further restricting one of the regressors to be a special one does *not* improve the identifying power for coefficients. For the model with conditionally symmetric errors, we find that using a special regressor *does* add to the identifying power for coefficients in the sense that it leads to an additional set of (paired states) that can be used for recovering the true coefficients. This is a useful additional insight, because (Manski 1988) only showed that, in the absence of a special regressor, the stronger restriction of conditional symmetry adds no identifying power relative to the weaker restriction of median independence.

Second, we show how the presence of a special regressor contributes to the information for coefficients in these two semiparametric binary regressions with heteroskedastic errors. For models with median-independent errors, we find the information for coefficients remains zero even after one of the regressors is required to be special. In comparison, for models with conditionally symmetric errors, the presence of a special regressor does yield positive information for coefficients. We provide some intuition for such positive information in this case, and propose a new two-step extremum estimator. Asymptotic properties of the estimator are derived and some monte carlo evidence for its performance is reported.

Our third set of results (Section 3.3) provides a more positive perspective on the role of special regressors in structural analyses. We argue that, even though a special regressor does not add to identifying power or information for coefficients when heteroskedastic errors are only required to be median independent, it is instrumental for recovering the distribution of the heteroskedastic error. This in turn can be used to predict counterfactual choice probabilities; and helps to recover the average structural function as defined in (Blundell and Powell 2003) as long as the support of the special regressor is large enough.

This paper contributes to a broad econometrics literature on identification, inference and information of semiparametric limited response models with heteroskedastic errors. A partial list of other papers that discussed related topics include (Chamberlain 1986), (Chen and Khan 2003), (Cosslett 1987), (Magnac and Maurin 2007), (Manski 1988) and (Zheng 1995) (which all studied semiparametric binary regressions with various specifications of heteroskedastic errors); as well as (Newey and McFadden 1994), (Powell 1994) and (Ichimura and Lee 2010) (which discussed estimation and inference methods in a more general framework that subsumes the models considered here).

## 2 Preliminaries

Consider a binary regression:

$$Y = 1\{\epsilon \leq X\beta - V\}. \quad (1)$$

where  $X \in \mathbb{R}^K, V \in \mathbb{R}$  and  $\epsilon \in \mathbb{R}^1$  and the first coordinate in  $X$  is a constant. We use upper cases for random variables and lower cases for their realizations. Let  $F_R, f_R, \Omega_R$  denote the distribution, the density and the support of a random vector  $R$  respectively, and let  $F_{R_1|R_2}, f_{R_1|R_2}$  and  $\Omega_{R_1|R_2}$  denote conditional distributions, densities and supports in the data-generating process (DGP). Assume the marginal effect of  $V$  is known to be negative, and is set to  $-1$  as a scale normalization. We maintain the following exclusion restriction throughout the paper.

**CI** (*Conditional Independence*)  $V$  is independent from  $\epsilon$  given any  $x \in \Omega_X$ .

For the rest of the paper, we also refer to this condition as an “exclusion restriction”, and use the terms “special regressors” and “excluded regressors” interchangeably. Let  $\Theta$  be the parameter space for  $F_{\epsilon|X}$  (i.e.  $\Theta$  is a collection of all conditional distributions of errors that satisfy the model restrictions imposed on  $F_{\epsilon|X}$ ). The distribution  $F_{V|X}$  is directly identifiable from data and considered known in the identification exercise. Let  $Z \equiv (X, V)$ , and let  $p(z)$  denote  $\Pr(Y = 1|z)$  (which is directly identifiable from data). Let  $(Z, Z')$  be a pair of independent draws from the same marginal distribution  $F_Z$ . Assume the distribution of  $Z$  has positive density with respect to a  $\sigma$ -finite measure, which consists of the counting measure for discrete coordinates and the Lebesgue measure for continuous coordinates.

For a generic pair of coefficients and the nuisance distribution  $(b, G_{\epsilon|X}) \in \mathbb{R}^K \otimes \Theta$ , define  $\xi(b, G_{\epsilon|X}) \equiv \{z : p(z) \neq \int 1(\epsilon \leq xb - v) dG_{\epsilon|x}\}$  and  $\tilde{\xi}(b, G_{\epsilon|X}) \equiv$

$$\left\{ (z, z') : (p(z), p(z')) \neq \left( \int 1(\epsilon \leq xb - v) dG_{\epsilon|x}, \int 1(\epsilon \leq x'b - v') dG_{\epsilon|x'} \right) \right\}. \quad (2)$$

In words, the set  $\xi(b, G_{\epsilon|X})$  consists of states for which propensity scores implied by  $(b, G_{\epsilon|X})$  differ from that generated in (and directly identifiable from) the true data-generating process (DGP) characterized by  $(\beta, F_{\epsilon|X})$ . In comparison, the set  $\tilde{\xi}$  in (2) consists of pairs of states such that the implied pairs of propensity scores differ from those generated by the true DGP. We say  $\beta$  is *identified relative to*  $b \neq \beta$  if

$$\int 1\{z \in \xi(b, G_{\epsilon|X})\} dF_Z > 0 \text{ or } \int 1\{(z, z') \in \tilde{\xi}(b, G_{\epsilon|X})\} dF_{(Z, Z')} > 0$$

for all  $G_{\epsilon|X} \in \Theta$ .

Of course, identification of  $\beta$  hinges on the restrictions imposed on  $\Theta$ , the parameter space for the error distribution given  $X$ . In the following sections, we discuss identification of  $\beta$  when CI is paired with *one* of the following two stochastic restrictions on error distributions.

**MI** (*Median Independence*) For all  $x$ ,  $\epsilon$  is continuously distributed with  $\text{Med}(\epsilon|x) = 0$  and the conditional density is strictly positive in an open neighborhood around 0.

**CS** (*Conditional Symmetry*) For all  $x$ ,  $\epsilon$  is continuously distributed with positive densities over the support  $\Omega_{\epsilon|x}$  and  $F_{\epsilon|x}(t) = 1 - F_{\epsilon|x}(-t)$  for all  $t \in \Omega_{\epsilon|x}$ .

We also discuss the semiparametric efficiency bound for  $\beta$  under these two sets of assumptions in Sections 3.2 and 4.2. In essence, this amounts to finding smooth parametric submodels which are nested in the semiparametric models and which have the least Fisher information for  $\beta$ . Formally, the semiparametric efficiency bound is defined as follows. Let  $\mu$  denote some measure on  $\{0, 1\} \otimes \Omega_Z$  such that  $\mu(\{0\} \otimes \omega) = \mu(\{1\} \otimes \omega) = F_Z(\omega)$ , where  $\omega$  is a Borel subset of  $\Omega_Z$ . A *path* that goes through  $F_{\epsilon|X}$  is a function  $\lambda(\epsilon, x; \delta')$  such that  $\lambda(\epsilon, x; \delta) = F_{\epsilon|x}(\epsilon)$  for some  $\delta \in \mathbb{R}$ , and  $\lambda(\cdot, \cdot; \delta')$  satisfies restrictions on the conditional distribution of errors in  $\Theta$  for all  $\delta'$  in an open neighborhood around  $\delta$ . Let  $f_\lambda(y|z; b, \delta')$  denote the probability mass function of  $Y$  conditional on  $z$  and given coefficients  $b$  as well as a conditional distribution of errors  $\lambda(\cdot, \cdot; \delta')$ . A *smooth* parametric submodel is characterized by a path  $\lambda$  such that there exists  $\{(\psi_k)_{k \leq K}, \psi_\lambda\}$  such that

$$f_\lambda^{1/2}(y|z; b, \delta') - f_\lambda^{1/2}(y|z; \beta, \delta) = \sum_k \psi_k(y, z) (b - \beta) + \psi_\lambda(y, z) (\delta' - \delta) + r(y, z; b, \delta') \quad (3)$$

with

$$(\|b - \beta\| + \|\delta' - \delta\|)^{-2} \int r^2(y, z; b, \delta') d\mu \rightarrow 0 \text{ as } b \rightarrow \beta \text{ and } \delta' \rightarrow \delta. \quad (4)$$

The *path-wise partial information* for the  $k$ -th coordinate in  $\beta$  is

$$I_{\lambda, k} \equiv \inf_{(\{\alpha_j\}_{j \neq k}, \alpha_\lambda)} 4 \int \left( \psi_k - \sum_{j \neq k} \alpha_j \psi_j - \alpha_\lambda \psi_\lambda \right)^2 d\mu. \quad (5)$$

The *information* for  $\beta_k$  is the infimum of  $I_{\lambda, k}$  over all  $\lambda$  in smooth parametric submodels.

### 3 Exclusion plus Median Independence

This section discusses the identification and information for  $\beta$  in a binary regression under CI and MI. The model differs from that in (Manski 1988) and (Horowitz 1992) due to the presence of a special regressor  $V$ , which satisfies the exclusion restriction by being conditionally independent from the error term. It also differs from that considered in (Lewbel 2000) and (Khan and Tamer 2010), for the stochastic location restriction on the error distribution is median independence rather than mean independence. To the best of our knowledge, there exists no previous work that discusses the identification and the Fisher information of coefficients in such a model.

#### 3.1 Identification

This subsection shows the exclusion restriction in CI contributes no identifying power for the recovery of  $\beta$  under MI. We formalize this result in Proposition 1 below by noting that the set of states in  $\Omega_Z$  that could help econometricians detect a vector  $b \neq \beta$  from  $\beta$  under MI remains unchanged, regardless of whether an exclusion restriction is added to one of the regressors.

**Proposition 1** *Suppose CI and MI hold in (1). Then  $\beta$  is identified relative to  $b$  if and only if  $\Pr\{z \in Q_b\} > 0$ , where  $Q_b \equiv \{z : x\beta \leq v < xb \text{ or } xb \leq v < x\beta\}$ .*

**Remark 1.1.** For a model under MI but without CI (i.e.  $F_{\epsilon|X,V}(0) = 1/2$  where the distribution of  $\epsilon$  depends on  $V$  as well as  $X$ ) (Manski 1988) showed  $Q_b$  is the set of states that can be used to detect  $b \neq \beta$  from  $\beta$ , based on observed propensity scores. Thus Proposition 1 suggests adding the exclusion restriction (CI) to a model with median independence does not improve the identifying power for recovery of  $\beta$ , as the set of states that help identify  $\beta$  relative to  $b$  remains unchanged. The intuition for such an equivalence builds on two observations. First, if states in  $Q_b$  can help identify  $\beta$  relative to  $b$  under the weaker assumption of median independence alone in (Manski 1988), then they must also do so when an additional exclusion restriction is invoked. Second, if  $\Pr\{Z \in Q_b\} = 0$ , then certain distribution of structural errors  $\tilde{G}_{\epsilon|X} \neq F_{\epsilon|X}$  can be constructed to satisfy CI and MI and, when paired with  $b \neq \beta$ , can also rationalize the propensity scores as generated by the DGP. The proof of Proposition 1 differs qualitatively from that of Manski's result in that the construction of such a distribution  $\tilde{G}_{\epsilon|X}$  needs to respect additional exclusion restrictions.

**Remark 1.2.** While not helping with identification of  $\beta$ , the exclusion restriction in CI on the other hand does help with the identification of the error distribution  $F_{\epsilon|X}$ , which in turn is useful for predicting propensity scores under counterfactual changes, and for

recovering average structural marginal effects of  $X$ . We provide more detailed discussions on the role of excluded regressors in such structural analyses in Section 3.3.

**Remark 1.3.** The result in Proposition 1 is also related to (Khan and Tamer 2010). To see this, suppose  $X$  consists of continuous coordinates only. Then  $\Pr\{Z \in Q_b\} \rightarrow 0$  as  $b$  converges to  $\beta$ . That is,  $Q_b$  becomes a “thin set” as  $b$  approaches  $\beta$ .

It also follows from Proposition 1 that following conditions for point identification of  $\beta$  in (Manski 1988) are also sufficient for exactly recovering  $\beta$  in our current model with CI and MI.

**SV** (*Sufficient Variation*) For all  $x$ ,  $V$  is continuously distributed with positive densities over  $\Omega_{V|x}$ , which includes  $x\beta$  in the interior.

**FR** (*Full Rank*),  $\Pr\{X\gamma \neq 0\} > 0$  for all nonzero vector  $\gamma \in \mathbb{R}^K$ .

It is worth noting that SV can be satisfied if the support  $\Omega_V$  is bounded. It differs from the large support conditions needed to point identify  $\beta$  under CI and mean independence of errors. FR is the typical full-rank condition analogous to that used in (Manski 1988). Recall the first coordinate in  $X$  is a constant. Hence FR implies that there exists no nonzero  $\tilde{\gamma}$  in  $\mathbb{R}^{K-1}$  and  $c \in \mathbb{R}$  with  $\Pr\{X_{-1}\tilde{\gamma} = c\} = 1$ . For any  $b \neq \beta$ , FR implies  $\Pr\{X(\beta - b) \neq 0\} > 0$ . Without loss of generality, suppose  $\Pr\{X\beta < Xb\} > 0$ . Under SV, for any  $x$  with  $x\beta < xb$ , there exists an interval of  $v$  with  $x\beta \leq v < xb$ . This implies  $\Pr\{Z \in Q_b\} > 0$  and thus  $\beta$  is identified relative to all  $b \neq \beta$ .

For estimation, we propose a new extremum estimator for  $\beta$  that differs qualitatively from the Maximum Score estimator in (Manski 1985). Our estimator builds on the following identification argument.

**Corollary 1** (*Proposition 1*) Suppose CI, MI, SV and FR hold in (1), and  $\Pr(X\beta = V) = 0$ . Then

$$\beta = \arg \min_b \mathbb{E}_Z[1\{p(Z) \geq \frac{1}{2}\}(Xb - V)_- + 1\{p(Z) < \frac{1}{2}\}(Xb - V)_+] \quad (6)$$

where  $(\cdot)_+ \equiv \max\{\cdot, 0\}$  and  $(\cdot)_- \equiv -\min\{\cdot, 0\}$ .

Let  $n$  denote the sample size and let  $\hat{p}_i$  denote kernel estimator for  $\mathbb{E}(Y|Z = z_i)$ . An alternative estimator is

$$\tilde{\beta} \equiv \arg \min \sum_i \kappa(\hat{p}_i - \frac{1}{2})(x_i b - v_i)_- + \kappa(\frac{1}{2} - \hat{p}_i)(x_i b - v_i)_+ \quad (7)$$

where the weight function  $\kappa : \mathbb{R} \rightarrow [0, 1]$  satisfies:  $\kappa(t) = 0$  for all  $t \leq 0$ ;  $\kappa(t) > 0$  for all  $t > 0$ ; and  $\kappa$  is increasing over  $[0, +\infty)$ .

A few remarks about the asymptotic properties of the estimator as well as its comparison with the Maximum Score estimator are in order. If either SV or FR fails, then  $\beta$  is only set-identified and the objective function in (6) have multiple minimizers.<sup>6</sup> The estimator in (7), which is a random set of coefficients, can be shown to be consistent for the identified set (with the metric between set being the Hausdorff metric), provided the sample objective function in (7) converges to its population counterpart in (6) uniformly over the parameter space for coefficients. Furthermore, we conjecture it is possible to show cubic rate of convergence under conditions cited in (Chernozhukov, Hong, and Tamer 2007). We may also want to require  $\kappa$  to be twice continuously differentiable with bounded derivatives in an open neighborhood around 0 for technical convenience in the proof of the asymptotic properties of  $\tilde{\beta}$ .

Compared with the Maximum Score estimator,  $\tilde{\beta}$  in (7) appears to have computational advantages once the propensity scores are estimated, as the argument  $b$  enters the estimand continuously through  $(\cdot)_-$  and  $(\cdot)_+$ , as opposed to in the indicator function in Maximum Score estimator. The flip side of our estimator is that it does require the choice of smoothing parameters in the estimation of  $\hat{p}_i$ . We leave further investigation of the asymptotic properties for future work.

### 3.2 Zero Fisher Information

This subsection shows information for  $\beta$  under CI and MI is zero, provided  $Z = (X, V)$  has finite second moments and some regularity condition on the coefficient and the error distribution below holds. In addition to Section 3.1, our finding herein provides an alternative way to formalize the equivalence between the two models (i.e. binary regressions with "MI alone" versus "MI and CI") when it comes to estimating  $\beta$ .

**RG** (*Regularity*) For each  $(b, G_{\epsilon|X})$  in the parameter space, there exists a measurable function  $q : \{0, 1\} \otimes \Omega_Z \rightarrow \mathbb{R}$  such that  $|\partial f^{1/2}(y, z; \eta, G_{\epsilon|X}) / \partial b| \leq q(y, z)$  for all  $\eta$  in an neighborhood around  $b$ ; and  $\int q^2(y, z) d\mu < \infty$ .

RG is needed to establish mean-square differentiability of the square-root likelihood of  $(y, x)$  with respect to  $b$  for each  $G_{\epsilon|X}$ . Let  $\Theta$  denote the parameter space for the distribution of  $\epsilon$  given  $X$ , which needs to satisfy CI, MI and RG now. We show that a set of paths similar to those considered in (Chamberlain 1986) makes the information for  $\beta$  under CI, MI and RG zero. Let  $\Lambda$  consist of paths

$$\lambda(\varepsilon, x; \delta') \equiv F_{\epsilon|x}(\varepsilon) [1 + (\delta' - \delta) h(\varepsilon, x)], \quad (8)$$

---

<sup>6</sup>The identified set for  $\beta$  is defined as the set of all coefficients that could generate propensity scores identical to that in the DGP for *all*  $z$  when paired with some nuisance parameter  $F_{\epsilon|X}$  that satisfies CI and MI.



where  $F_{\epsilon|X}$  is the true conditional distribution in DGP from  $\Theta$ ; and  $h : \mathbb{R}^{K+1} \rightarrow \mathbb{R}$  is continuously differentiable, is zero outside of some compact set; and satisfies  $h(0, x) = 0$  for all  $x \in \mathbb{R}^K$ . Such a set of paths differs from those leading to zero information of  $\beta$  in a model without exclusion restrictions. In the latter case, the paths that lead to zero information is  $\lambda(\epsilon, x; \delta') \equiv F_{\epsilon|x, v}(\epsilon) \left[ 1 + (\delta' - \delta) \tilde{h}(\epsilon, x, v) \right]$  with  $\tilde{h} : \mathbb{R}^{K+2} \rightarrow \mathbb{R}$  continuously differentiable; is zero outside of some compact set; and satisfying  $\tilde{h}(0, x, v) = 0$ . (See (Chamberlain 1986) for details.)

Using arguments similar to (Chamberlain 1986), we can show  $\lambda(\cdot, \cdot; \delta')$  in (8) is in  $\Theta$  for  $\delta'$  close enough to  $\delta$ . Besides,  $f_\lambda^{1/2}(\cdot; b, \delta')$  is mean-square differentiable at  $(b, \delta') = (\beta, \delta)$  with:

$$\psi_k(y, z) \equiv \frac{1}{2} \left\{ y F_{\epsilon|x}(w)^{-1/2} - (1-y) [1 - F_{\epsilon|x}(w)]^{-1/2} \right\} f_{\epsilon|x}(w) x_k \quad (9)$$

$$\psi_\lambda(y, z) \equiv \frac{1}{2} \left\{ y F_{\epsilon|x}(w)^{-1/2} - (1-y) [1 - F_{\epsilon|x}(w)]^{-1/2} \right\} F_{\epsilon|x}(w) h(w, x) \quad (10)$$

where  $w$  is a shorthand for  $x\beta - v$ . Again note the excluded regressor  $v$  is dropped from  $F_{\epsilon|x}$  and  $f_{\epsilon|x}$  due to CI.

**Proposition 2** *Suppose CI, MI, SV, FR and RG hold in (1);  $Z$  has finite second moments; and  $\Pr(X\beta = V) = 0$ . Then the information for  $\beta_k$  is zero for all  $k \leq K$ .*

Proof of Proposition 2 is similar to that of Theorem 5 in (Chamberlain 1986) for the case of binary regressions under MI alone, and therefore is omitted for brevity. It suffices to note that the main difference the proof of Proposition 2 and that of Theorem 5 in (Chamberlain 1986) is that the path leading to zero information for  $\beta$  in the current model (with an additional assumption CI) has to respect the additional exclusion restriction.

A few remarks related to this zero information result are in order. First, zero information for  $\beta_k$  under CI and MI is related to two observations: there is no incremental identifying power for  $\beta$  from CI given MI; and there is zero information for  $\beta_k$  under the MI restriction alone. Second, root-n estimator for  $\beta$  is possible when the second moments for regressors are infinite. In such a case, (Khan and Tamer 2010) showed that parametric rate can be attained for  $\beta$  under mean independence and exclusion restriction. A similar result can be shown for the current case with median independence plus exclusion restriction as well. Third, if there are multiple special regressors (i.e.  $V$  is a vector rather than scalar) and if one of the coefficients for  $V$  is normalized to have absolute value 1, then the information for coefficients of the other regressors satisfying the exclusion restriction is positive, and root-n estimation of coefficients for  $V$  can be defined (e.g. using average-derivative-type of estimators).

### 3.3 Further Discussions

We have shown the presence of a special regressor satisfying the exclusion restriction, when combined with median independence, does not help to improve the identification or information of coefficients for other regressors (either in the sense of augmenting the set of states usable for identification, or in the sense of yielding positive information). This subsection concludes with a positive perspective on the role of excluded regressors in identifying counterfactual choice probabilities as well as average structural functions.

First, exclusion restrictions *does* help recover the heteroskedastic error distributions, which in turn helps recover counterfactual choice probabilities. To see how to recover  $F_{\epsilon|X}$  under median independence, note  $\mathbb{E}(Y|x, v) = F_{\epsilon|x}(x\beta - v)$  by construction. With  $\beta$  identified,  $F_{\epsilon|x}(t)$  can be recovered for all  $t$  over the support of  $X\beta - V$  give  $X = x$  as  $\mathbb{E}(Y|X = x, V = x\beta - t)$ . To see how this helps with counterfactual predictions, consider a stylized model of retirement decisions. Let  $Y = 1$  if the individual decides to retire and  $Y = 0$  otherwise. The decision is given by:

$$Y = 1\{X_1\beta_1 + X_2\beta_2 - V \geq \epsilon\}$$

where  $X_1, X_2$  are log age and health status respectively and  $V$  denotes the total market value of individual's cumulated assets. Further suppose that conditional on age and health status, asset values are uncorrelated with idiosyncratic elements (e.g. unobserved family factors such as money and energy spent on offsprings). Suppose a researcher is interested in predicting retirement patterns among another population of senior workers, not observed in data, with the same  $\beta_1$  and  $F_{\epsilon|X_1, X_2}$  but the tastes for (or weights assigned to) health status is changed to  $\tilde{\beta}_2$  so that in this new population  $\tilde{\beta}_2 > \beta_2$ . Then knowledge of  $F_{\epsilon|X_1, X_2}$  as well as  $\beta_1, \beta_2$  helps at least bound the counterfactual retirement probabilities conditional on  $(X_1, X_2, V)$ . If the magnitude of the difference between  $\tilde{\beta}_2$  and  $\beta_2$  is also known, then point-identification of such a counterfactual conditional retirement probability is also attained for  $(x_1, x_2, v)$ , provided the support  $\Omega_{V|x_1, x_2}$  is large enough. (That is, the index  $x_1\beta_1 + x_2\tilde{\beta}_2 - v$  is within the support of  $X_1\beta_1 + X_2\beta_2 - V$  given  $(x_1, x_2)$ .)

Second, exclusion restrictions *does* help identify the average structural function, as defined in (Blundell and Powell 2003), under the large support condition of  $V$ . To see this, note the average structural function is defined as  $G(x.v) \equiv \int 1\{\epsilon \leq x\beta - v\}dF_{\epsilon}(\epsilon) = \Pr(\epsilon \leq x\beta - v)$ . If  $\Omega_{V|x} = \mathbb{R}^1$  for all  $x \in \Omega_X$ , then

$$G(x.v) = \int \varphi(s, x, v)dF_X(s),$$

where

$$\varphi(s, x, v) \equiv \mathbb{E}[Y|X = s, V = v + (s - x)\beta] = F_{\epsilon|s}(x\beta - v).$$

With  $\beta$  identified previously,  $\varphi(s, x, v)$  can be constructed as long as the support of  $V$  spans the real line for all  $x$ . If this large support condition fails, then identification of

$G(x, v)$  is lost at any  $(x, v)$  such that there exists  $s \in \Omega_X$  where  $v + (s - x)\beta$  falls outside of the support  $\Omega_{V|s}$ .

Based on the analog principle, we propose a natural estimator for the average structural function as follows:

$$\hat{G}(x, v) \equiv \sum_{i=1}^n \hat{\varphi}(x_i, x, v)$$

where

$$\hat{\varphi}(x_i, x, v) \equiv \frac{\sum_{j \neq i} y_j \mathcal{K}_\sigma \left( x_j - x_i, v_j - (v + (x_i - x)\tilde{\beta}) \right)}{\sum_{j \neq i} \mathcal{K}_\sigma \left( x_j - x_i, v_j - (v + (x_i - x)\tilde{\beta}) \right)}$$

with  $\mathcal{K}_\sigma(\cdot) \equiv \sigma^{-(k+1)} \mathcal{K}(\cdot/\sigma^{k+1})$  where  $\mathcal{K}$  is a product kernel; and  $\tilde{\beta}$  being some first-stage preliminary estimator such as the one defined in (7), or the maximum score estimator as proposed in (Manski 1985).

## 4 Exclusion plus Conditional Symmetry

This section discusses identification and information of  $\beta$  under CI while the location restriction of median independence (MI) is replaced by the stronger location and shape restriction of conditional symmetry (CS). To motivate the use of CS in binary regressions, consider the latent utility associated with binary actions are  $h_j(x) + \epsilon_j$  for  $j \in \{0, 1\}$ . Then the action  $Y \in \{0, 1\}$  is given by  $Y = 1\{h_1(X) + \epsilon_1 \geq h_0(X) + \epsilon_0\} = 1\{h^*(X) + \epsilon^* \geq 0\}$ , where  $h^* \equiv h_1 - h_0$  and  $\epsilon^* \equiv \epsilon_1 - \epsilon_0$ . As long as  $\epsilon_1$  and  $\epsilon_0$  are independent draws from the same marginal  $F_{\epsilon_j|X}$ , the normalized error  $\epsilon^*$  must be symmetrically distributed around 0 given  $X$ .

In Section 4.1, we characterize a subset in the support of paired states  $(Z, Z')$  that help distinguish  $\beta$  from some  $b \neq \beta$  based on observed propensity scores. Building on this result, we then specify sufficient conditions for the point identification of  $\beta$ . In Section 4.2 we show the Fisher information for  $\beta$  is zero under mild regularity conditions. We then conclude this section with the introduction of a root-N estimator for  $\beta$ .

### 4.1 Identification

Our first finding is that replacing MI with CS while maintaining CI *does* help with the identification of  $\beta$ . Let  $X$  consist of subvectors of continuous and discrete coordinates, denoted by  $X_c$  and  $X_d$  respectively. Let  $\Theta_{CS}$  denote parameter space for the distribution of  $\epsilon$  given  $X$  under the restrictions of CI and CS. We need further smoothness restrictions on  $\Theta_{CS}$  due to continuous coordinates in  $X_c$ .

**UC** (*Uniform Continuity*) For any  $\eta > 0$  and  $(x, \varepsilon)$ , there exists  $\delta_\eta(x, \varepsilon) > 0$  such that for all  $G_{\varepsilon|X} \in \Theta_{CS}$ ,

$$|G_{\varepsilon|\tilde{x}}(\tilde{\varepsilon}) - G_{\varepsilon|x}(\varepsilon)| \leq \eta \text{ whenever } \|\tilde{x} - x\|^2 + \|\tilde{\varepsilon} - \varepsilon\|^2 \leq \delta_\eta(x, \varepsilon).$$

This condition requires the pointwise continuity in  $(x, \varepsilon)$  to hold uniformly over  $\Theta_{CS}$ , in the sense that the same  $\delta_\eta(x, \varepsilon)$  is used to satisfy the “ $\delta$ - $\eta$ -neighborhood” definition of pointwise continuity at  $(x, \varepsilon)$  for all elements in  $\Theta_{CS}$ .<sup>7</sup> Such a condition is needed because the identification of  $\beta$  relative to  $b \neq \beta$  states that  $b$  cannot be paired with *any*  $G_{\varepsilon|X} \neq F_{\varepsilon|X}$  in  $\Theta_{CS}$  to generate propensity scores that are identical to the true ones in DGP at any paired states  $(z, z')$ . It is a technicality that arises as we try to implement definition of identification in (Manski 1988) in the presence of continuous coordinates in  $X$ . A sufficient condition for UC is that all  $G_{\varepsilon|X}$  in  $\Theta_{CS}$  are Lipschitz-continuous with their modulus uniformly bounded by a finite constant.

To formally quantify the incremental identifying power due to CS, define:

$$R_b(x) \equiv \left\{ (v_i, v_j) : x\beta < \frac{v_i + v_j}{2} < xb \text{ or } x\beta > \frac{v_i + v_j}{2} > xb \right\} \quad (11)$$

for any  $x$ . Let  $F_{V_i, V_j|X}$  denote the joint distribution of  $V_i$  and  $V_j$  drawn independently from the same marginal distribution  $F_{V_i|X}$ . In addition we also need the joint distribution of  $V$  and  $X_c$  given  $X_d$  to be continuous.

**CT** (*Continuity*) For any  $x_d$ , the distribution  $F_{V, X_c|x_d}$  is continuous with positive densities a.e. with respect to the Lebesgue measure.

Under CT, if  $\Pr\{V_i \in \mathcal{A}|(x_c, x_d)\} > 0$  for any set  $\mathcal{A}$  then  $\Pr\{V_i \in \mathcal{A}|(\tilde{x}_c, x_d)\} > 0$  for  $\tilde{x}_c$  close enough to  $x_c$ .

**Proposition 3** Under CI, CS, UC and CT,  $\beta$  is identified relative to  $b$  if and only if either (i)  $\Pr\{Z \in Q_b\} > 0$ ; or (ii) there exists a set  $\omega$  open in  $\Omega_X$  such that for all  $x \in \omega$ ,

$$\int 1\{(v_i, v_j) \in R_b(x)\} dF_{V_i, V_j|x} > 0. \quad (12)$$

Proof of Proposition 3 is presented in the appendix. To illustrate the intuition behind this result, focus on the simpler case where  $X$  only consists of discrete regressors (i.e.  $X = X_d$ ). For a fixed  $b \neq \beta$ , consider a pair  $(z_i, z_j) \in \tilde{Q}_{b,S}$ , where

$$\tilde{Q}_{b,S} \equiv \{(z_i, z_j) : x_i = x_j \text{ and } (v_i, v_j) \in R_b(x_i)\}.$$

<sup>7</sup>An alternative way to formulate UC is that for any  $\eta > 0$  and  $(x, \varepsilon)$ , the infimum of  $\delta_\eta(x, \varepsilon; G_{\varepsilon|X})$  (i.e. the radius of neighborhood around  $x$  in the definition of pointwise continuity) over  $G_{\varepsilon|X} \in \Theta_{CS}$  is bounded away from zero by a positive constant.

Then either " $x_i\beta - v_i < -(x_j\beta - v_j)$  and  $x_ib - v_i > -(x_jb - v_j)$ " or " $x_i\beta - v_i > -(x_j\beta - v_j)$  and  $x_ib - v_i < -(x_jb - v_j)$ " for  $(z_i, z_j) \in \tilde{Q}_{b,S}$ . In the former case, the actual propensity scores are such that  $p(z_i) + p(z_j) < 1$  while in contrast the propensity scores at  $z_i$  and  $z_j$  as implied by  $b \neq \beta$  and any  $G_{\epsilon|X} \in \Theta_{CS}$  must add up to be greater than 1. This suggests any such pair  $(z_i, z_j)$  should help distinguish  $\beta$  from  $b \neq \beta$ , as the sign of  $p(z_i) + p(z_j) - 1$  differs from that of  $(x_ib - v_i) + (x_jb - v_j)$ . Thus if condition (ii) in Proposition 3 holds for  $b$  and all coordinates in  $X$  are discrete, then  $\Pr\{(Z_i, Z_j) \in \tilde{\xi}(b, G_{\epsilon|X})\} > 0$  for all  $G_{\epsilon|X} \in \Theta_{CS}$ . When both conditions (i) and (ii) fail,  $\beta$  is not identified to  $b$ , because some  $G_{\epsilon|X} \neq F_{\epsilon|X}$  can be constructed so that  $(b, G_{\epsilon|X})$  is observationally equivalent to the true parameters  $(\beta, F_{\epsilon|X})$  (in the sense of generating the same pairs of propensity scores  $p(z)$ ,  $p(z')$  as the latter almost surely).

Conditions UC and CT are necessary for extending this intuition to the more general case with continuous coordinates in  $X$ . In such a case, new challenges arise because, with  $z_i$  and  $z_j$  being i.i.d. draws from the marginal  $F_Z$  and with  $X$  containing continuous components, the paired states in  $\tilde{Q}_{b,S}$  occur with zero probability for all  $b \neq \beta$ . With the smoothness restrictions in UC and CT, the inequalities leading to identification of  $\beta$  in the previous paragraph (where we consider the "all-discrete" case with  $X = X_d$ ) also hold for paired states in a small " $\delta$ -expansion" (or a superset) of  $\tilde{Q}_{b,S}$  defined as:

$$\tilde{Q}_{b,S}^\delta \equiv \{(z, \tilde{z}) : x_d = \tilde{x}_d \wedge \|\tilde{x}_c - x_c\| \leq \delta \wedge (v, \tilde{v}) \in R_b(x_c, x_d)\},$$

provided  $\delta > 0$  is small enough. To identify  $\beta$  relative from  $b$ , it then suffices to require  $\tilde{Q}_{b,S}^\delta$  to have positive probability for such small  $\delta$ , which is possible with continuous coordinates in  $X$ .

(Manski 1988) showed that, in the absence of excluded regressors, replacing median independence with conditional symmetry does not add to the identifying power for recovering  $\beta$ . In particular, he showed that in both cases the set of states that could help distinguish  $\beta$  from  $b \neq \beta$  remains the same. Our finding in Proposition 3 shows this equivalence breaks down when the vector of states contain a special regressor that satisfies the exclusion restriction in CI. In this case, and replacing MI with CS does lead to an additional set  $R_b$  of paired states which could help identify  $\beta$  relative to  $b$ .

Note that by construction the conditions that lead to identification of  $\beta$  under CI and MI must also be sufficient for its identification under the stronger conditions of CI and CS. To see this, note under FR, for all  $b \neq \beta$ , there exists an open set  $\omega \subseteq \Omega_X$  with  $x\beta \neq xb$  for all  $x \in \omega$ . SV then implies either  $\int 1\{x\beta < \frac{v_i+v_j}{2} < xb\}dF_{V_i, V_j|x} > 0$  or  $\int 1\{xb < \frac{v_i+v_j}{2} < x\beta\}dF_{V_i, V_j|x} > 0$  for all  $x \in \omega$ . This is because under SV,  $V_i$  and  $V_j$  are independent draws from  $F_{V|x}$  and both fall in an open neighborhood around  $x\beta$  with positive probability. Identification of  $\beta$  then follows from Proposition 3.

## 4.2 Positive Fisher Information

We now show how CS combined with CI leads to positive information for  $\beta$ . (Zheng 1995) showed that a binary regression under CS but with no excluded regressor has zero information for  $\beta$ . In contrast, we show here that with excluded regressors, conditional symmetry of errors does yield positive information for  $\beta$  under mild regularity conditions added below. Our finding suggests there exists root-N consistent estimators for  $\beta$ .

**CS'** *CS holds and there exists  $c > 0$  such that  $f_{\epsilon|x}(\epsilon) \geq c$  for all  $\epsilon \in (-\epsilon^*, \epsilon^*)$  for all  $x \in \Omega_X$ .*

**RG'** *RG holds and for any  $\bar{w}$  such that  $\Pr(X \in \bar{w}) > 0$ , there exists no nonzero  $\alpha \in \mathbb{R}^K$  such that  $\Pr\{X\alpha = 0 | X \in \bar{w}\} = 1$ .*

Let  $\Lambda$  consist of paths  $\lambda : \Omega_{\epsilon, X} \otimes \mathbb{R} \rightarrow [0, 1]$  such that (i) for some  $\delta \in \mathbb{R}^1$ ,  $\lambda(\epsilon, x; \delta) = F_{\epsilon|x}(\epsilon)$  for all  $\epsilon, x$ ; (ii) for  $\eta$  in an neighborhood around  $\delta$ ,  $\lambda(\epsilon, x; \eta)$  is a conditional distribution of  $\epsilon$  given  $X$  that satisfies:

$$\lambda(\epsilon, x; \eta) = 1 - \lambda(-\epsilon, x; \eta) \text{ for all } \epsilon, x \in \Omega_{\epsilon, X}; \quad (13)$$

and (iii) the square-root density  $f_\lambda^{1/2}(y, z; b, \eta)$  is mean-square differentiable at  $(b, \eta) = (\beta, \delta)$ , with the pathwise derivative with respect to  $\eta$  being:

$$\psi_\lambda(y, z) \equiv \frac{1}{2} \left\{ y F_{\epsilon|x}(w)^{-1/2} - (1 - y) [1 - F_{\epsilon|x}(w)]^{-1/2} \right\} \lambda_\eta(w, x; \delta) \quad (14)$$

where  $w \equiv x\beta - v$  and  $\lambda_\eta(\epsilon, x; \delta) \equiv \partial \lambda(\epsilon, x; \eta) / \partial \eta |_{\eta=\delta}$ .

**Proposition 4** *Under CI, CS', UC, CT, FR, SV and RG', the information for  $\beta_k$  is positive for all  $k$ .*

Proof of Proposition 4 is presented in the appendix. We sketch the heuristics of the idea here. Exploiting properties of  $\mu$  (the measure on  $\{0, 1\} \otimes \Omega_Z$  defined in Section 2), we can show the Fisher information for  $\beta_k$  takes the form of

$$\inf_{\lambda \in \Lambda} 4 \int \phi(z) \left[ f_{\epsilon|x}(w) \left( x_k - \sum_{j \neq k} \alpha_j^* x_j \right) - \alpha_\lambda^* \lambda_\eta(w, x; \delta) \right]^2 dF_Z \quad (15)$$

where  $\phi(z) \equiv [F_{\epsilon|x}(w)(1 - F_{\epsilon|x}(w))]^{-1} \geq 0$ ; and  $(\alpha_j^*)_{j \neq k}$  and  $\alpha_\lambda^*$  constitute a solution to the minimization problem in (5) that defines path-wise information  $I_{\lambda, k}$ . To begin with, note that if  $I_{\lambda, k}$  were to be zero for any  $\lambda \in \Lambda$ , it must be the case that  $\alpha_\lambda^* \neq 0$ .<sup>8</sup>

<sup>8</sup>Otherwise the pathwise information  $I_{\lambda, k}$  under  $\lambda$  would equal that of a parametric model where true error distribution  $F_{\epsilon|X}$  is known, and be positive. This would contradict the claim that  $I_{\lambda, k} = 0$ .

Since each path  $\lambda$  in  $\Lambda$  needs to satisfy conditional symmetry assumption for  $\eta$  close to  $\delta$ ,  $\lambda_\eta(w, x; \delta)$  and consequently its product with  $\alpha_\lambda^* \neq 0$  must be odd functions in  $w$  once  $x$  is fixed. At the same time  $f_{\epsilon|x}(w)$  is an even function of  $w$  (i.e. symmetric in  $w$  around 0) given  $x$ . Then the pathwise information for  $\beta_k$  under  $\lambda$  amounts to a weighted integral of squared distances between an odd and an even function. Thus positive information for  $\beta_k$  is achieved because with the true index  $W$  falls to both sides of zero with positive probabilities, the even function cannot approximate the odd function well enough to reduce  $I_{\lambda,k}$  to zero.

Some discussions relating Proposition 4 to existing literature are in order. Recall the model in (Zheng 1995) where  $\epsilon$  is symmetric around 0 given  $Z = (X, V)$  with all coordinates in  $Z$  correlated with  $\epsilon$ . The information for  $\beta_k$  is zero in that case because the scores  $\psi_k$  and  $\psi_\lambda$  are both flexible in the sense of depending on  $V$  as well as  $X$ . Hence one can construct linear combinations of  $(\psi_j)_{j \neq k}$  and  $\psi_\lambda$  that are arbitrarily good approximation to  $\psi_k$  in  $L^2(\mu)$ -norm, provided for the path  $\lambda$  is appropriately a chosen. To see this, note  $I_{\lambda,k} \equiv$

$$\inf_{\alpha_\lambda, (\alpha_j)_{j \neq k}} \int \left( \psi_k - \alpha_\lambda \psi_\lambda - \sum_{j \neq k} \alpha_j \psi_j \right)^2 d\mu \leq 4 \int \phi(z) [f_{\epsilon|z}(w)x_k - \lambda_\eta(w, z; \delta)]^2 dF_Z. \quad (16)$$

Indeed the same path used for showing zero information for  $\beta_k$  under MI with no excluded regressors (see Theorem 5 in (Chamberlain 1986)) also drives information for  $\beta_k$  to zero in (Zheng 1995) as well. Specifically, the path is  $\lambda(\epsilon, z; \eta) = F_{\epsilon|z}(\epsilon) [1 + (\eta - \delta)h(\epsilon, z)]$ , where  $h$  is continuously differentiable, equals zero outside of some compact set, and  $h(0, z) = 0$  for all  $z$  so that  $\lambda_\eta(\epsilon, z; \delta) = h(\epsilon, z)$ . Since there is no restriction on how the vector  $z$  enters  $\lambda_\eta$ , one can exploit such flexibility to make the approximation on the right-hand side of (16) arbitrarily good and establish zero information in (Zheng 1995)'s case. In contrast, in our model under CI as well as CS, the excluded regressor  $v$  can only enter  $\lambda_\eta(w, x; \delta)$  through the index  $w = x\beta - v$ . This additional form restriction is what delivers the positive information for  $\beta_k$ .

(Magnac and Maurin 2007) considered a binary regression model with CI, mean independence ( $E(\epsilon|X) = 0$ ) and certain tail condition that restricts the truncated expectation of  $F_{\epsilon|X}$  outside of the support of  $V$  given  $X$ .<sup>9</sup> They showed the information for  $\beta_k$  is positive in such a model and derived the semiparametric efficiency bound. The tail condition in (Magnac and Maurin 2007) is a joint restriction on the location of the support of  $V$  as well as tail behaviors outside of the support of  $V$ , while the CS condition we consider in this article is a restriction on the shape of  $F_{\epsilon|X}$  over the full support. These two sets of conditions are non-nested. (See Appendix B for details.)

---

<sup>9</sup>See equation (5) in Proposition 5 of Magnac and Maurin (2007) for the tail restriction. Essentially, this is sufficient and necessary for to extending the proof of identification of  $\beta$  in Lewbel (2000), a model with CI and mean independence, when the support of the excluded regressor  $V$  is bounded between  $v_L > -\infty$  and  $v_H < \infty$ .

### 4.3 Root-N Estimation: Extremum Estimator

Our findings in Proposition 4 suggest root-N regular estimators for  $\beta$  can be constructed. To our knowledge, (Chen 2005) was the first to propose an estimator that attains the parametric rate under the assumptions considered in this section. We now conclude this section with an alternative estimator to (Chen 2005). Our estimator and its asymptotic properties are qualitative different from that in (Chen 2005). We consider the case where all coordinates in  $Z$  are continuously distributed. Extensions to cases with mixed covariates are straightforward and omitted for brevity.

Let  $(\cdot)_- \equiv -\min\{\cdot, 0\}$  and  $(\cdot)_+ \equiv \max\{\cdot, 0\}$ . Our estimator is

$$\hat{\beta} \equiv \arg \min_{b \in \mathcal{B}} \hat{H}_n(b), \quad (17)$$

where:

$$\begin{aligned} \hat{H}_n(b) &\equiv \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) [\kappa(\hat{w}_{i,j} - 1) \varphi^-(Z_i, Z_j; b) + \kappa(1 - \hat{w}_{i,j}) \varphi^+(Z_i, Z_j; b)]; \\ \varphi^-(z_i, z_j; b) &\equiv \left( \frac{(x_i + x_j)'}{2} b - \frac{v_i + v_j}{2} \right)_- \text{ and } \varphi^+(z_i, z_j; b) \equiv \left( \frac{(x_i + x_j)'}{2} b - \frac{v_i + v_j}{2} \right)_+; \\ \hat{w}_{i,j} &\equiv \hat{p}_i + \hat{p}_j; \text{ and } \hat{p}_l \equiv \hat{p}(z_l) \equiv \frac{\sum_{s \neq l} y_s \mathcal{K}_\sigma(z_s - z_l)}{\sum_{s \neq l} \mathcal{K}_\sigma(z_s - z_l)} \text{ for } l = i, j. \end{aligned}$$

where  $K_h(\cdot) \equiv h^{-k} K(\cdot/h^k)$  and  $\mathcal{K}_\sigma(\cdot) \equiv \sigma^{-(k+1)} \mathcal{K}(\cdot/\sigma^{k+1})$ , with  $K, \mathcal{K}$  and  $h_n, \sigma_n$  being kernel functions and bandwidths whose properties are to be specified below. The weighting function  $\kappa$  satisfies the following properties.

**WF** (*Weighting Function*)  $\kappa : \mathbb{R} \rightarrow [0, 1]$  satisfies:  $\kappa(t) = 0$  for all  $t \leq 0$ ;  $\kappa(t) > 0$  for all  $t > 0$ ;  $\kappa$  is increasing over  $[0, +\infty)$  and twice continuously differentiable with bounded derivatives in an open neighborhood around 0.

The weight function, evaluated at  $\hat{w}_{i,j}$ , could be intuitively interpreted as a smooth replacement for the indicator function  $1\{\hat{w}_{i,j} \geq 1\}$ . To derive asymptotic properties of  $\hat{\beta}$ , we first show  $\hat{H}_n$  converges in probability to a limiting function  $H_0$  uniformly over the parameter space, where

$$H_0(b) = \mathbb{E} \left\{ f(X) \mathbb{E} \left[ \kappa(W_{i,j} - 1) \varphi^-(Z_i, Z_j; b) + \kappa(1 - W_{i,j}) \varphi^+(Z_i, Z_j; b) \mid X_j = X, X_i = X \right] \right\} \quad (18)$$

where  $f$  is the true density for non-special regressors  $X$  in the data-generating process; and  $w_{i,j}$  is the sum of true propensity scores  $p(z_i)$  and  $p(z_j)$ . The inner expectation of (18) is taken with respect to  $V_i, V_j$  given  $X_j = X_i = X$  while the outer expectation is taken w.r.t.  $X$  (distributed according to  $f$ ). The next proposition shows  $\beta$  is identified as the unique minimizer of  $H_0$  in  $\mathcal{B}$ .

**Proposition 5** *Suppose CI, CS, UC, CT, SV, FR and WF hold. Then  $H_0(b) > 0$  for all  $b \neq \beta$  and  $H_0(\beta) = 0$ .*



Proof of this proposition follows from arguments similar to that of Proposition 3, and is included in Appendix C. We now collect conditions for our estimator to be consistent.

**PS** (*Parameter Space*)  $\beta$  lies in the interior of the parameter space  $\mathcal{B}$ , which is a compact subset of  $\mathbb{R}^k$ .

**SM1** (*Smoothness*) (i) The density of  $Z = (X, V)$  is bounded away from zero by some positive constant over its compact support. (ii) The density of  $Z$  and the propensity score  $p(Z)$  is  $m_{\mathcal{K}}$ -times continuously differentiable (where  $m_{\mathcal{K}} \geq k + 2$ ); and the derivatives are all Lipschitz continuous. (iii)  $\mathbb{E}\{[Y - p(z)]^2 | z\}$  is continuous in  $z$ . (iv)  $H_0(b)$  is continuous in  $b$  in an open neighborhood around  $\beta$ . (v) For all  $x_i$ ,  $\mathbb{E}[\tilde{\varphi}(Z_i, Z_j; b) | X_i = x_i, X_j = x_j] f(x_j)$  is twice continuously differentiable in  $x_j$  around  $x_j = x_i$ , where

$$\tilde{\varphi}(z_i, z_j; b) \equiv \kappa(w_{i,j} - 1)\varphi^-(z_i, z_j; b) + \kappa(1 - w_{i,j})\varphi^+(z_i, z_j; b).$$

**KF1** (*Kernel Function for Estimating Propensity Scores*) (i)  $\mathcal{K}$  is the product of  $k + 1$  univariate kernel functions (denoted  $\tilde{K}$ ), each of which is symmetric around 0, bounded over a compact support, and integrates to 1. (ii) The order of  $\tilde{K}$  is  $m_{\mathcal{K}}$ . (iii)  $\|t\|^l \tilde{K}(t)$  is Lipschitz continuous for  $0 \leq l \leq m_{\mathcal{K}}$ .

**BW1** (*Bandwidth for Estimating Propensity Scores*)  $\sigma_n$  is proportional to  $n^{-\rho_\sigma}$ , where  $\rho_\sigma \in \left(\frac{1}{2m_{\mathcal{K}}}, \frac{1}{2(k+1)}\right)$ .

**FM1** (*Finiteness*)  $\mathbb{E}\{[\mathcal{C}(X_i, X_j) - (V_i + V_j)/2]^2\}$  and  $\mathbb{E}\{[\mathcal{D}(X_i, X_j) - (V_i + V_j)/2]^2\}$  are finite, where  $\mathcal{C}(X_i, X_j) \equiv \inf_{b \in \mathcal{B}} (X_i + X_j)'b/2$  and  $\mathcal{D}(X_i, X_j) \equiv \sup_{b \in \mathcal{B}} (X_i + X_j)'b/2$ .

**KF2** (*Kernel Functions for Matching*)  $K$  is the product of  $k$  univariate kernel functions (denoted  $\tilde{K}(\cdot)$ ) such that (i)  $\tilde{K}(\cdot)$  is bounded over a compact support, symmetric around 0 and integrates to 1. (ii) The order of  $\tilde{K}(\cdot)$  is  $m_\varphi$ , where  $m_\varphi > 2k$ .

**BW2** (*Bandwidths for Matching*)  $h_n$  is proportional to  $n^{-\rho_h}$  with  $\rho_h \in \left(\frac{1}{4k}, \frac{1}{3k}\right)$ .

**Proposition 6** Suppose conditions for Proposition 5 hold. Under PS, SM1, FM1, KF1,2 and BW1,2,  $\hat{\beta} \xrightarrow{P} \beta$ .

Proof of Proposition 6 amounts to checking conditions for basic consistency theorems for extreme estimators, such as Theorem 4.1 in (Amemiya 1985) and Theorem 2.1 in (Newey and McFadden 1994). A key step of the proof is to show that our objective function  $\hat{H}_n$  converges the limiting function  $H_0$  uniformly over the parameter space. Our approach is to first show the difference between the objective function to an infeasible version, where estimates for propensity scores  $\hat{p}(z)$  are replaced by the truth  $p(z)$ , is negligible in a uniform sense. Since the infeasible objective function takes the form of a

second-order U-process indexed by  $b \in \mathcal{B}$ , it can be decomposed by the H-decomposition into the sum of an unconditional expectation involving the matching kernel; and two degenerate U-processes with orders one and two respectively. We then use known results from (Sherman 1994b) to show the two U-processes converge to 0 uniformly over  $\mathcal{B}$  given our choices of kernels and bandwidths; and show the unconditional expectation is  $H_0(b) + o(1)$  for all  $b$  by a standard approach of changing variables.

The kernel and bandwidth conditions in KF1 and BW1, together with smoothness conditions (i)-(iii) in SM1, ensures the preliminary estimates of propensity scores converge uniformly to the true population propensity score. This is useful for showing that replacing  $\hat{p}$  with the true propensity score only results in negligible differences. The choice of  $\rho_\sigma$  in BW1 ensures that: (a) the order of the part of mean-square error due to bias is dominated by that ascribed to variance (i.e.  $1/\sqrt{n\sigma_n^{k+1}} > \sigma_n^{m\kappa}$ ); (b) the resulted rates of uniform converge of  $\hat{p}$  is faster than  $n^{-1/4}$  (i.e.  $1/\sqrt{n\sigma_n^{k+1}} < n^{-1/4}$ ); and (c) the order of  $\sigma_n^{m\kappa}$  is smaller than  $o(n^{-1/2})$ . The requirements (b) and (c) are sufficient but not necessary for consistency. As is explained later, (b) and (c) help to show a quadratic approximation of the objective function is accurate enough in a uniform sense over certain shrinking neighborhood around the true  $\beta$  to lead to the parametric rate. BW2 is also sufficient but not necessary for consistency. This is because the uniform convergence of U-processes over  $\mathcal{B}$  in the H-decomposition only require  $n^{-1/2}h_n^{-k}$  (and therefore  $n^{-1}h_n^{-k}$ ) to be  $o(1)$ ; and the convergence of the unconditional expectation only requires  $h_n \rightarrow 0$ . Nonetheless, just as with  $\sigma_n$ , the specific range of magnitude for  $h_n$  is needed for showing the quadratic approximation  $H_0$  uniformly over  $\mathcal{B}$  is fast enough to induce the parametric rate. Continuity of  $H_0$  in SM1-(iv) is a necessary condition for applying the consistency theorem for extremum estimators. The other conditions in SM1 are also useful for showing uniform convergence of  $\hat{H}_n$  to  $H_0$ . The finiteness condition in FM1 is instrumental as it is need for applying the results on uniform convergence of degenerate U-processes in (Sherman 1994b).

To establish that  $\hat{\beta}$  attains the parametric rate with normal limiting distribution, we need the following additional restriction on smoothness and finiteness of some population moments. To simplify notations, let  $\Delta\varphi_{i,j}^-(b) \equiv \varphi^-(Z_i, Z_j; b) - \varphi^-(Z_i, Z_j; \beta)$  and likewise define  $\Delta\varphi_{i,j}^+$ . Let  $\kappa_-(W_{i,j}) \equiv \kappa(W_{i,j} - 1)$  and  $\kappa_+(W_{i,j}) \equiv \kappa(1 - W_{i,j})$ ; and let  $\kappa'_-(W_{i,j}) \equiv \kappa'(W_{i,j} - 1)$  and  $\kappa'_+(W_{i,j}) \equiv \kappa'(1 - W_{i,j})$ .

**SM2** (*Smoothness of Population Moments*) (i)  $H_0(b)$  is twice continuously differentiable in an open neighborhood around  $\beta$ . (ii) For all  $x$  and  $x'$ ,  $\bar{\varphi}^-(x, x'; b)$  and  $\bar{\varphi}^+(x, x'; b)$  are twice continuously differentiable in  $b$  in an open neighborhood around  $\beta$  where for superscripts  $\diamond \in \{+, -\}$ ,

$$\bar{\varphi}^\diamond(x, x'; b) \equiv \mathbb{E}[\kappa_\diamond(W_{i,j})\Delta\varphi_{i,j}^\diamond(b)|X_i = x, X_j = x'].$$

For all  $x'$ ,  $\nabla_b\bar{\varphi}^-(x, x'; \beta)f(x)$  and  $\nabla_b\bar{\varphi}^+(x, x'; \beta)f(x)$  are  $m_\varphi$ -times continuously differentiable in  $x$  at  $x = x'$  with bounded derivatives; and  $\nabla_{bb}\bar{\varphi}^-(x, x'; \beta)f(x)$  and  $\nabla_{bb}\bar{\varphi}^+(x, x'; \beta)f(x)$  are both continuously differentiable in  $x$  at  $x = x'$  with bounded

derivatives. (iii) For all  $x, x'$ ,  $\varpi^-(x, x'; b)$  and  $\varpi^+(x, x'; b)$  are continuously differentiable in an open neighborhood around  $\beta$ , where for superscripts  $\diamond \in \{+, -\}$ ,

$$\varpi^\diamond(x, x'; b) \equiv \mathbb{E} [|\Delta\varphi_{i,j}^\diamond(b)| | X_i = x, X_j = x'] .$$

For all  $x'$ ,  $\nabla_b \varpi^-(x, x'; \beta)f(x)$  and  $\nabla_b \varpi^+(x, x'; \beta)f(x)$  are continuously differentiable in  $x$  around  $x = x'$  with bounded derivatives. (iv) For all  $z \equiv (x, v)$  and all  $b$  in an open neighborhood around  $\beta$ ,  $\tilde{\mu}^-(z, x'; b)f(x')$  and  $\tilde{\mu}^+(z, x'; b)f(x')$  are  $m_\varphi$ -times continuously differentiable w.r.t.  $X'$  around  $X' = x$ , where for superscripts  $\diamond \in \{+, -\}$

$$\tilde{\mu}^\diamond(z, x'; b) \equiv E[\kappa'_\diamond(W_{i,j})\Delta\varphi_{i,j}^\diamond(b) | Z_i = z, X_j = x'] .$$

The derivatives are all bounded over support of  $z$ . (v) For all  $z$ ,  $m_-^*(z; b)$  and  $m_+^*(z; b)$  are continuously differentiable in  $b$  around  $\beta$  with bounded derivatives, where for super- and subscripts  $\diamond \in \{+, -\}$

$$m_\diamond^*(z; b) \equiv \nabla w(z)f(x)\tilde{\mu}^\diamond(z, x; b), \text{ with } \nabla w(z) \equiv [1/f(z), -p(z)f(z)/f(z)^2] .$$

Besides,  $\nabla_b m_-^*(z; \beta)f(z)$  and  $\nabla_b m_+^*(z; \beta)f(z)$  are both  $m_\mathcal{K}$ -times continuously differentiable with bounded derivatives in  $z$  over its full support.

**FM2** (Finiteness of Population Moments) (i) There exists an open neighborhood around  $\beta$  in  $B$ , denoted  $N(\beta)$ , such that

$$\int \sup_{b \in \mathcal{N}(\beta)} \|\nabla_{bb} \bar{\varphi}^\diamond(x, x; b)\| f(x) dF(x) < \infty \text{ and } \int \sup_{b \in \mathcal{N}(\beta)} \|\nabla_b \varpi^\diamond(x, x; b)\| f(x) dF(x) < \infty,$$

for superscripts  $\diamond \in \{+, -\}$ . (ii) For subscripts  $\diamond \in \{+, -\}$ ,  $\int \|\nabla_b m_\diamond^*(z; \beta)f(z)\| < \infty$  and there exists  $\tilde{\varepsilon} > 0$  such that

$$E[\sup_{\|\tilde{\varepsilon}\| \geq 0} \|\nabla_b m_\diamond^*(Z + \tilde{\varepsilon}; \beta)f(Z + \tilde{\varepsilon})\|^4] < \infty.$$

(iii) For subscripts  $\diamond \in \{+, -\}$ ,  $\int \|\nabla_b m_\diamond^*(z; \beta)f(z)\| dz < \infty$ .

Let  $Q \equiv (Y, 1)$ . Define  $\delta^* = \delta_-^* + \delta_+^*$ , where for subscripts  $\diamond \in \{+, -\}$ ,

$$\delta_\diamond^*(y, z) \equiv q \nabla_b m_\diamond^*(z; \beta)f(z) - \mathbb{E}[Q \nabla_b m_\diamond^*(Z; \beta)f(Z)]$$

**Proposition 7** Suppose conditions for Proposition 6 hold. Under additional conditions SM2 and FM2,

$$\sqrt{n} \left( \hat{\beta} - \beta \right) \xrightarrow{d} \mathcal{N}(0, \Sigma^{-1} \Omega (\Sigma^{-1})')$$

where

$$\Sigma \equiv \nabla_{bb} H_0(\beta) \text{ and } \Omega \equiv 4\mathbb{E}[\delta^*(Y, Z)\delta^*(Y, Z)'] .$$

The proof follows steps similar to Khan (2001). The continuity of  $H_0$  under SM1-(v) is strengthened to SM2-(i), which helps showing that the limiting function  $H_0$  to have quadratic approximation that is sufficiently precise over an open neighborhood around  $\beta$ .

#### 4.4 Root-N Estimation: Close-Form Estimator

There exists an alternative estimator that has a close form and is easier to compute:

$$\hat{\beta}_{CF} \equiv \left[ \sum_{i,j} K_1 \left( \frac{x_i - x_j}{h_{1,n}} \right) K_2 \left( \frac{\hat{p}_i + \hat{p}_j - 1}{h_{2,n}} \right) (x_i + x_j)' (x_i + x_j) \right]^{-1} \times \left[ \sum_{i,j} K_1 \left( \frac{x_i - x_j}{h_{1,n}} \right) K_2 \left( \frac{\hat{p}_i + \hat{p}_j - 1}{h_{2,n}} \right) (x_i + x_j)' (v_i + v_j) \right] \quad (19)$$

where  $K_1$  is a product kernel;  $K_2$  is a univariate kernel;  $\hat{p}_i, \hat{p}_j$  are kernel estimates of propensity scores as before; and  $h_{1,n}, h_{2,n}$  are sequences of bandwidths. The intuition for this estimator is as follows: Suppose one can collect pairs of observations  $z_i, z_j$  with  $i \neq j$  such that  $x_i = x_j$  and  $p_i + p_j = 1$ . Then CI and CS imply for any such pair,  $v_i + v_j = (x_i + x_j)'\beta$ . The estimator in (19) implements this intuition by using kernel smoothing to collect such matched pairs of  $z_i, z_j$  and then estimates the coefficient by finding a vector that provides the best linear fit of  $v_i + v_j$  as a function of  $x_i + x_j$ .

A couple of further remarks regarding close-form and extremum estimators are in order. The close-form estimator has an obvious computational advantage in that it does not require any optimization routine for finding the minimizer of a non-linear objective function. It does require choosing three bandwidths in the two matching kernels  $K_1, K_2$  and the kernel  $\mathcal{K}$  for estimating  $\hat{p}_i, \hat{p}_j$ , as compared to two bandwidths in the extremum estimator (one in the matching kernel  $K$  and one in  $\mathcal{K}$ ). Nevertheless, we do not expect the additional choice of bandwidth in the univariate kernel  $K_2$  to pose much additional computational problem due to its low dimension. We also conjecture the close-form estimator has a smaller asymptotic variance than the extremum estimator, and leave derivation of its asymptotic properties to future research.

The extremum estimator, on the other hand, has an advantage of being robust to the loss of point identification in the following sense: In case the conditions for point identifying  $\beta$  under CI and CS (e.g. SV and FR) fail, the estimator in (17) provides a consistent estimator for  $\{b : H_0(b) = 0\}$ , which is the identified set for  $\beta$  under CI and CS according to Proposition 3.<sup>10</sup> This is in part due to the fact that, as shown in our proof in Appendix C, the objective function  $\hat{H}_n$  in (17) converges uniformly to  $H_0$  uniformly over the parameter space for coefficients.

Comparison between close-form and extremum estimators in the current model under CI and CS is reminiscent of comparison between Ichimura's two-step estimator in (Ichimura 1993) and the Maximum Score estimator in (Manski 1985). The latter two are both estimators for  $\beta$  in binary regressions under a weaker assumption of MI alone. Ichimura's estimator has a close-form and involves an additional choice of bandwidth

<sup>10</sup>In the context of set estimators, consistency can be defined using the Hausdorff set metric as in (Manski and Tamer 2002).

in the estimation of propensity scores in the first-step, while Manski’s Maximum Score estimator has no close form but does not require any choice of bandwidth.

We conclude this subsection with a technical note on the number of paired observations used in both types of estimators. The close-form estimator uses  $n^2$  pairs of observations in total, including  $n$  pairs with  $i = j$ . For such pairs with  $i = j$ , the sums in the square brackets in (19) are reduced to  $4K_1(0) \sum_i K_2\left(\frac{\hat{p}_i - 1/2}{h_{2,n}}\right) x_i' x_i$  and  $4K_1(0) \sum_i K_2\left(\frac{\hat{p}_i - 1/2}{h_{2,n}}\right) x_i' v_i$  respectively. That is, if we were to use unpaired observations only (as opposed to pairs) in the summand of (19), then the estimator would be reduced to the two-step close-form estimator proposed by (Ichimura 1993) for the case of heteroskedastic binary regressions under MI, which is known to converge at a rate slower than the parametric rate. By the same token, the extremum estimator in the preceding subsection could in principle also be modified to include  $n$  pairs with  $i = j$ . Indeed, if  $H_n$  in (17) were to be defined only using pairs with  $i = j$ , then it would lead to an estimator numerically equivalent to the one proposed under CI and MI in (7), which is known to converge at a rate slower than root- $N$  because of zero information for  $\beta$  under CI and MI (Proposition 2). While the inclusion of “ $i = j$ ” pairs in the definition of extremum and close-form estimators has asymptotically negligible impact on these estimators, we expect it to improve the finite sample performance of both estimators.

## 4.5 Monte Carlo

We now present some simulation evidence for performance of our extremum estimator and the two-step, inverse-density-weighted estimator in (Lewbel 2000). The estimator in (Lewbel 2000) was introduced under CI and mean independence, which is a weaker set of assumptions than CI and CS.

The data-generating process is specified as follows:  $Y = 1\{\alpha + X\beta + V + \epsilon \geq 0\}$  where  $V$  is a scalar variable following the standard normal distribution. Both  $X$  and  $\epsilon$  are scalar variables. The triplet  $(X, V, \epsilon)$  are mutually independent. We experiment with three sets of parametric specifications of marginal distributions for  $(X, \epsilon)$ , where both of them are either (a) standard normal; (b) standard logistic; or (c) standard Laplace. We choose these distributional designs with a view to studying performance of these estimators when the errors have different thickness of tails. Among the three parametric classes, the normal distribution has the thinnest tail while the logistic distribution has the thickest.

The true values for  $\alpha$  and  $\beta$  in the data-generating process (DGP) are set to 0.2 and 0.5 respectively. For each choice of sample sizes ( $N = 50, 100, 200, 400$  and  $800$ ), we simulate 1000 data sets and apply our extremum estimator (labeled as “Pairwise”) and the inverse-density weighted estimator in (Lewbel 2000). The bandwidths in the matching kernel  $K$  and the kernel for estimating propensity scores are chosen according to conditions listed in Section 4.3. We report descriptive statics from sampling distributions of these

estimators out of 1000 simulations. These include bias, standard deviation, square-root of mean-square errors, and median absolute deviations).

Table 1(a):  $X \sim Normal(0, 1)$ ,  $\epsilon \sim Normal(0, 1)$

			$N = 50$	$N = 100$	$N = 200$	$N = 400$	$N = 800$
<i>Bias</i>	Pairwise	$\alpha$	0.0242	0.0235	0.0215	0.0147	0.0156
		$\beta$	0.0142	-0.0012	-0.0140	-0.0107	-0.0130
	Inverse-DW	$\alpha$	-0.0308	-0.0159	-0.0180	-0.0160	-0.0093
		$\beta$	-0.1262	-0.1081	-0.0952	-0.0788	-0.0651
<i>Std</i>	Pairwise	$\alpha$	0.4384	0.2953	0.2146	0.1475	0.1019
		$\beta$	0.5537	0.3705	0.2521	0.1747	0.1238
	Inverse-DW	$\alpha$	0.2095	0.1538	0.1130	0.0810	0.0616
		$\beta$	0.1966	0.1418	0.1041	0.0770	0.0579
<i>RMSE</i>	Pairwise	$\alpha$	0.4389	0.2961	0.2155	0.1482	0.1030
		$\beta$	0.5536	0.3703	0.2524	0.1749	0.1245
	Inverse-DW	$\alpha$	0.2116	0.1545	0.1144	0.0825	0.0623
		$\beta$	0.2335	0.1782	0.1411	0.1101	0.0871
<i>MAD</i>	Pairwise	$\alpha$	0.3037	0.2020	0.1515	0.1024	0.0696
		$\beta$	0.3325	0.2390	0.1702	0.1251	0.0837
	Inverse-DW	$\alpha$	0.1439	0.1038	0.0776	0.0528	0.0438
		$\beta$	0.1625	0.1283	0.1044	0.0858	0.0691

Table 1(b):  $X \sim Laplace(0, 1)$ ,  $\epsilon \sim Laplace(0, 1)$

			$N = 50$	$N = 100$	$N = 200$	$N = 400$	$N = 800$
<i>Bias</i>	Pairwise	$\alpha$	-0.0420	-0.0173	-0.0090	-0.0046	-0.0055
		$\beta$	-0.0371	-0.0570	-0.0554	-0.0421	-0.0320
	Inverse-DW	$\alpha$	-0.0939	-0.0890	-0.0698	-0.0628	-0.0567
		$\beta$	-0.3150	-0.3027	-0.2892	-0.2712	-0.2528
<i>Std</i>	Pairwise	$\alpha$	0.5117	0.3880	0.2572	0.1846	0.1339
		$\beta$	0.4705	0.3306	0.2275	0.1689	0.1171
	Inverse-DW	$\alpha$	0.2392	0.1787	0.1398	0.1083	0.0866
		$\beta$	0.0944	0.0684	0.0507	0.0401	0.0313
<i>RMSE</i>	Pairwise	$\alpha$	0.5132	0.3882	0.2572	0.1846	0.1339
		$\beta$	0.4717	0.3353	0.2340	0.1740	0.1213
	Inverse-DW	$\alpha$	0.2568	0.1996	0.1562	0.1251	0.1035
		$\beta$	0.3289	0.3103	0.2936	0.2742	0.2548
<i>MAD</i>	Pairwise	$\alpha$	0.3438	0.2484	0.1659	0.1234	0.0899
		$\beta$	0.2804	0.2138	0.1569	0.1217	0.0856
	Inverse-DW	$\alpha$	0.1700	0.1414	0.1096	0.0880	0.0726
		$\beta$	0.3192	0.3031	0.2902	0.2721	0.2535

Table 1(c):  $X \sim Logistic(0, 1)$ ,  $\epsilon \sim Logistic(0, 1)$

			$N = 50$	$N = 100$	$N = 200$	$N = 400$	$N = 800$
<i>Bias</i>	Pairwise	$\alpha$	-0.0329	-0.0096	-0.0075	-0.0070	0.0035
		$\beta$	-0.0386	-0.0682	-0.0523	-0.0336	-0.0193
	Inverse-DW	$\alpha$	-0.0855	-0.0716	-0.0644	-0.0602	-0.0504
		$\beta$	-0.2828	-0.2679	-0.2488	-0.2292	-0.2102
<i>Std</i>	Pairwise	$\alpha$	0.5206	0.3674	0.2560	0.1900	0.1358
		$\beta$	0.4859	0.3050	0.2196	0.1555	0.1177
	Inverse-DW	$\alpha$	0.2360	0.1787	0.1419	0.1041	0.0851
		$\beta$	0.1140	0.0828	0.0642	0.0479	0.0378
<i>RMSE</i>	Pairwise	$\alpha$	0.5214	0.3673	0.2560	0.1900	0.1358
		$\beta$	0.4872	0.3123	0.2257	0.1590	0.1192
	Inverse-DW	$\alpha$	0.2509	0.1924	0.1558	0.1202	0.0989
		$\beta$	0.3049	0.2804	0.2570	0.2342	0.2136
<i>MAD</i>	Pairwise	$\alpha$	0.3092	0.2284	0.1658	0.1268	0.0897
		$\beta$	0.2581	0.2138	0.1529	0.1079	0.0800
	Inverse-DW	$\alpha$	0.1671	0.1364	0.1088	0.0816	0.0658
		$\beta$	0.2872	0.2697	0.2508	0.2290	0.2110

In all three designs, both the pairwise extremum estimator and the inverse-density-weighted estimator are shown to converge to the true parameter values as sample sizes increase to the moderate size of  $N = 800$ . The pairwise estimator converges at approximately the root-n rate regardless of parametrization of error distributions. The inverse-density-weighted estimator appears to converge faster under the normal errors than under logistic and Laplace errors. This conforms with earlier observations in (Khan and Tamer 2010) that the performance of the inverse-density-weighted estimator could be sensitive to the thickness of the tails of error distributions relative to that of the special regressor.

Besides, under all three specifications, the inverse-density-weighted estimator seems to outperform the pairwise estimator in terms of RMSE when sample sizes are as small as  $N = 50$ . Nevertheless, it is shown to converge more slowly than the pairwise estimator under all designs. The inverse-density-weighted estimator demonstrates smaller variances than the pairwise estimator uniformly across all designs and sample sizes. On the other hand, the pairwise estimator shows lower bias than the inverse-density-weighted estimator in almost all designs and sample sizes. The figures in the appendix show both our estimator  $(\alpha_1, \beta_1)$  and the inverse-density-weighted estimators  $(\alpha_2, \beta_2)$  appear to be approximately normally distributed in the simulated samples.

## 5 Concluding Remarks

In semiparametric heteroskedastic binary regressions, we study how some regressors, which are additively separable in the latent payoff and independent from errors given all other regressors, contribute to the identifying power of the model and Fisher information for coefficients. We consider two classes of models where identification of coefficients

do not depend on “large support” of the special regressors: one with median independent errors; and one with conditionally symmetric errors.

We find that with median-independent errors, using a special regressor does not directly add to the identifying power or information for coefficients. Nonetheless it does help recover error distributions and average structural functions. In contrast, with conditionally symmetric errors, the presence of a special regressor improves the identifying power by the criterion in (Manski 1988), and the Fisher information for coefficients is strictly positive under mild conditions. We also propose root-n estimators for a binary regressions with heteroskedastic but conditionally symmetric errors. Directions of future investigations could include similar exercises for other limited dependent variable models such as censored or truncated regressions.



## Appendix A: Proofs

**Proof of Proposition 1.** (Sufficiency) Under CI and MI,  $p(x, v) \leq 1/2$  if and only if  $x\beta \leq v$ . Consider  $b \neq \beta$  with  $\Pr\{Z \in Q_b\} > 0$ . Without loss of generality, consider some  $(x, v) \in Q_b$  with  $x\beta \leq v < xb$ . Then for any  $G_{\epsilon|X} \in \Theta$  (where  $\Theta$  here in Section 3.1 is the set of conditional distributions that satisfy CI and MI), we have  $\int 1(\epsilon \leq xb - v) dG_{\epsilon|x} > 1/2$ , which implies  $(x, v) \in \xi(b, G_{\epsilon|X})$ . Therefore,  $\Pr\{Z \in \xi(b, G_{\epsilon|X})\} > 0$  for such a  $b$  and all  $G_{\epsilon|X} \in \Theta$ . Since  $(Z, \tilde{Z})$  is a pair of states drawn independently from the same marginal, this also implies  $\Pr\{(Z, \tilde{Z}) \in \tilde{\xi}(b, G_{\epsilon|X})\} > 0$  for such a  $b$  and all  $G_{\epsilon|X} \in \Theta$ . Thus  $\beta$  is identified relative to  $b$ .

(Necessity) Consider some  $b \neq \beta$ . Suppose  $\Pr\{Z \in Q_b\} = 0$  so that  $\text{sign}(V - X\beta) = \text{sign}(V - Xb)$  with probability one. Construct a  $\tilde{G}_{\epsilon|x}$  so that  $\tilde{G}_{\epsilon|x}(t; b) = \mathbb{E}(Y|x, V = xb - t)$  for all  $t$  on the support of  $V - xb$  given  $x$ . For  $t$  outside the support of  $V - xb$  given  $x$ , define  $\tilde{G}_{\epsilon|x}(t; b)$  arbitrarily subject to the requirement that  $\tilde{G}_{\epsilon|x}(t; b)$  is monotone in  $t$  over the support  $\Omega_{\epsilon|x}$ . By construction,  $\tilde{G}_{\epsilon|x}(xb - v; b) = \mathbb{E}(Y|x, V = v) \equiv p(z)$  for all  $z \equiv (x, v)$ . If  $xb \in \Omega_{V|x}$ , then  $\tilde{G}_{\epsilon|x}(0; b) = 1/2$  by construction. Otherwise (i.e. zero is outside the support of  $V - xb$  given  $x$ ), construct  $\tilde{G}_{\epsilon|x}(\cdot; b)$  outside the support of  $V - xb$  given  $x$  subject to the requirement that  $\tilde{G}_{\epsilon|x}(0; b) = 1/2$ . This can be done, because  $\Pr\{Z \in Q_b\} = 0$  implies that  $p(x, v) \geq 1/2$  for all  $(x, v)$  if and only if  $v - xb \geq 0$  for all  $(x, v)$ . Hence as long as  $\Pr\{Z \in Q_b\} = 0$  there exists  $\tilde{G}_{\epsilon|X} \in \Theta$  satisfying CI and MI such that  $\Pr\{Z \in \xi(b, \tilde{G}_{\epsilon|X})\} = 0$ . Furthermore, with any pair of  $Z$  and  $\tilde{Z}$  that are drawn independently from the same marginal, that  $\Pr\{Z \in Q_b\} = 0$  implies “ $\text{sign}(X\beta - V) = \text{sign}(Xb - V)$  and  $\text{sign}(\tilde{X}\beta - \tilde{V}) = \text{sign}(\tilde{X}b - \tilde{V})$ ” with probability 1. Thus the distribution  $\tilde{G}_{\epsilon|X}$  constructed as above is in  $\Theta$  and also satisfies  $\Pr\{(Z, \tilde{Z}) \in \tilde{\xi}(b, \tilde{G}_{\epsilon|X})\} = 0$ . Thus  $\beta$  is not identified relative to  $b$ . *Q.E.D.*

**Proof of Corollary 1.** The objective function in (6) is non-negative by construction. We show it is positive for all  $b \neq \beta$ , and 0 for  $b = \beta$ . Consider  $b \neq \beta$ . Then  $\Pr(Xb \neq X\beta) = \Pr(Xb > X\beta \text{ or } Xb < X\beta) > 0$  under FR. W.L.O.G. suppose  $\Pr(Xb > X\beta) > 0$ . SV implies for any  $x$  with  $xb > x\beta$ , there exists an interval of  $v$  with  $xb > v \geq x\beta$ . Hence  $\Pr(Xb - V > 0 \geq X\beta - V) = \Pr(p(Z) \leq 1/2 \text{ and } Xb - V > 0) > 0$ . With  $\Pr(X\beta = V) = 0$  (and hence  $\Pr(p(Z) = 1/2) = 0$ ), this implies  $1\{p(Z) \leq 1/2\}(Xb - V)_+ > 0$  with positive probability. Thus the objective function in (6) is positive for  $b \neq \beta$ . On the other hand, CI and MI implies  $p(Z) \geq 1/2$  if and only if  $X\beta - V \geq 0$ , and the objective function in (6) is 0 for  $b = \beta$ . *Q.E.D.*

**Proof of Proposition 3.** Proposition 1 shows  $\beta$  is identified relative to  $b$  under CI and MI whenever (i) holds. It follows immediately that (i) also implies identification of  $\beta$  relative to  $b$  under the stronger assumptions of CI and CS. To see how (ii) is also sufficient

for identification of  $\beta$  relative to  $b$ , define  $\tilde{Q}_{b,S} \equiv \{(z, \tilde{z}) : \tilde{x} = x \text{ and } (v, \tilde{v}) \in R_b(x)\}$ . By construction, for any  $(z, \tilde{z}) \in \tilde{Q}_{b,S} \subseteq \Omega_Z \otimes \Omega_Z$ , either " $x\beta - v < \tilde{v} - x\beta$  and  $xb - v > \tilde{v} - xb$ " or " $x\beta - v > \tilde{v} - x\beta$  and  $xb - v < \tilde{v} - xb$ ". Under CI and CS, this implies for any  $G_{\epsilon|X} \in \Theta_{CS}$  and any  $(z, \tilde{z}) \in \tilde{Q}_{b,S}$ , either

$$"F_{\epsilon|x}(x\beta - v) + F_{\epsilon|\tilde{x}}(\tilde{x}\beta - \tilde{v}) < 1 \text{ and } G_{\epsilon|x}(xb - v) + G_{\epsilon|\tilde{x}}(\tilde{x}b - \tilde{v}) > 1" \quad (20)$$

or

$$"F_{\epsilon|x}(x\beta - v) + F_{\epsilon|\tilde{x}}(\tilde{x}\beta - \tilde{v}) > 1 \text{ and } G_{\epsilon|x}(xb - v) + G_{\epsilon|\tilde{x}}(\tilde{x}b - \tilde{v}) < 1".$$

Thus  $\tilde{Q}_{b,S} \subseteq \tilde{\xi}(b, G_{\epsilon|X})$  for any  $G_{\epsilon|X} \in \Theta_{CS}$ . Next, for any  $\delta > 0$ , define a " $\delta$ -expansion" of  $\tilde{Q}_{b,S}$  as:

$$\tilde{Q}_{b,S}^\delta \equiv \{(z, \tilde{z}) : \tilde{x}_d = x_d \text{ and } (v, \tilde{v}) \in R_b(x) \text{ and } \|\tilde{x}_c - x_c\| \leq \delta\}.$$

Without loss of generality, suppose all  $(z, \tilde{z}) \in \tilde{Q}_{b,S}$  satisfies (20) for all  $G_{\epsilon|X} \in \Theta_{CS}$ . Then UC implies when  $\delta > 0$  is small enough,  $\|\tilde{x}_c - x_c\|^2$  and  $\|(\tilde{x} - x)\beta\|^2$  and  $\|(\tilde{x} - x)b\|^2$  are also small enough so that (20) holds for all  $(z, \tilde{z})$  in  $\tilde{Q}_{b,S}^\delta$  and all  $G_{\epsilon|X} \in \Theta_{CS}$ . Thus with such a small  $\delta$ , we have  $\tilde{Q}_{b,S}^\delta \subseteq \tilde{\xi}(b, G_{\epsilon|X})$  for all  $G_{\epsilon|X} \in \Theta_{CS}$ . Finally, suppose condition (ii) in Proposition 3 holds for some  $b \neq \beta$  and a set  $\omega$  open in  $\Omega_X$ . Then CT implies

$$\int 1\{(v_i, v_j) \in R_b(x)\} dF_{V_i|\tilde{x}}(v_j) dF_{V_i|x}(v_i) > 0 \quad (21)$$

for all  $(x, \tilde{x})$  with  $x \equiv (x_c, x_d) \in \omega$ ,  $\tilde{x}_d = x_d$  and  $\|\tilde{x}_c - x_c\| \leq \tilde{\delta}$  where  $\tilde{\delta} > 0$  is small enough. Apply the law of total probability to integrate out  $(\tilde{X}, X)$  on the left-hand side of (21) then implies  $\Pr\{(Z, \tilde{Z}) \in \tilde{Q}_{b,S}^\delta\} > 0$  for such a small  $\tilde{\delta}$ . Hence for such a  $b \neq \beta$ ,  $\Pr\{(Z_i, Z_j) \in \tilde{\xi}(b, G_{\epsilon|X})\} > 0$  for all  $G_{\epsilon|X} \in \Theta_{CS}$ , and  $\beta$  is identified relative to  $b$ . The necessity of these two conditions for identifying  $\beta$  relative to  $b$  follows from constructive arguments similar to that in the proof of (1), and is hence omitted for brevity. *Q.E.D.*

**Proof of Proposition 4.** Under assumptions CI, CS', UC, CT, SV and FR,  $\beta$  is point identified. With  $\mu$  consisting of counting measure for  $y \in \{0, 1\}$  and probability measure for  $Z$ , we can show path-wise information for  $\beta_k$  under a path  $\lambda \in \Lambda$  (denoted by  $I_{\lambda,k}$ ) takes the form

$$4 \int \left( \psi_k - \alpha_\lambda^* \psi_\lambda - \sum_{j \neq k} \alpha_j^* \psi_j \right)^2 d\mu = 4 \int_{\Omega_Z} \frac{[f_{\epsilon|x}(w)(x_k - \sum_{j \neq k} \alpha_j^* x_j) - \alpha_\lambda^* \lambda_\eta(w, x; \delta)]^2}{F_{\epsilon|x}(w)[1 - F_{\epsilon|x}(w)]} dF_Z \quad (22)$$

where  $(\alpha_j^*)_{j \neq k}$  and  $\alpha_\lambda^*$  are constants that solve the minimization problem in the definition of  $I_{\lambda,k}$  in (5).

We prove the proposition through contradiction. Suppose  $I_{\lambda,k} = 0$  for some  $\lambda \in \Lambda$ . First off, note  $\alpha_\lambda^*$  must be nonzero for such a  $\lambda$ , because otherwise the path-wise

information  $I_{\lambda,k}$  would equal the Fisher information for  $\beta$  in a parametric model where the true error distribution  $F_{\epsilon|X}$  is known, which is positive. This would lead to a contradiction.

Suppose  $I_{\lambda,k} = 0$  for some  $\lambda \in \Lambda$  with  $\alpha_\lambda^* \neq 0$ . Condition SV states the support  $\Omega_{V|x}$  must include  $x\beta$  in its interior for all  $x$ . Thus there exists an open interval  $(-\varepsilon^*, \varepsilon^*)$  such that  $W \equiv X\beta - V$  is continuously distributed with positive densities over  $(-\varepsilon^*, \varepsilon^*)$  given any  $x$ . Note the integrand in (22) is non-negative by construction. Thus the right-hand side of (22) is bounded below by

$$4 \int_{\Omega_X} \int_{-\varepsilon^*}^{\varepsilon^*} \frac{[f_{\epsilon|x}(w)(x_k - \sum_{j \neq k} \alpha_j^* x_j) - \alpha_\lambda^* \lambda_\eta(w, x; \delta)]^2}{F_{\epsilon|x}(w)[1 - F_{\epsilon|x}(w)]} dF_{W,X}.$$

Differentiating both sides of (13) with respect to  $\eta$  at  $\delta$  suggests  $\lambda_\eta(-\varepsilon, x; \delta) = -\lambda_\eta(\varepsilon, x; \delta)$  for all  $x$  and  $\varepsilon$ . This implies  $\alpha_\lambda^* \lambda_\eta(w, x; \delta)$  is an odd function in  $w$  given any  $x$ . On the other hand, conditional symmetry of errors implies that  $f_{\epsilon|x}(w) \left(x_k - \sum_{j \neq k} \alpha_j^* x_j\right)$  is even in  $w$  (i.e. symmetric in  $w$  around 0) given any  $x$ . Due to CS',  $F_{\epsilon|x}(t)^{-1} [1 - F_{\epsilon|x}(t)]^{-1}$  is uniformly bounded between positive constants for all  $t \in (-\varepsilon^*, \varepsilon^*)$  and  $x \in \Omega_X$ . It follows that for any constant  $\varphi > 0$ ,

$$\int_{\Omega_X} \int_{-\varepsilon^*}^{\varepsilon^*} \left[ f_{\epsilon|x}(w) \left(x_k - \sum_{j \neq k} \alpha_j^* x_j\right) - \alpha_\lambda^* \lambda_\eta(w, x; \delta) \right]^2 dF_{W|x}(w) dF_X(x) < \varphi.$$

Thus for any  $\varphi > 0$ , there exist  $\mathcal{I} \subset [0, \varepsilon^*) \otimes \Omega_X$  or  $\mathcal{I} \subset (-\varepsilon^*, 0] \otimes \Omega_X$  with  $\Pr\{(W, X) \in \mathcal{I}\} > 0$  and

$$\left| f_{\epsilon|x}(t) \left(x_k - \sum_{j \neq k} \alpha_j^* x_j\right) - \alpha_\lambda^* \lambda_\eta(t, x; \delta) \right| < \varphi \quad (23)$$

for all  $(t, x) \in \mathcal{I}$ . Without loss of generality, suppose  $\mathcal{I} \subset [0, \varepsilon^*) \otimes \Omega_X$ , and define  $\bar{\omega} \equiv \{x : \exists t \text{ with } (t, x) \in \mathcal{I}\}$ .

The new condition RG' implies  $\Pr\{X_k - \sum_{j \neq k} \alpha_j^* X_j > 0 \neq 0 | X \in \bar{\omega}\} > 0$ . Consider  $\bar{x} \in \bar{\omega}$  with  $a(\bar{x}) \equiv \bar{x}_k - \sum_{j \neq k} \alpha_j^* \bar{x}_j > 0$ . Thus  $f_{\epsilon|\bar{x}}(t) \left(\bar{x}_k - \sum_{j \neq k} \alpha_j^* \bar{x}_j\right)$  is positive and bounded below by  $a(\bar{x})c > 0$  for all  $t$  such that  $(t, \bar{x}) \in \mathcal{I}$ . Pick  $\varphi \leq \frac{a(\bar{x})c}{2}$ . Then (23) implies  $\alpha_\lambda^* \lambda_\eta(t, \bar{x}; \delta) \geq \frac{a(\bar{x})c}{2} > 0$  for all  $t$  with  $(t, \bar{x}) \in \mathcal{I}$ . By symmetry of  $f_{\epsilon|x}$  and oddness of  $\lambda_\eta(t, x; \delta)$  in  $t$  given any  $x$ ,  $\left| f_{\epsilon|\bar{x}}(-t) \left(\bar{x}_k - \sum_{j \neq k} \alpha_j^* \bar{x}_j\right) - \alpha_\lambda^* \lambda_\eta(-t, \bar{x}; \delta) \right| \geq \frac{3}{2} a(\bar{x})c > 0$  for all  $t$  with  $(t, \bar{x}) \in \mathcal{I}$ . A symmetric argument applies to show such a distance is also bounded below by positive constants for any  $\bar{x} \in \bar{\omega}$  with  $a(\bar{x}) < 0$  and any  $t$  such that  $(t, \bar{x}) \in \mathcal{I}$ . Due to SV,  $\Pr\{(W, X) \in \mathcal{I}^-\} > 0$  where  $\mathcal{I}^- \equiv \{(t, x) : (-t, x) \in \mathcal{I}\}$ . Thus  $\left| f_{\epsilon|x}(t) \left(x_k - \sum_{j \neq k} \alpha_j^* x_j\right) - \alpha_\lambda^* \lambda_\eta(t, x; \delta) \right|$  is bounded away from zero by some positive constant over  $\mathcal{I}^-$ . It then follows that

$$\int_{\mathcal{I}^-} \left[ f_{\epsilon|x}(w) \left(x_k - \sum_{j \neq k} \alpha_j^* x_j\right) - \alpha_\lambda^* \lambda_\eta(w, x; \delta) \right]^2 dF_{W,X}$$

is bounded away from zero by some positive constant. This contradicts the claim that  $I_{\lambda,k} = 0$  for  $\lambda \in \Lambda$  where  $\alpha_\lambda^* \neq 0$ . *Q.E.D.*

## Appendix B: CS and Tail Conditions in Magnac and Maurin (2007)

We now give an example of how some  $F_{\epsilon|X}$  that satisfies CS can fail the tail requirements in (Magnac and Maurin 2007). Suppose the distribution of a continuous random variable  $W$  is such that  $\lim_{t \rightarrow -\infty} tF_W(t) = 0$ . Then for any  $c$ ,

$$\mathbb{E}[(W - c)1(W < c)] = \int_{-\infty}^c (s - c) dF_W(s) = 0 - 0 - \int_{-\infty}^c F_W(s) ds$$

and  $\mathbb{E}[(W - c)1(W > c)] = \mathbb{E}(W - c) - \mathbb{E}[(W - c)1(W < c)] = \mu_W - c + \int_{-\infty}^c F_W(w) dw$ . Let  $Y_H \equiv -(X\beta + \epsilon + v_H)$  and  $Y_L \equiv X\beta + \epsilon + v_L$ . Therefore, for any given  $x$ ,

$$\mathbb{E}[Y_H 1(Y_H > 0)|x] = \int_{-\infty}^{-v_H} F_{X\beta + \epsilon|X=x}(s) ds \quad (24)$$

$$\mathbb{E}[Y_L 1(Y_L > 0)|x] = x\beta + v_L + \int_{-\infty}^{-v_L} F_{X\beta + \epsilon|X=x}(s) ds \quad (25)$$

so that the difference of (25) minus (24) is given by

$$x\beta + v_L + \int_{-v_H}^{-v_L} F_{X\beta + \epsilon|X=x}(s) ds. \quad (26)$$

Suppose  $F_{\epsilon|X}$  satisfies CS, then  $F_{X\beta + \epsilon|X=x}$  is symmetric around  $x\beta$  for all  $x$ . If  $x\beta = \frac{-v_L - v_H}{2}$ , then (26) equals

$$v_L - \frac{1}{2}(v_H + v_L) + \frac{1}{2}(v_H - v_L) = 0.$$

If  $x\beta < \frac{-v_L - v_H}{2}$ , then (26) is strictly less than 0. Likewise if  $x\beta > \frac{-v_L - v_H}{2}$ , then (26) is strictly greater than 0. Now suppose  $x\beta < \frac{-v_L - v_H}{2}$  for all  $x$  on the support  $\Omega_X \subseteq \mathbb{R}_{++}^K$ . Then  $\mathbb{E}[X'Y_H 1(Y_H > 0)] < \mathbb{E}[X'Y_L 1(Y_L > 0)]$ , and the tail condition in Proposition 5 of (Magnac and Maurin 2007) does not hold.

## 6 Appendix C: Asymptotic Properties of $\hat{\beta}$

Our proof follows steps similar to those in Sherman (1994b), Khan (2001), Khan and Tamer (2010) and Abrevaya, Hausman and Khan (2010).

### 6.1 Consistency

Define the objective function of an "infeasible" estimator as follows:

$$H_n(z_i, z_j; b) = \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) [\kappa(w_{i,j} - 1) \varphi^-(z_i, z_j; b) + \kappa(1 - w_{i,j}) \varphi^+(z_i, z_j; b)]$$

where  $w_{i,j}$  is the sum of the true propensity scores (i.e.  $w_{i,j} \equiv p_i + p_j$  with  $p_l \equiv p(z_l)$ ).

**Proof of Proposition 5.** Consider any  $b \neq \beta$ . Under FR,  $\Pr(X\beta - Xb \neq 0) > 0$ . Without loss of generality, suppose  $\Pr(X\beta - Xb > 0) > 0$  and let  $\omega \equiv \{x : x\beta > xb\}$ . Then under SV,

$$\int 1\{2x\beta > v_i + v_j > 2xb\} dF_{V_i, V_j|x}(v_i, v_j) > 0$$

for all  $x \in \omega$ . By construction, whenever  $x_i = x_j$ ,  $p(x_i, v_i) + p(x_j, v_j) > 1$  if and only if  $v_i + v_j < 2x_i\beta = 2x_j\beta$ . Thus for all  $x \in \omega$ , properties of  $\kappa$  in WF imply that:

$$\begin{aligned} & \mathbb{E} [\kappa(W_{i,j} - 1)\varphi^-(Z_i, Z_j; b) + \kappa(1 - W_{i,j})\varphi^+(Z_i, Z_j; b) | X_j = X_i = x] \\ & \geq \mathbb{E} [\kappa(W_{i,j} - 1)\varphi^-(Z_i, Z_j; b) | V_i + V_j \leq 2x\beta, X_j = X_i = x] \Pr(V_i + V_j \leq 2x\beta | X_j = X_i = x) > (27) \end{aligned}$$

By construction, the conditional expectation on the left-hand side can never be negative for any  $x$ . Multiply both sides of (27) by  $f(x)$  and then integrate out  $x$  over its full support (including  $\omega$ ) with respect to the distribution of non-special regressors. Thus we get  $H_0(b) > 0$  for all  $b \neq \beta$ . Likewise, if  $b \neq \beta$  and  $\Pr(X\beta < Xb) > 0$ , then for any  $x$  with  $x\beta < xb$ , SV implies

$$\begin{aligned} & \mathbb{E} [\kappa(W_{i,j} - 1)\varphi^-(Z_i, Z_j; b) + \kappa(1 - W_{i,j})\varphi^+(Z_i, Z_j; b) | X_j = X_i = x] \\ & \geq \mathbb{E} [\kappa(1 - W_{i,j})\varphi^+(Z_i, Z_j; b) | V_i + V_j > 2x\beta, X_j = X_i = x] \Pr(V_i + V_j > 2x\beta | X_j = X_i = x) > 0. \end{aligned}$$

Then  $H_0(b) > 0$  for all  $b \neq \beta$  by the same argument as above.

Next, consider  $b = \beta$ . For any  $x$ ,

$$\begin{aligned} H_0(\beta) &= \mathbb{E} \{f(X)\mathbb{E} [\kappa(W_{i,j} - 1)\varphi^-(Z_i, Z_j; \beta) + \kappa(1 - W_{i,j})\varphi^+(Z_i, Z_j; \beta) | X_j = X_i = X]\} \\ &= \mathbb{E} \{f(X)\mathbb{E} [\kappa(W_{i,j} - 1)\varphi^-(Z_i, Z_j; \beta) | W_{i,j} \geq 1, X_j = X_i = X] \Pr(W_{i,j} \geq 0 | X_j = X_i = X)\} \\ &+ \mathbb{E} \{f(X)\mathbb{E} [\kappa(1 - W_{i,j})\varphi^+(Z_i, Z_j; \beta) | W_{i,j} < 1, X_i = X_i = X] \Pr(W_{i,j} < 0 | X_j = X_i = X)\}. \end{aligned} \quad (28)$$

The first conditional expectation on the right-hand side of (28) is 0, because whenever  $x_i = x_j$ , we have  $w_{i,j} \geq 1$  if and only if  $v_i + v_j \leq 2x_i\beta$ . Likewise the second conditional expectation is also 0. Thus  $H_0(\beta) = 0$ . *Q.E.D*

**Proof of Proposition 6.** The first step of the proof is to establish that

$$\sup_{b \in \mathcal{B}} |\hat{H}_n(b) - H_n(b)| = o_p(1). \quad (29)$$

Let  $\varphi_{i,j}^-(b)$  be a shorthand for  $\varphi^-(z_i, z_j; b)$  and likewise for  $\varphi_{i,j}^+(b)$ . Applying the Taylor's expansion around  $w_{i,j}$  and using the boundedness conditions in FM1 and KF2, we have:

$$\begin{aligned} & \sup_{b \in \mathcal{B}} \left| \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \varphi_{i,j}^-(b) [\kappa(\hat{w}_{i,j} - 1) - \kappa(w_{i,j} - 1) - \kappa'(w_{i,j} - 1)(\hat{w}_{i,j} - w_{i,j})] \right| \\ &= \sup_{b \in \mathcal{B}} \left| \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \varphi_{i,j}^-(b) \kappa_1''(\tilde{w}_{i,j} - 1) \|\hat{w}_{i,j} - w_{i,j}\|^2 \right| \\ &\leq a \sup_z \|\hat{p}(z) - p(z)\|^2 \sup_{b \in \mathcal{B}} \left\{ \frac{1}{n(n-1)} \sum_{j \neq i} |K_h(x_i - x_j) \varphi_{i,j}^-(b) \kappa_1''(\tilde{w}_{i,j})| \right\} \end{aligned} \quad (30)$$

where  $\kappa', \kappa''$  are first- and second-order derivatives of  $\kappa$ ;  $\tilde{w}_{i,j}$  is a random variable between  $\hat{w}_{i,j}$  and  $w_{i,j}$ ; and  $a > 0$  is some finite constant. Under KF2-(iii), FM1-(i) and WF, the second term on the right-hand side (i.e. the supreme of the term in the braces) is  $O_p(1)$ . Under SM1 and KF1,  $\sup_z |\hat{p}(z) - p(z)| = O_p\left(\frac{(\log n)}{\sqrt{n\sigma_n^{k+1}}} + (k+1)\sigma_n^{m\kappa}\right)$  almost surely by Theorem 2.6 of Li and Racine (2007). Our choice of bandwidth in BW1 implies this term is  $o_p(n^{-1/4})$ . Hence the remainder term of the approximation (l.h.s. of (30)) is  $o_p(1)$ . Next, note:

$$\begin{aligned} & \sup_{b \in \mathcal{B}} \left| \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \varphi_{i,j}^-(b) \kappa'(w_{i,j} - 1) (\hat{w}_{i,j} - w_{i,j}) \right| \\ & \leq 2 \sup_z \|\hat{p}(z) - p(z)\| \sup_{b \in \mathcal{B}} \left\{ \frac{1}{n(n-1)} \sum_{j \neq i} |K_h(x_i - x_j) \varphi_{i,j}^-(b) \kappa'(w_{i,j} - 1)| \right\}. \end{aligned}$$

By similar arguments, the second term is bounded in probability, and the first term is  $o_p(n^{-1/4})$ . Thus (29) holds.

Next, decompose  $H_n(z_i, z_j; b)$  as

$$H_n(z_i, z_j; b) = \mathbb{E}[g_n(Z_i, Z_j; b)] + \frac{2}{n} \sum_{i \leq n} g_{n,1}(z_i; b) + \frac{2}{n(n-1)} \sum_{j \neq i} g_{n,2}(z_i, z_j; b) \quad (31)$$

where

$$\begin{aligned} g_n(z_i, z_j; b) & \equiv K_h(x_i - x_j) [\kappa(w_{i,j} - 1) \varphi_{i,j}^-(b) + \kappa(1 - w_{i,j}) \varphi_{i,j}^+(b)]; \\ g_{n,1}(z_i; b) & \equiv \mathbb{E}[g_n(Z, Z'; b) | Z = z_i] + \mathbb{E}[g_n(Z, Z'; b) | Z' = z_i] - 2\mathbb{E}[g_n(Z, Z'; b)]; \text{ and} \\ g_{n,2}(z_i, z_j; b) & \equiv g_n(z_i, z_j; b) - \mathbb{E}[g_n(Z, Z'; b) | Z = z_i] - \mathbb{E}[g_n(Z, Z'; b) | Z' = z_j] + \mathbb{E}[g_n(Z, Z'; b)]. \end{aligned}$$

By construction,  $\mathbb{E}[g_{n,1}(Z_i; b)] = 0$  and  $\mathbb{E}[g_{n,2}(Z_i, Z_j; b) | Z_i = z_i] = \mathbb{E}[g_{n,2}(Z_i, Z_j; b) | Z_j = z_j] = 0$  for all  $z_i, z_j$ .

We now show the second and third term in (31) are  $o_p(1)$  under our conditions. Under KF2 and PS, we get

$$\sup_{n, b \in \mathcal{B}} |h_n^k g_n(z_i, z_j; b)| \leq \mathcal{F}(z_i, z_j) \equiv a' \left[ \kappa(w_{i,j} - 1) \left( \mathcal{C}(x_i, x_j) - \frac{v_i + v_j}{2} \right)_- + \kappa(1 - w_{i,j}) \left( \mathcal{D}(x_i, x_j) - \frac{v_i + v_j}{2} \right)_+ \right]$$

for all  $(z_i, z_j)$ , where  $\mathcal{C}(\cdot)$  and  $\mathcal{D}(\cdot)$  are defined in FM1 and  $a' > 0$  is some finite constant. By arguments as in (Pakes and Pollard 1989), the class of functions:

$$\{h_n^k g_n(z_i, z_j; b) : b \in \mathcal{B}\}$$

is Euclidean with a constant envelop  $\mathcal{F}$ , which satisfies  $\mathbb{E}[\mathcal{F}(Z_i, Z_j)^2] < \infty$  under KF2 and FM1. Besides,  $\mathbb{E}[\sup_{b \in \mathcal{B}} h_n^{2k} g_n(Z_i, Z_j; b)^2] = O(1)$  under KF2 and FM1. It then follows from Theorem 3 in (Sherman 1994b) that the second and the third terms in the decomposition in (31) are  $O_p(n^{-1/2} h_n^{-k})$  and  $O_p(n^{-1} h_n^{-k})$  uniformly over  $b \in \mathcal{B}$  respectively. Under our choice of bandwidth in BW2, these two terms are both  $o_p(1)$ .

Next, we deal with the first term in the H-decomposition above. Let  $\kappa^-(z_i, z_j) \equiv \kappa(w_{i,j} - 1)$  and  $\kappa^+(z_i, z_j) \equiv \kappa(1 - w_{i,j})$  and

$$\tilde{\varphi}(z_i, z_j; b) \equiv \kappa(w_{i,j} - 1) \varphi_{i,j}^-(b) + \kappa(1 - w_{i,j}) \varphi_{i,j}^+(b)$$

to facilitate derivations. By definition,

$$\begin{aligned}
\mathbb{E}[g_n(Z_i, Z_j; b)] &= \int K_h(x_i - x_j) \tilde{\varphi}(z_i, z_j; b) dF(z_i, z_j) \\
&= \int K_h(x_i - x_j) \mathbb{E}[\tilde{\varphi}(Z_i, Z_j; b) | x_i, x_j] dF(x_i, x_j) \\
&= \int K(u) \mathbb{E}[\tilde{\varphi}(Z_i, Z_j; b) | X_i = x_i, X_j = x_i + h_n^k u] f(x_i + h_n^k u) du dF(x_i)
\end{aligned}$$

Changing variables between  $x_j$  and  $u \equiv (x_j - x_i)/h_n^k$  and applying the dominated convergence theorem, we can show that  $\mathbb{E}[g_n(Z_i, Z_j; b)] = H_0(b) + O(kh_n^2) = H_0(b) + o(1)$  for all  $b \in \mathcal{B}$ . Thus the sum of the three terms on the right-hand side of (31) is  $o_p(1)$  uniformly over  $b \in \mathcal{B}$ .

Combine this result with (29), we get:

$$\sup_{b \in \mathcal{B}} |\hat{H}_n(b) - H_0(b)| = o_p(1). \quad (32)$$

The limiting function  $H_0(b)$  is continuous under SM1 in an open neighborhood around  $\beta$ . Besides, Proposition 5 has established that  $H_0(b)$  is uniquely minimized at  $\beta$ . It then follows from Theorem 2.1 in (Newey and McFadden 1994) that  $\hat{\beta} \xrightarrow{p} \beta$ . *Q.E.D.*

## 6.2 Root-N and Asymptotic Normality

For convenience of proof in this section, define:

$$\hat{\mathcal{H}}_n(b) = \hat{H}_n(b) - \hat{H}_n(\beta) \text{ and } \mathcal{H}_n(b) = H_n(b) - H_n(\beta).$$

By construction, the optimizers of  $\hat{\mathcal{H}}_n$  and  $\mathcal{H}_n$  are the same as those for  $\hat{H}_n$  and  $H_n$ .

Having shown consistency, our strategy for deriving the limiting distribution of  $\hat{\beta}$  is to approximate  $\hat{\mathcal{H}}_n(\cdot)$  locally in a neighborhood of  $\beta$  by some function that is quadratic in  $b$ . The approximation needs to accommodate the fact that the objective function is not smooth in  $b$ . Quadratic approximation of such objective functions have been provided in, for example, (Pakes and Pollard 1989), and (Sherman 1994a), (Sherman 1994b) among others. A preliminary step is to show  $\|\hat{\beta} - \beta\|$  converges at a rate no slower than  $\sqrt{n}$ . Once established, this result allows us to focus on such a shrinking neighborhood around  $\beta$  where quadratic approximation mentioned above becomes more precise so that root-n consistency and asymptotic normality can be established in one step. A useful theorem that will be invoked for showing these results is Theorem 1 in (Sherman 1994b), which require the following conditions:

1.  $\hat{\beta} - \beta = O_p(\delta_n)$ ;

2. There exists a neighborhood of  $\beta$  and a constant  $\tilde{a} > 0$  such that  $H_0(b) - H_0(\beta) \geq \tilde{a}\|b - \beta\|^2$  for all  $b$  in this neighborhood of  $\beta$ ; and
3. Uniformly over an  $O_p(\delta_n)$  neighborhood of  $\beta$ :

$$\widehat{\mathcal{H}}_n(b) = H_0(b) + O_p(\|b - \beta\|/\sqrt{n}) + o_p(\|b - \beta\|^2) + O_p(\epsilon_n). \quad (33)$$

Under these three conditions, Theorem 1 in (Sherman 1994b) states  $\hat{\beta} - \beta_0 = O_p(\max\{\sqrt{\epsilon_n}, 1/\sqrt{n}\})$ .

**Lemma C1.** *Under SM2-(i), there exists an open neighborhood of  $\beta$  and some constant  $\tilde{a} > 0$  such that  $H_0(b) - H_0(\beta) \geq \tilde{a}\|b - \beta\|^2$  for all  $b$  in this neighborhood of  $\beta$ .*

**Proof of Lemma C1.** Under SM-(i), we can apply the Taylor's expansion to write:

$$H_0(b) = \frac{1}{2}(b - \beta)' \nabla_{bb} H_0(\tilde{b})(b - \beta)$$

where  $\tilde{b}$  is on the line segment linking  $b$  and  $\beta$ . Note we have used  $H_0(\beta) = 0$  and  $\nabla_b H_0(\beta) = 0$  due to the identification result in Proposition 5. The claim in this lemma then follows from the positive definiteness of  $\nabla_{bb} H_0(\beta)$  and its continuity at  $\beta$ . *Q.E.D.*

To simplify notations in what follows, we let

$$\widehat{\mathcal{H}}_{1,n}(b) \equiv \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \kappa(\hat{w}_{i,j} - 1) [\varphi_{i,j}^-(b) - \varphi_{i,j}^-(\beta)]$$

denote the first half of the "location-normalized" objective function  $\widehat{\mathcal{H}}_n$  (which only involve  $(\cdot)_-$ ); and likewise let  $\mathcal{H}_{1,n}(b)$  and  $H_{1,0}(b)$  denote the first halves of  $\mathcal{H}_n$  and  $H_0$  respectively. Similarly, define  $\widehat{\mathcal{H}}_{2,n}$ ,  $\mathcal{H}_{2,n}$  and  $H_{2,0}$  as the second halves involving  $(\cdot)_+$ . Recall that  $H_{1,0}(\beta) = H_{2,0}(\beta) = 0$  by construction.

**Lemma C2.** *Suppose conditions for Proposition 6 hold. Under additional conditions SM2 and FM2,*

$$\mathcal{H}_n(b) - H_0(b) = o_p(\|b - \beta\|^2) + o_p(\|b - \beta\|/\sqrt{n}) + O_p(n^{-1}h^{-k})$$

*uniformly over an  $o_p(1)$  neighborhood around  $\beta$  in  $\mathcal{B}$ ; and the  $O_p(n^{-1}h^{-k})$  term is further reduced to  $o_p(n^{-1})$  uniformly over an  $O_p(1/\sqrt{nh^k})$  neighborhood of  $\beta$ .*

**Proof of Lemma C2.** We analyze the order of magnitude of  $\mathcal{H}_{1,n} - H_{1,0}$  in this proof. The case with  $\mathcal{H}_{2,n} - H_{2,0}$  follows from the same arguments and is omitted for brevity. For any  $b$ , decomposed  $\mathcal{H}_{1,n}(b) - H_{1,0}(b)$  as:

$$\{\mathbb{E}[\tilde{g}_n(Z_i, Z_j; b)] - H_{1,0}(b)\} + \frac{1}{n} \sum_{i=1}^n \tilde{g}_{n,1}(z_i; b) + \frac{1}{n(n-1)} \sum_{j \neq i} \tilde{g}_{n,2}(z_i, z_j; b) \quad (34)$$



where

$$\begin{aligned}\tilde{g}_n(z_i, z_j; b) &\equiv K_h(x_i - x_j) \kappa^-(z_i, z_j) [\varphi^-(z_i, z_j; b) - \varphi^-(z_i, z_j; \beta)]; \\ \tilde{g}_{n,1}(z_i; b) &\equiv \mathbb{E}[\tilde{g}_n(Z, Z'; b)|Z = z_i] + \mathbb{E}[\tilde{g}_n(Z, Z'; b)|Z' = z_i] - 2\mathbb{E}[\tilde{g}_n(Z, Z'; b)]; \text{ and} \\ \tilde{g}_{n,2}(z_i, z_j; b) &\equiv \tilde{g}_n(z_i, z_j; b) - \mathbb{E}[\tilde{g}_n(Z, Z'; b)|Z = z_i] - \mathbb{E}[\tilde{g}_n(Z, Z'; b)|Z' = z_j] + \mathbb{E}[\tilde{g}_n(Z, Z'; b)].\end{aligned}$$

where  $\kappa^-(z_i, z_j)$  is a shorthand for  $\kappa(w_{i,j} - 1)$ .

We first deal with the first term in (34). With a slight abuse of notation, let  $F$  denote distributions and  $f$  denote densities. Let  $\Delta\varphi^-(z_i, z_j; b) \equiv \varphi^-(z_i, z_j; b) - \varphi^-(z_i, z_j; \beta)$ . Note by the Law of Iterated Expectation, we can write for all  $b$ :

$$\mathbb{E}[\tilde{g}_n(Z_i, Z_j; b)] = \int K_h(x_i - x_j) \bar{\varphi}(x_i, x_j; b) dF(x_i, x_j)$$

where

$$\bar{\varphi}(x, x'; b) \equiv \mathbb{E}\{\kappa^-(Z_i, Z_j) \Delta\varphi^-(Z_i, Z_j; b) | X_i = x, X_j = x'\}.$$

By construction,  $\bar{\varphi}(x_i, x_j; \beta) = 0$  and under SM2-(ii),

$$\bar{\varphi}(x_i, x_j; b) = \nabla_b \bar{\varphi}(x_i, x_j; \beta)(b - \beta) + \frac{1}{2}(b - \beta)' \nabla_{bb} \bar{\varphi}(x_i, x_j; \beta)(b - \beta) + o(\|b - \beta\|^2)$$

for all  $b$  in an  $o(1)$  neighborhood around  $\beta$ ; where  $\nabla_b \bar{\varphi}$  and  $\nabla_{bb} \bar{\varphi}$  are gradient and Hessian w.r.t.  $b$  respectively. Since the magnitude of the remainder is invariant in  $x_i, x_j$ , we can decompose  $\mathbb{E}[\tilde{g}_n(Z_i, Z_j; b)]$  as

$$\begin{aligned}(b - \beta)' \left[ \int K_h(x_i - x_j) \frac{1}{2} \nabla_{bb} \bar{\varphi}(x_i, x_j; \beta) dF(x_i, x_j) \right] (b - \beta) \\ + \left\{ \int K_h(x_i - x_j) \nabla_b \bar{\varphi}(x_i, x_j; \beta) dF(x_i, x_j) \right\} (b - \beta) + o(\|b - \beta\|^2).\end{aligned}\tag{35}$$

for all  $b$  in an  $o(1)$  neighborhood around  $\beta$ . Changing variables between  $x_i$  and  $u \equiv (x_i - x_j)/h_n^k$ , we can write the square bracket term in (35) as

$$\begin{aligned}&\int \left[ \int K(u) \frac{1}{2} \nabla_{bb} \bar{\varphi}(x_j + h^k u, x_j; \beta) f(x_j + h_n^k u) du \right] dF(x_j) \\ &= \frac{1}{2} \int [\nabla_{bb} \bar{\varphi}(x, x; \beta) f(x) + O(h_n^2)] dF(x) = \frac{1}{2} \int [\nabla_{bb} \bar{\varphi}(x, x; \beta) f(x)] dF(x) + o(1) = \frac{1}{2} \nabla_{bb} H_{1,0}(\beta) + o(1).\end{aligned}$$

The first equality above follows from a Taylor expansion of  $\nabla_{bb} \bar{\varphi}(x, x_j; \beta) f(x)$  around  $x = x_j$  under SM2; the order  $K$  in KF2; and the fact that  $\int K(u) du = 1$ . The second equality is due to the facts that the expansion applies for all  $x_j$ ; that the order of remainder is invariant in  $x_j$ ; and that  $O(h_n^2)$  is  $o(1)$  under BW2. The third equality follows from the fact that the order of differentiation and integration can be exchanged under SM2-(ii) and FM2-(i). Similarly we can show the term in the braces in (35) is

$$\begin{aligned}&\int \left[ \int K(u) \nabla_b \bar{\varphi}(x_j + h^k u, x_j; \beta) f(x_j + h_n^k u) du \right] dF(x_j) \\ &= \int [\nabla_b \bar{\varphi}(x, x; \beta) f(x) + O(h_n^{m_\varphi})] dF(x) = \int [\nabla_b \bar{\varphi}(x, x; \beta) f(x)] dF(x) + o(n^{-1/2})\end{aligned}$$

where  $\int [\nabla_b \bar{\varphi}(x, x; \beta) f(x)] dF(x) = \nabla_b H_{1,0}(\beta) = 0$  because the order of integration and differentiation can be exchanged and  $H_{1,0}(b)$  is uniquely minimized at  $b = \beta$ . To sum up, (35) is

$$\frac{1}{2}(b - \beta)' \nabla_{bb} H_{1,0}(\beta) (b - \beta) + o_p(\|b - \beta\| / \sqrt{n}) + o_p(\|b - \beta\|^2).$$

By standard Taylor expansion using SM-(i),  $H_{1,0}(b) = \frac{1}{2}(b - \beta)' \nabla_{bb} H_{1,0}(\beta) (b - \beta) + o_p(\|b - \beta\|^2)$  over an  $o(1)$  neighborhood of  $\beta$ , it then follows that the first term in (34) is  $o_p(\|b - \beta\| / \sqrt{n}) + o_p(\|b - \beta\|^2)$ .

Next, we turn to the second term in (34). By SM2-(ii), we can apply the Taylor expansion around  $\beta$  to the second term in (34) to get

$$\frac{1}{n} \sum_{i=1}^n \tilde{g}_{n,1}(z_i; b) = \frac{1}{n} \sum_{i=1}^n \tilde{g}_{n,1}(z_i; \beta) + \left( \frac{1}{n} \sum_{i=1}^n \nabla_b \tilde{g}_{n,1}(z_i; \tilde{b}) \right) (b - \beta)$$

where  $\tilde{b}$  is on the line segment between  $b$  and  $\beta$  (and possibly depends on  $z_i$ ). By construction,  $\tilde{g}_{n,1}(z_i; \beta) = 0$  for all  $n$  and  $z_i$ . Besides, for any given  $n$  and  $b$ ,  $\nabla_b \tilde{g}_{n,1}(z_i; b)$  has mean zero for all  $z_i$ . To see this, note for any fixed  $n$  and  $b$ :

$$\mathbb{E}[\nabla_b \tilde{g}_{n,1}(Z_i; b)] = \mathbb{E}\{\nabla_b \mathbb{E}[\tilde{g}_n(Z, Z'; b) | Z = Z_i]\} + \mathbb{E}\{\nabla_b \mathbb{E}[\tilde{g}_n(Z, Z'; b) | Z' = Z_i]\} - 2\nabla_b \mathbb{E}[\tilde{g}_n(Z, Z'; b)] = 0$$

because under FM1,2 and SM1,2 the order of integration and differentiation in the first two terms on the right-hand side can be exchanged. Also note that by definition,

$$\nabla_b \mathbb{E}[\tilde{g}_n(Z, Z'; b) | Z = z] = \int K_h(x - x') \nabla_b \hat{\varphi}(z, x'; b) f(x') dx' = \int K(u) \nabla_b \hat{\varphi}(z, x - h^k u; b) f(x - h^k u) du$$

where  $\hat{\varphi}^-(z, x'; b) \equiv \mathbb{E}[\kappa^-(Z_i, Z_j) \Delta \varphi_{i,j}^-(b) | Z_i = z, X_j = x']$ . The first equality follows from an interchange of integration and differentiation; and the second equality follows from a change of variables between  $x'$  and  $u \equiv (x - x')/h^k$ . Likewise we can derive a similar expression for  $\nabla_b \mathbb{E}[\tilde{g}_n(Z, Z'; b) | Z' = z]$ . It then follows from boundedness of  $K$  in KF and the finite moment condition in FM2 that  $\mathbb{E}[\nabla_b \tilde{g}_{n,1}(Z; b) \nabla_b \tilde{g}_{n,1}(Z_i; b)'] = O(1)$  for all  $b$  in an open neighborhood around  $\beta$ . Thus for any fixed  $b$  in an open neighborhood around  $\beta$ , the Liapunov Central Limit Theorem applies and  $\frac{1}{n} \sum_{i=1}^n \nabla_b \tilde{g}_{n,1}(z_i; b) = O_p(n^{-1/2})$  under FM1,2. With  $\tilde{b}$  between  $b$  and  $\beta$ , and with  $b \xrightarrow{p} \beta$ , an application of Lemma 2.17 in (Pakes and Pollard 1989) shows  $\frac{1}{n} \sum_{i=1}^n \nabla_b \tilde{g}_{n,1}(z_i; \tilde{b})$  is  $o_p(n^{-1/2})$  uniformly over an  $o_p(1)$  neighborhood around  $\beta$ . Thus the second term in the decomposition in (34) is  $o_p(\|b - \beta\| / \sqrt{n})$  uniformly over an  $o_p(1)$  neighborhood of  $\beta$ .

Next, arguments similar to Proposition 6 suggest conditions for Theorem 3 in (Sherman 1994b) hold for the third term in the decomposition in (34) multiplied with  $h^k$ , which is a second-order degenerate U-process. Hence the third term is  $O_p(n^{-1} h^{-k})$  uniformly over  $o_p(1)$  neighborhood around  $\beta$ . Furthermore, this term is reduced to  $O_p(n^{-3/2} h^{-3k/2})$  uniformly over an  $O_p(1/\sqrt{nh^k})$  neighborhood around  $\beta$ , which is  $o_p(n^{-1})$  due to our choice of bandwidth in BW2. *Q.E.D.*

Next, we show the difference between  $\widehat{\mathcal{H}}_{1,n}(b)$  and  $\mathcal{H}_{1,n}(b)$  can be expressed in terms of a simple sample average plus some negligible approximation errors over a shrinking neighborhood of  $\beta$  in  $\mathcal{B}$ .

**Lemma C3.** Suppose conditions for Proposition 6 hold. Under additional conditions in SM2 and FM2,

$$|\widehat{\mathcal{H}}_{1,n}(b) - \mathcal{H}_{1,n}(b)| = \frac{2}{n} \sum_{i=1}^n \delta_1^*(y_i, z_i)(b - \beta) + o_p(\|b - \beta\|/\sqrt{n}) + o_p(\|b - \beta\|^2) + O_p(n^{-1}h^{-k}) \quad (36)$$

uniformly over an  $o_p(1)$  neighborhood of  $\beta$  in  $\mathcal{B}$ ; where

$$\delta_1^*(y, z) \equiv q \nabla_b m_-^*(z; \beta) f(z) - \mathbb{E}[Q \nabla_b m_-^*(Z; \beta) f(Z)] \text{ with } q \equiv (y, 1)';$$

Besides, the  $O_p(n^{-1}h^{-k})$  term in (36) is further reduced to  $o_p(n^{-1})$  uniformly over an  $O_p(1/\sqrt{nh^k})$  neighborhood around  $\beta$ .

**Proof of Lemma C3.** Let  $\Delta \varphi_{i,j}^-(b) \equiv \Delta \varphi^-(z_i, z_j; b) \equiv \varphi_{i,j}^-(b) - \varphi_{i,j}^-(\beta)$ . By smoothness of  $\kappa$  in WF, we can use the Taylor's expansion to decompose  $\widehat{\mathcal{H}}_{1,n}(b) - \mathcal{H}_{1,n}(b)$  into:

$$\Delta_{1,n} \equiv \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \kappa'(w_{i,j} - 1)(\hat{w}_{i,j} - w_{i,j})$$

and

$$R_{1,n} \equiv \frac{1}{n(n-1)} \sum_{j \neq i} K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \kappa''(\tilde{w}_{i,j} - 1)(\hat{w}_{i,j} - w_{i,j})^2$$

where  $\tilde{w}_{i,j}$  is between  $\hat{w}_{i,j}$  and  $w_{i,j}$ . Using the triangular inequality and by the fact that the second-order derivative  $\kappa''$  is bounded, we have:

$$|R_{1,n}| \leq \hat{a} \left\{ \frac{1}{n(n-1)} \sum_{j \neq i} |K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b)| \right\} \sup_{z, z'} \|\hat{p}(z) + \hat{p}(z') - p(z) - p(z')\|^2 \quad (37)$$

for some finite constant  $\hat{a} > 0$ . The second term on the right-hand side of (37) is  $o_p(n^{-1/2})$  since under our conditions of SM1, KF1 and BW1,

$$\sup_z |\hat{p}(z) - p(z)| = o_p(n^{-1/4}). \quad (38)$$

As for the first term in the braces of (37), we use the H-decomposition to break it down into the sum of an unconditional expectation and two degenerate U-processes:

$$\mathbb{E}[\varpi_n(Z_i, Z_j; b)] + \frac{1}{n} \sum_{i=1}^n \varpi_{n,1}(z_i; b) + \frac{1}{n(n-1)} \sum_{j \neq i} \varpi_{n,2}(z_i, z_j; b) \quad (39)$$

where  $\varpi_n(z_i, z_j; b) \equiv |K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b)|$ ; and

$$\begin{aligned} \varpi_{n,1}(z_i; b) &\equiv \mathbb{E}[\varpi_n(Z, Z'; b) | Z = z_i] + \mathbb{E}[\varpi_n(Z, Z'; b) | Z' = z_i] - 2\mathbb{E}[\varpi_n(Z, Z'; b)]; \text{ and} \\ \varpi_{n,2}(z_i, z_j; b) &\equiv \varpi_n(z_i, z_j; b) - \mathbb{E}[\varpi_n(Z, Z'; b) | Z = z_i] - \mathbb{E}[\varpi_n(Z, Z'; b) | Z' = z_j] + \mathbb{E}[\varpi_n(Z, Z'; b)]. \end{aligned}$$

By standard arguments in (Pakes and Pollard 1989), the class of functions  $\{h^k \varpi_n(z_i, z_j; b) : b \in \mathcal{B}\}$  is Euclidean with a constant envelop that has finite second moments. Our conditions in FM1, the boundedness of  $K$  over its compact support in KF2, and boundedness of derivatives in WF all imply that both conditions for Theorem 3 of (Sherman 1994b) hold with  $\delta_n$  and  $\gamma_n$  therein being  $o(1)$  and  $O(1)$  respectively. Hence

$$\frac{1}{n} \sum_{i=1}^n \varpi_{n,1}(z_i; b) = O_p(n^{-1/2}h^{-k}); \text{ and } \frac{1}{n(n-1)} \sum_{j \neq i} \varpi_{n,2}(z_i, z_j; b) = O_p(n^{-1}h^{-k}) \quad (40)$$

uniformly over an  $o_p(1)$  neighborhood around  $\beta$ .

As for the unconditional expectation in (39), by definition, it equals

$$\mathbb{E}[\varpi_n(Z_i, Z_j; b)] = \int |K_h(x_i - x_j)| \varpi(x_i, x_j; b) dF(x_i, x_j) \quad (41)$$

where  $\varpi(x, x'; b) \equiv \mathbb{E}\{|\Delta\varphi_{i,j}^-(b)| \mid X_i = x, X_j = x'\}$ . By construction,  $\varpi(x_i, x_j; \beta) = 0$  and under SM2,

$$\varpi(x_i, x_j; b) = \nabla_b \varpi(x_i, x_j; \beta)(b - \beta) + o(\|b - \beta\|) \quad (42)$$

for all  $b$  in an  $o(1)$  neighborhood around  $\beta$ ; where  $\nabla_b \varpi$  is a gradient w.r.t.  $b$ . Change variables between  $x_i$  and  $u \equiv (x_i - x_j)/h^k$  given any  $x_j$  on the right-hand side of (41), and we get

$$\begin{aligned} & \int |K_h(x_i - x_j)| \nabla_b \varpi(x_i, x_j; \beta) dF(x_i, x_j) \\ &= \int \left[ \int |K(u)| \nabla_b \varpi(x_j + h^k u, x_j; \beta) f(x_j + h_n^k u) du \right] dF(x_j) \\ &= \tilde{\kappa}_1 \nabla_b \mathbb{E}[f(X) \varpi(X, X; \beta)] + o(1) \end{aligned} \quad (43)$$

where  $\tilde{\kappa}_1 \equiv \int |K(u)| du$  is finite under KF2. The second equality follows from an application of a first-order Taylor expansion of  $\nabla_b \varpi(x_i, x_j; \beta) f(x_i)$  around  $x_i = x_j$ ; and from changing the order of integration and differentiation allowed under SM2 and FM2. Note that

$$\nabla_b \mathbb{E}[f(X) \varpi(X, X; \beta)] = 0 \quad (44)$$

because  $\mathbb{E}[f(X) \varpi(X, X; \beta)]$  is minimized to 0 at  $b = \beta$ . Hence combining results from (41), (42), (43) and (44), we have:

$$\mathbb{E}[\varpi_n(Z_i, Z_j; b)] = o(\|b - \beta\|). \quad (45)$$

Combining results from (38), (40) and (45), we know the order of  $|R_{1,n}|$  is bounded above by

$$o_p(\|b - \beta\|/\sqrt{n}) + o_p(n^{-1}h^{-k}) + o_p(n^{-3/2}h^{-k}) \quad (46)$$

uniformly over an  $o_p(1)$  neighborhood of  $\beta$ . The third term in (46) is  $o_p(n^{-1})$  due to choice of bandwidth in BW2. The second term is due to the product of  $\sup_z |\hat{p}(z) - p(z)|$  and a degenerate empirical process  $\frac{1}{n} \sum_{i=1}^n \varpi_{n,1}(z_i; b)$ . Next, let  $\delta_n = O(n^{-1/2}h^{-k/2})$ . Following the same arguments in (Khan 2001), another application of Theorem 3 in (Sherman 1994b) implies that over an  $O_p(\delta_n)$  neighborhood of  $\beta$ , the magnitude of this product would be  $h^{-k} O_p(\delta_n n^{-1/2}) o_p(n^{-1/2}) = o_p(h^{-3k/2} n^{-3/2})$ , which is  $o_p(n^{-1})$  given our choice of bandwidth in BW2.

We now deal with  $\Delta_{1,n}$ . We first derive the correction term due to estimation errors in  $\hat{p}(z_i)$ . Let  $\gamma_0 \equiv (\gamma_{0,1}, \gamma_{0,2})'$  denote  $\mathbb{E}[Y_i | z_i] f(z_i)$  and density  $f(z_i)$  in the population and let  $\hat{\gamma} \equiv (\hat{\gamma}_1, \hat{\gamma}_2)'$  denote their kernel estimates respectively so that  $\hat{\gamma}_1/\hat{\gamma}_2 = \hat{p}$ . With

a slight abuse of notation, let  $\kappa'(z_i, z_j)$  be a shorthand for  $\kappa'(w_{i,j} - 1)$ , and write the first half of  $\Delta_{1,n}$  as:

$$\begin{aligned} & \frac{1}{n(n-1)} \sum_{j \neq i} \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) [\hat{p}(z_i) - p(z_i)] \\ &= \frac{1}{n(n-1)} \sum_{j \neq i} \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) [\hat{\gamma}_1(z_i) / \hat{\gamma}_2(z_i) - \gamma_{0,1}(z_i) / \gamma_{0,2}(z_i)] \\ &= \frac{1}{n(n-1)} \sum_{j \neq i} \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \nabla w(z_i) [\hat{\gamma}(z_i) - \gamma_0(z_i)]' + \tilde{R}_{1,n} \end{aligned} \quad (47)$$

where  $\nabla w(z_i) \equiv [1/\gamma_{0,2}(z_i), -\gamma_{0,1}(z_i)/\gamma_{0,2}^2(z_i)]$ ; and  $\tilde{R}_{1,n}$  is of order  $o_p(\|b - \beta\|/\sqrt{n}) + O_p(n^{-1}h^{-k}) + o_p(n^{-1})$  uniformly over an  $o_p(1)$  neighborhood around  $\beta$  due to Taylor-expansion-based arguments similar to those applied to  $R_{1,n}$ . Also similar to the case with  $R_{1,n}$ , the second term in  $\tilde{R}_{1,n}$ , which is of the order  $O_p(n^{-1}h^{-k})$  uniformly over  $o_p(1)$  neighborhood of  $\beta$ , is further reduced to  $o_p(n^{-1})$  over an  $O_p(n^{-1/2}h^{-k/2})$  neighborhood around  $\beta$ , due to a repeated application of Theorem 3 in (Sherman 1994b) and our choice of bandwidth in BW2.

Next, let  $q \equiv (y, 1)'$ . Write the first term on the last line in (47) as

$$\begin{aligned} & \frac{1}{n(n-1)} \sum_{j \neq i} \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \nabla w(z_i) [\hat{\gamma}(z_i) - \mathbb{E}(\hat{\gamma}(z_i))] + \hat{R}_{1,n}; \\ & \text{where } \hat{R}_{1,n} \equiv \frac{1}{n(n-1)} \sum_{j \neq i} \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) [\mathbb{E}(\hat{\gamma}(z_i)) - \gamma_0(z_i)]. \end{aligned} \quad (48)$$

By triangular inequality, we have:

$$\hat{R}_{1,n} \leq \left\{ \frac{1}{n(n-1)} \sum_{j \neq i} |\kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \nabla w(z_i)| \right\} \sup_z |\mathbb{E}(\hat{\gamma}(z_i)) - \gamma_0(z_i)|. \quad (49)$$

By arguments similar to those apply to the first term in  $R_{1,n}$ , the first term in the product on the right-hand side of (49) is

$$o_p(\|b - \beta\|) + o_p(n^{-1/2}h^{-k}) + o_p(n^{-1}h^{-k})$$

Furthermore, the second term on the right-hand side of (49) is  $O(\sigma_n^m \kappa)$  due to Lemma 8.9 in (Newey and McFadden 1994), which is  $o_p(n^{-1/2})$  by our choice of  $\sigma_n$  in BW1 and smoothness condition in SM1. Thus the order of  $\hat{R}_{1,n}$  is no greater than  $o_p(\|b - \beta\|/\sqrt{n}) + o_p(n^{-1}h^{-k}) + o_p(n^{-3/2}h^{-k})$  uniformly over an  $o_p(1)$  neighborhood around  $\beta$ . Again, similar to the case with  $R_{1,n}$ , the  $o_p(n^{-1}h^{-k})$  term in  $\hat{R}_{1,n}$  above is further reduced to  $o_p(n^{-1})$  over an  $O_p(n^{-1/2}h^{-k/2})$  neighborhood around  $\beta$  by a repeated application of Theorem 3 in (Sherman 1994b).

We now write the first term in (48) as a third-order U-statistic:

$$\frac{1}{n(n-1)(n-2)} \sum_{i \neq j \neq s} \phi_n(\tilde{z}_i, \tilde{z}_j, \tilde{z}_s; b) \quad (50)$$

where  $\tilde{z} \equiv (y, z) \equiv (y, x, v)$  and

$$\phi_n(\tilde{z}_i, \tilde{z}_j, \tilde{z}_s; b) \equiv \kappa'(z_i, z_j) K_h(x_i - x_j) \Delta \varphi_{i,j}^-(b) \nabla w(z_i) \{q_s \mathcal{K}_\sigma(z_i - z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(z_i - Z_s)]\}$$

with  $\mathcal{K}_\sigma$  being a shorthand for  $\sigma^{-(k+1)} \mathcal{K}(\cdot/\sigma^{k+1})$ ; and the expectation is taken w.r.t.  $\tilde{Z}_s$  while  $z_i$  is some realized value of  $Z_i$ . Note  $\phi_n$  is not symmetric in the three arguments,

for it depends on  $y_s$  but not  $y_i$  and  $y_j$ . Let  $\tilde{Z}_{i,j,s}$  be a shorthand for  $(\tilde{Z}_i, \tilde{Z}_j, \tilde{Z}_s)$ . Then apply the H-decomposition to write this third-order U-statistic in (50) as

$$\mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b)] + \frac{1}{n} \sum_{i=1}^n \phi_n^{(1)}(\tilde{z}_i; b) + U^2 \phi_n^{(2)}(b) + U^3 \phi_n^{(3)}(b) \quad (51)$$

where

$$\phi_n^{(1)}(\tilde{z}_i) \equiv \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_i = \tilde{z}_i] + \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_j = \tilde{z}_i] + \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_s = \tilde{z}_i] - 3\mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b)] \quad (52)$$

and  $U^2 \phi_n^{(2)}(b)$  and  $U^3 \phi_n^{(3)}(b)$  are second- and third-order degenerate U-statistics as defined in (Sherman 1994b).

To deal with the second- and third-order processes  $U^2 \phi_n^{(2)}(b)$  and  $U^3 \phi_n^{(3)}(b)$ , we use the same arguments as in (Khan 2001). It follows from our conditions on BW2 and KF2, FM1,2 that the two classes  $\{h^k \phi_n^{(2)}(b) : b \in \mathcal{B}\}$  and  $\{h^k \phi_n^{(3)}(b) : b \in \mathcal{B}\}$  are both Euclidean. Besides, these conditions ensure condition (ii) of Theorem 3 in (Sherman 1994b) holds with the " $\gamma_n$ " therein being  $O(1)$  for any sequence of  $\delta_n$  converging to 0. Hence uniformly over an  $o_p(1)$  neighborhood of  $\beta$  in  $\mathcal{B}$ , the third-order term  $U^3 \phi_n^{(3)}(b)$  is  $h^{-k} O_p(n^{-3/2})$ , which is  $o_p(n^{-1})$  under our choice of bandwidth in BW2. The second-order term is  $O_p(h^{-k} n^{-1})$  over an  $o_p(1)$  neighborhood of  $\beta$ . Let  $\delta_n = O(h^{-k/2} n^{-1/2})$ . Furthermore, following the same arguments in (Khan 2001), another application of Theorem 3 implies that over an  $O_p(\delta_n)$  neighborhood of  $\beta$ , the second-order term is  $O_p(\delta_n n^{-1}) = O_p(h^{-k/2} n^{-3/2})$ , which is  $o_p(n^{-1})$  given our choice of bandwidth in BW2.

Next, we deal with the first-order term  $\phi_n^{(1)}$ . By definition,

$$\begin{aligned} & \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_i = (y, z)] \\ & \equiv \mathbb{E} \{ \kappa'(z, Z_j) K_h(x - X_j) \Delta \varphi_{i,j}^-(b) \nabla w(z) [Q_s \mathcal{K}_\sigma(z - Z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(z - Z_s)]] | Z_i = z \} \\ & = \mathbb{E} \{ \kappa'(z, Z_j) \nabla w(z) K_h(x - X_j) \Delta \varphi_{i,j}^-(b) | Z_i = z \} \{ \mathbb{E}[Q_s \mathcal{K}_\sigma(z - Z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(z - Z_s)]] | Z_i = z \} \end{aligned}$$

where the second term on the right-hand side is 0 by construction. Besides,

$$\begin{aligned} & \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_j = (y, z)] \\ & = \mathbb{E}_{Z_i} \{ \kappa'(Z_i, z) \nabla w(Z_i) K_h(X_i - x) \Delta \varphi_{i,j}^-(b) \mathbb{E}_{Z_s} \{ Q_s \mathcal{K}_\sigma(Z_i - Z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(Z_i - Z_s)] | Z_i, Z_j = z \} | Z_j = z \} \\ & = \mathbb{E}_{Z_i} \{ \kappa'(Z_i, z) \nabla w(Z_i) K_h(X_i - x) \Delta \varphi_{i,j}^-(b) \mathbb{E}_{Z_s} \{ Q_s \mathcal{K}_\sigma(Z_i - Z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(Z_i - Z_s)] | Z_i \} | Z_j = z \} \end{aligned}$$

where  $\mathbb{E}_{Z_s} [Q_s \mathcal{K}_\sigma(Z_i - Z_s) - \mathbb{E}[Q_s \mathcal{K}_\sigma(Z_i - Z_s)] | Z_i] = 0$  conditional on all  $Z_i$ . Hence this term is also degenerate at 0 for all  $Z_i$  and  $b$ . It then follows that the unconditional expectation  $\mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b)] = 0$  for all  $b$ . Hence the first two terms in the H-decomposition in (51) are reduced to:

$$\begin{aligned} & \frac{1}{n} \sum_{l=1}^n \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_s = (y_l, z_l)] \\ & = \frac{1}{n} \sum_{l=1}^n \mathbb{E}_{Z_i} \left\{ \tilde{m}_n(Z_i; b) [Q_s \mathcal{K}_\sigma(Z_i - Z_s) - \mathbb{E}_{Q', Z'} (Q' \mathcal{K}_\sigma(Z_i - Z'))] | \tilde{Z}_s = (y_l, z_l) \right\} \quad (53) \end{aligned}$$

where  $q_l \equiv (y_l, 1)'$ ; and

$$\tilde{m}_n(z; b) \equiv \mathbb{E}_{Z_j}[\nabla w(Z_i)\kappa'(Z_i, Z_j)K_h(X_i - X_j)\Delta\varphi_{i,j}^-(b)|Z_i = z].$$

To clarify notations, note on the second line of (53),  $\mathbb{E}[Q'\mathcal{K}_\sigma(Z_i - Z')]$  is a function of  $Z_i$  and the expectation is taken w.r.t.  $Q', Z'$ .

It remains to show that we can write the right-hand side of (53) as a sample average of some function of  $(z_l, y_l)$  plus a term that is smaller than  $o_p(\|b - \beta\|/\sqrt{n}) + o_p(\|b - \beta\|^2) + o_p(n^{-1})$  uniformly over  $o_p(1)$  neighborhood around  $\beta$ .

By changing variables between  $x_j$  and  $u \equiv h^{-k}(x_i - x_j)$  while fixing  $z_i$ , we have

$$\begin{aligned} \tilde{m}_n(z_i; b) &= \nabla w(z_i) \int \kappa'(z_i, z_j)K_h(x_i - x_j)\Delta\varphi^-(z_i, z_j; b)dF(z_j) \\ &= \nabla w(z_i) \int K_h(x_i - x_j)\tilde{\mu}^-(z_i, x_j; b)dF(x_j) \\ &= \nabla w(z_i) \int K(u)\tilde{\mu}^-(z_i, x_i - h^k u; b)f(x_i - h^k u)du \\ &= \nabla w(z_i)\tilde{\mu}^-(z_i, x_i; b)f(x_i) + O(h_n^{m_\varphi}) \end{aligned}$$

where  $\tilde{\mu}^-(z_i, x_j; b) \equiv \mathbb{E}[\kappa'(Z_i, Z_j)\Delta\varphi_{i,j}^-(b)|Z_i = z_i, X_j = x_j]$ . The first equality follows from independence between  $Z_i$  and  $Z_j$ ; the second from the law of iterated expectation; the third from changing variables between  $u$  and  $x_j$ ; and the last from applying a Taylor expansion of  $x_j$  around  $x_i$ , and using the boundedness of the derivatives under SM2 and WF and the order of  $K$  in KF2.

Let  $m_-^*(z; b) \equiv \nabla w(z)f(x)\tilde{\mu}^-(z, x; b)$ . Then (53) can be written as:

$$\begin{aligned} &\int \tilde{m}_n(z; b) \left( \frac{1}{n} \sum_{l=1}^n q_l \mathcal{K}_\sigma(z - z_l) - \int \mathbb{E}(Q'|x') \mathcal{K}_\sigma(z - z') f(z') dz' \right) dF(z) \\ &= \frac{1}{n} \sum_{l=1}^n \int q_l [m_-^*(z; b) + O(h_n^{m_\varphi})] \mathcal{K}_\sigma(z - z_l) f(z) dz \\ &- \int f(z') \mathbb{E}(Q'|x') \left( \int [m_-^*(z; b) + O(h_n^{m_\varphi})] \mathcal{K}_\sigma(z - z') f(z) dz \right) dz'. \end{aligned}$$

Thus this suggests  $\frac{1}{n} \sum_{l=1}^n \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_s = (y_l, z_l)]$  can be decomposed into the sum of the following two terms:

$$\frac{1}{n} \sum_{l=1}^n \int q_l m_-^*(z; b) \mathcal{K}_\sigma(z - z_l) f(z) dz - \int f(z') \mathbb{E}(Q'|x') \left( \int m_-^*(z; b) \mathcal{K}_\sigma(z - z') f(z) dz \right) dz' \quad (54)$$

and

$$O(h_n^{m_\varphi}) \left\{ \frac{1}{n} \sum_{l=1}^n \int q_l \mathcal{K}_\sigma(z - z_l) f(z) dz - \int f(z') \mathbb{E}(Q'|x') \left( \int \mathcal{K}_\sigma(z - z') f(z) dz \right) dz' \right\}. \quad (55)$$

We first examine the term in (55). Note  $\int \mathcal{K}_\sigma(z - z') f(z) dz = f(z') + O(\sigma_n^{m_\kappa})$  for all  $z'$  under our conditions, and the term in the braces above can be written as

$$\begin{aligned} & \frac{1}{n} \sum_{l=1}^n \int q_l \mathcal{K}_\sigma(z - z_l) f(z) dz - \int f(z') \mathbb{E}(Q'|x') [f(z') + O(\sigma_n^{m_\kappa})] dz' \\ &= \frac{1}{n} \sum_{l=1}^n \int f(z) q_l \mathcal{K}_\sigma(z - z_l) dz - \mathbb{E}[f(Z)Q] + O(\sigma_n^{m_\kappa}) \\ &= \frac{1}{n} \sum_{l=1}^n \{f(z_l)q_l - \mathbb{E}[f(Z)Q]\} + o_p(n^{-1/2}) + O(\sigma_n^{m_\kappa}) \end{aligned}$$

where the last equality follows from arguments identical to Theorem 8.11 in (Newey and McFadden 1994) and our choice of bandwidth in BW1. Also by our choice of bandwidth in BW1,2, both  $O(h_n^{m_\varphi})$  and  $O(\sigma_n^{m_\kappa})$  are  $o(n^{-1/2})$ . Note  $\frac{1}{n} \sum_{l=1}^n \{f(z_l)q_l - \mathbb{E}[f(Z)Q]\}$  is  $O_p(n^{-1/2})$  by the Central Limit Theorem. Thus the term in (55) is  $o_p(n^{-1})$  uniformly over an  $o_p(1)$  neighborhood around  $\beta$ .

Next, to deal with (54), for any  $z$ , we can apply a Taylor expansion of  $m_-^*$  around  $b = \beta$  to get:

$$m_-^*(z; b) = 0 + \nabla_b m_-^*(z; \beta)(b - \beta) + o(\|b - \beta\|).$$

Substituting this into (54) above, we decompose it into the sum of

$$(b - \beta) \left\{ \frac{1}{n} \sum_{l=1}^n \int q_l \nabla_b m_-^*(z; \beta) \mathcal{K}_\sigma(z - z_l) f(z) dz - \int f(z') \mathbb{E}(Q'|x') \left( \int \nabla_b m_-^*(z; \beta) \mathcal{K}_\sigma(z - z') f(z) dz \right) dz' \right\} \quad (56)$$

and

$$o(\|b - \beta\|) \left\{ \frac{1}{n} \sum_{l=1}^n \int q_l \mathcal{K}_\sigma(z - z_l) f(z) dz - \int f(z') \mathbb{E}(Q'|x') \left( \int \mathcal{K}_\sigma(z - z') f(z) dz \right) dz' \right\}, \quad (57)$$

where the latter term is of order smaller than  $o_p(\|b - \beta\| / \sqrt{n})$  as the term in the braces in (57) is  $O_p(n^{-1/2})$  by the same arguments above.

As for the term in the braces in (56), first note by standard arguments using change of variables, we have:

$$\int \nabla_b m_-^*(z; \beta) f(z) \mathcal{K}_\sigma(z - z') dz = \nabla_b m_-^*(z'; \beta) f(z') + O(\sigma_n^{m_\kappa}).$$

Thus the term in the braces of (56) is

$$\begin{aligned} & \frac{1}{n} \sum_{l=1}^n \int q_l \nabla_b m_-^*(z; \beta) f(z) \mathcal{K}_\sigma(z - z_l) dz - \int \mathbb{E}(Q'|x') \nabla_b m_-^*(z'; \beta) f(z')^2 dz' + O(\sigma_n^{m_\kappa}) \\ &= \left\{ \frac{1}{n} \sum_{l=1}^n \int q_l \nabla_b m_-^*(z; \beta) f(z) \mathcal{K}_\sigma(z - z_l) dz - \mathbb{E}[Qf(Z)\nabla_b m_-^*(Z; \beta)] \right\} + O(\sigma_n^{m_\kappa}) \quad (58) \end{aligned}$$

Again by arguments similar to above and citing same arguments from Theorem 8.11 in (Newey and McFadden 1994) under SM2 and FM2, the term in the braces of (58) is

$$\frac{1}{n} \sum_{l=1}^n \delta_-^*(y_l, z_l) + o_p(n^{-1/2}) \text{ where } \delta_-^*(y, z) \equiv q \nabla_b m_-^*(z; \beta) f(z) - \mathbb{E}[Q \nabla_b m_-^*(Z; \beta) f(Z)];$$



while  $O(\sigma_n^{m\kappa}) = o(n^{-1/2})$  under our choice of bandwidth. To sum up, we have shown

$$\frac{1}{n} \sum_{l=1}^n \mathbb{E}[\phi_n(\tilde{Z}_{i,j,s}; b) | \tilde{Z}_s = (y_l, z_l)] = \frac{1}{n} \sum_{l=1}^n \delta_-^*(y_l, z_l)(b - \beta) + o_p(\|b - \beta\|/\sqrt{n}) + o_p(n^{-1})$$

uniformly over an  $o_p(1)$  neighborhood around  $\beta$  in  $\mathcal{B}$ , where  $\delta_-^*$  is the correction term due to  $\hat{p}(z_i)$  in  $\hat{H}_n$ .

Because  $\hat{p}_i$  and  $\hat{p}_j$  enter the objective function in the same way, and  $p_i$  and  $p_j$  are additively separable in the first-order expansion, we can apply identical arguments above to derive another identical correction term due to the use of  $\hat{p}(z_j)$  in  $\hat{H}_n$ . This proves the claim of the lemma. Q.E.D.

Replicating the arguments in the preceding lemma we can prove a result similar to Lemma C3 holds for the other half of the difference between "feasible" and "infeasible" objective function  $|\hat{\mathcal{H}}_{2,n}(b) - \mathcal{H}_{2,n}(b)|$ , except that  $\delta_-^*$  needs to be replaced by

$$\delta_+^*(y, z) \equiv q \nabla_b m_+^*(z; \beta) f(z) - \mathbb{E}[Qf(Z) \nabla_b m_+^*(Z; \beta)]$$

where

$$\begin{aligned} m_+^*(z) &\equiv \nabla w(z) f(x) \tilde{\mu}^+(z, x; b); \text{ with} \\ \tilde{\mu}^+(z_i, x_j; b) &\equiv \mathbb{E}[\kappa'(Z_i, Z_j) \Delta \varphi^+(Z_i, Z_j; b) | Z_i = z_i, X_j = x_j]. \end{aligned}$$

Building on the preceding Lemmas, we are now ready to prove the final result about the limiting distribution of  $\hat{\beta}$ .

**Proof of Proposition 7.** By Lemma C2 and Lemma C3,

$$\hat{H}_n(b) = H_0(b) + O_p(\|b - \beta\|/\sqrt{n}) + o_p(\|b - \beta\|^2) + O_p(n^{-1}h^{-k}) \quad (59)$$

uniformly over an  $o_p(1)$  neighborhood around  $\beta$  in  $\mathcal{B}$ . Recall  $H_0(b)$  is minimized at  $b = \beta$  due to Proposition 5. Hence it follows from (59), Lemma C1 above and Theorem 1 in (Sherman 1994b) that  $\hat{\beta}$ , as the minimizer of  $\hat{H}_n(b)$  over  $b \in \mathcal{B}$ , converges to  $\beta$  at a rate of  $1/\sqrt{nh^k}$ . As stated in Lemma C2 and Lemma C3, the  $O_p(n^{-1}h^{-k})$  term in (59) is further reduced to  $o_p(n^{-1})$  under conditions of the proposition. Hence another application of Theorem 1 in (Sherman 1994b) suggests  $\|\hat{\beta} - \beta\| = O_p(n^{-1/2})$ .

Recall that by a second-order Taylor expansion,  $H_0(b) = \frac{1}{2}(b - \beta)' \nabla_{bb} H_0(\beta)(b - \beta) + o_p(\|b - \beta\|^2)$  over an  $o(1)$  neighborhood of  $\beta$ , for  $\nabla_b H_0(\beta) = 0$  by construction. This, together with Lemma C2 and Lemma C3 and the root-n convergence shown in the previous paragraph, suggests that

$$\hat{H}_n(b) = \frac{1}{2}(b - \beta)' \nabla_{bb} H_0(\beta)(b - \beta) + \frac{1}{n} \sum_{i=1}^n 2[\delta_-^*(\tilde{z}_i) + \delta_+^*(\tilde{z}_i)](b - \beta) + o_p(n^{-1})$$

uniformly over an  $O_p(n^{-1/2})$  neighborhood around  $\beta$ . The limiting distribution then follow from Theorem 2 in (Sherman 1994b) and that  $\mathbb{E}[\delta^*(\tilde{Z})\delta^*(\tilde{Z})'] < \infty$  under FM2. Q.E.D.

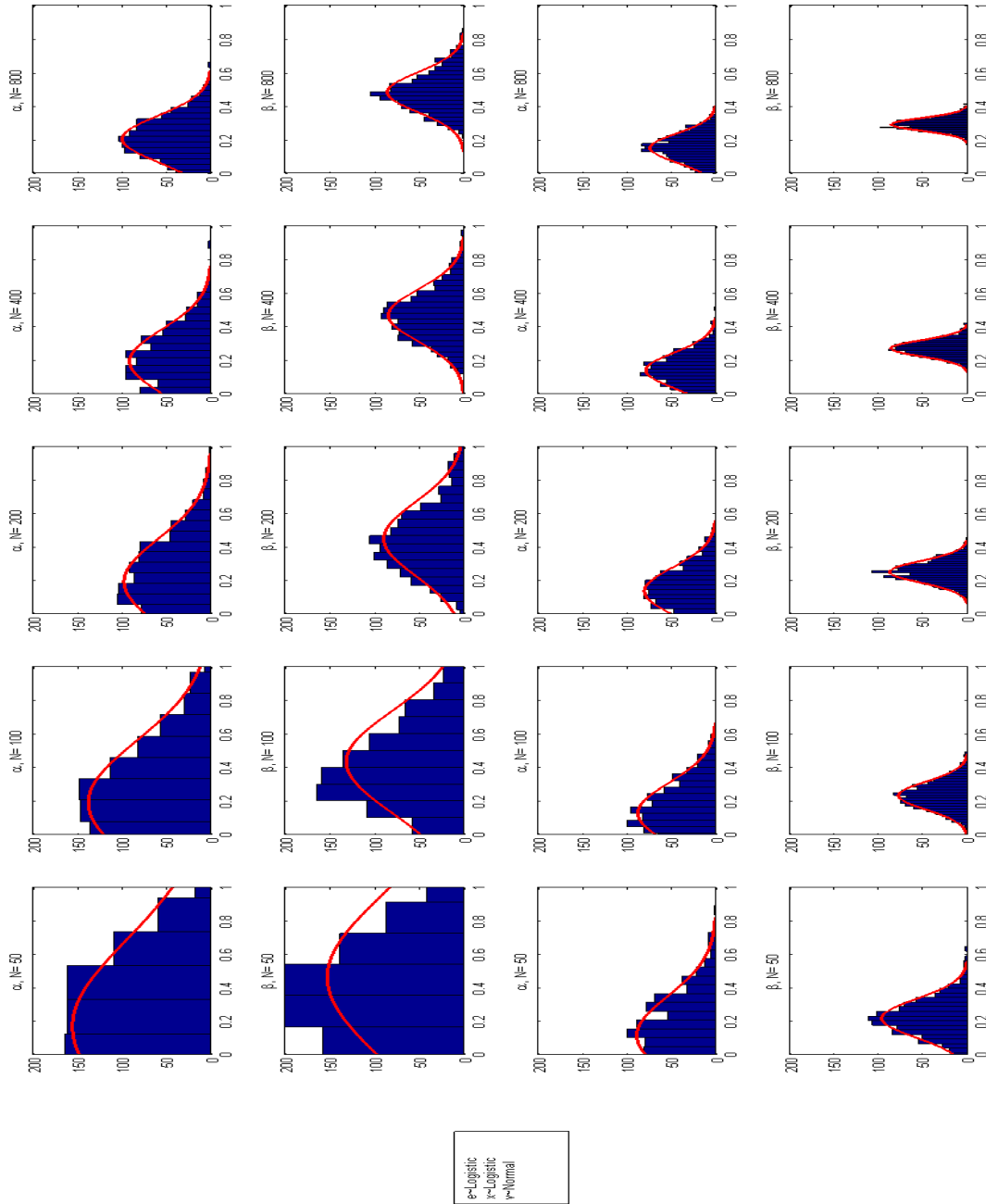


Figure 1:  $(X, \epsilon) \sim (\text{logistic}, \text{logistic})$ . First two rows: Pairwise extremum estimator. Last two rows: Inverse-density-weighted estimator.

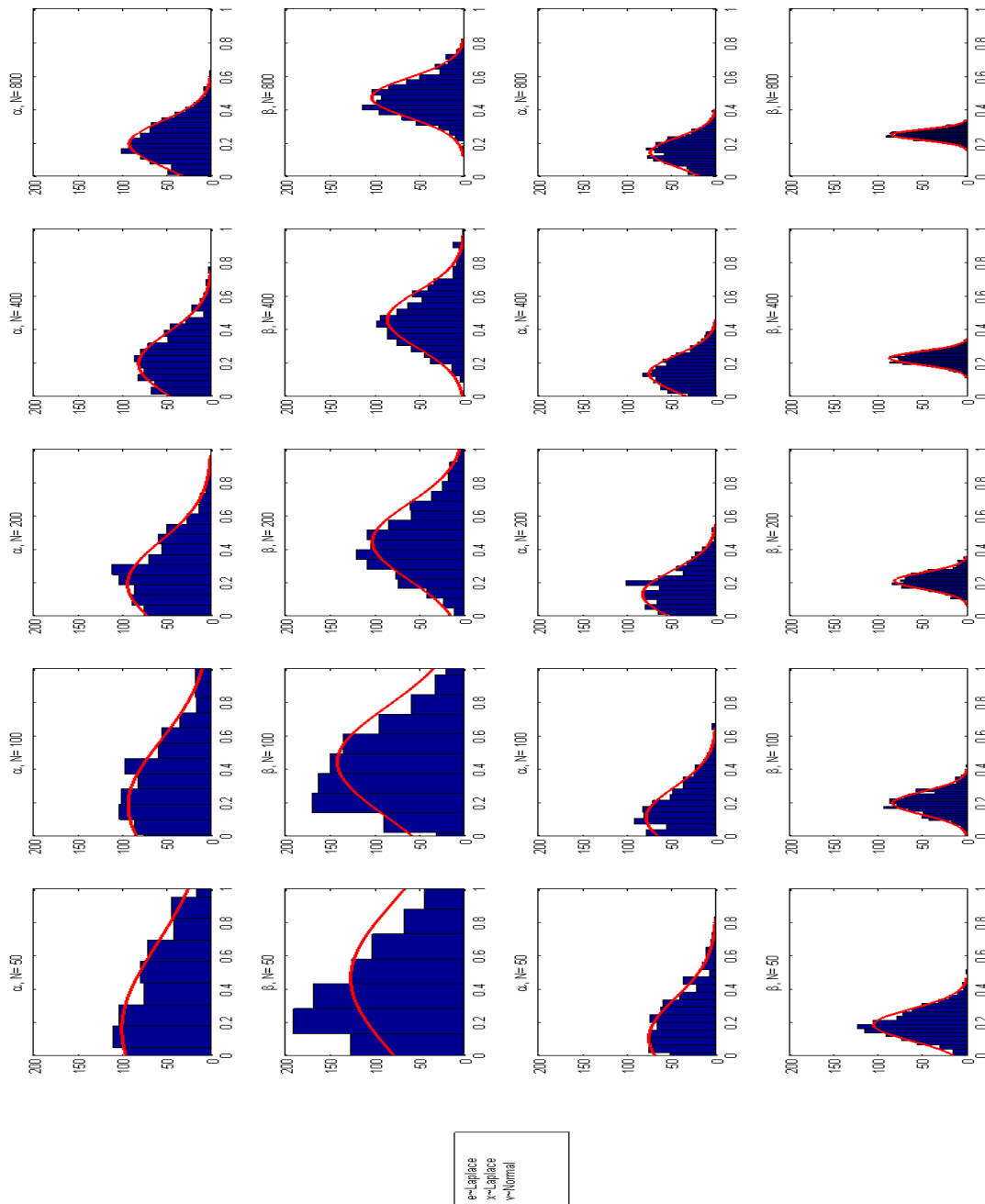
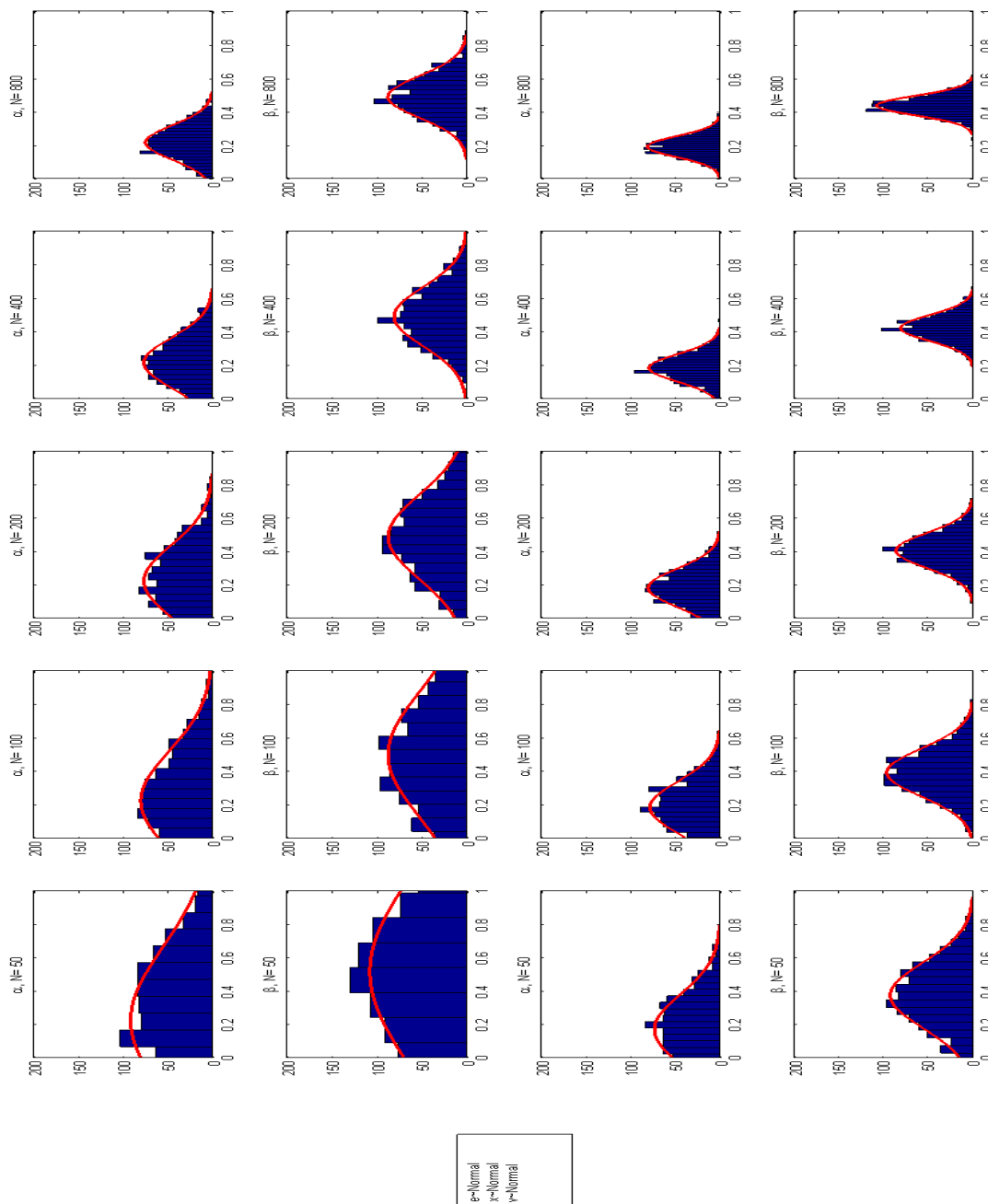


Figure 2:  $(X, \epsilon) \sim (\text{Laplace}, \text{Laplace})$ .

Figure 3:  $(X, \epsilon) \sim (\text{Normal}, \text{Normal})$

## References

- AMEMIYA, T. (1985): *Advanced Econometrics*. Harvard University Press.
- BERRY, S., AND P. HAILE (2010): “Nonparametric Identification of Multinomial Choice Demand Models with Heterogeneous Consumers,” Discussion paper no. 1718, Yale University.
- BLUNDELL, R., AND J. POWELL (2003): “Endogeneity in Nonparametric and Semiparametric Regression Models,” in *Advances in Economics and Econometrics: Theory and Applications, Eighth World Congress*, ed. by M. Dewatripont, L. Hansen, and S. Turnovsky, vol. II. Cambridge University Press.
- CHAMBERLAIN, G. (1986): “Asymptotic Efficiency in Semiparametric Models with Censoring,” *Journal of Econometrics*, 32, 189–218.
- CHEN, S. (2005): “Semiparametric Estimation of a Heteroscedastic Binary Choice Model,” Working paper, HKUST.
- CHEN, S., AND S. KHAN (2003): “Rates of Convergence for Estimating Regression Coefficients in Heteroskedastic Discrete Response Models,” *Journal of Econometrics*, 117, 245–278.
- CHERNOZHUKOV, V., H. HONG, AND E. TAMER (2007): “Estimation and Confidence Regions for Parameter Sets in Econometric Models,” *Econometrica*, 75, 1243–1284.
- COSSLETT, S. (1987): “Efficiency Bounds for Distribution Free Estimators of the Binary Choice Model,” *Econometrica*, 51, 765–782.
- FOX, J., AND C. YANG (2012): “Unobserved Heterogeneity in Matching Games,” Working paper, University of Michigan.
- HOROWITZ, J. (1992): “A Smoothed Maximum Score Estimator for the Binary Response Model,” *Econometrica*, 60, 505–531.
- ICHIMURA, H. (1993): “Local quantile regression estimation of binary response models with conditional heteroskedasticity,” Working paper, University of Minnesota.
- ICHIMURA, H., AND S. LEE (2010): “Characterizing Asymptotic Distributions of Semiparametric M-Estimators,” *Journal of Econometrics*, 159, 252–266.
- KHAN, S. (2001): “Two Stage Rank Estimation of Quantile Index Models,” *Journal of Econometrics*, 100, 319–355.
- KHAN, S., AND E. TAMER (2010): “Irregular Identification, Support Conditions, and Inverse Weight Estimation,” *Econometrica*, 78, 2021–2042.

- LEWBEL, A. (2000): “Semiparametric qualitative response model estimation with unknown heteroskedasticity or instrumental variables,” *Journal of Econometrics*, 97, 145–177.
- LEWBEL, A., AND X. TANG (2012): “Identification and Estimation of Games with Incomplete Information Using Excluded Regressors,” Working paper, Boston College and U Penn.
- MAGNAC, T., AND E. MAURIN (2007): “Identification and information in monotone binary models,” *Journal of Econometrics*, 139, 76–104.
- MANSKI, C. (1985): “Semiparametric Analysis of Discrete Response: Asymptotic Properties of the Maximum Score Estimator,” *Journal of Econometrics*, 27, 313–333.
- (1988): “Identification of binary response models,” *Journal of the American Statistical Association*, 83, 729–738.
- MANSKI, C., AND E. TAMER (2002): “Inference on Regressions with Interval Data on a Regressor or Outcome,” *Econometrica*, 70(2), 519–546.
- NEWBY, W., AND D. MCFADDEN (1994): “Large sample estimation and hypothesis testing,” in *Handbook of Econometrics*, ed. by R. F. Engle, and D. L. McFadden, vol. 4 (36), pp. 2111–2245. Elsevier.
- PAKES, A., AND D. POLLARD (1989): “Simulation and the Asymptotics of Optimization Estimators,” *Econometrica*, 57(5), 1027–57.
- POWELL, J. (1994): “Estimation of Semiparametric Models,” in *Handbook of Econometrics*, ed. by R. F. Engle, and D. L. McFadden, vol. 4 (41), pp. 2443–2521. Elsevier.
- SHERMAN, R. (1994a): “Maximal inequalities for degenerate U-processes with applications to optimization estimators,” *Annals of Statistics*, 22, 439–459.
- (1994b): “U-processes in the analysis of a generalized semiparametric regression estimator,” *Econometric Theory*, 10, 372–395.
- ZHENG, X. (1995): “Semiparametric efficiency bounds for the binary choice and sample selection models under conditional symmetry,” *Economics Letters*, 47, 249–253.