

# Semiparametric Selection Models with Binary Outcomes

Roger Klein

Chan Shen

Francis Vella

Rutgers University

The University of Texas

Georgetown University

MD Anderson Cancer Center

## Abstract

This paper addresses the estimation of a semiparametric sample selection index model where both the selection rule and the outcome variable are binary. Since the marginal effects are often of primary interest and are difficult to recover in a semiparametric setting, we develop estimators for both the marginal effects and the underlying model parameters. The marginal effect estimator only uses observations which are members of a high probability set in which the selection problem is not present. A key innovation is that this high probability set is data dependent. The model parameter estimator is a quasi-likelihood estimator based on regular kernels with bias corrections. We establish their large sample properties and provide simulation evidence confirming that these estimators perform well in finite samples.

## 1 Introduction

Despite the substantial literature extending the sample selection model of Heckman (1974, 1979) there is no semiparametric estimation of marginal effects in a model where both the outcome variable and the selection rule are binary.<sup>1</sup> This represents a significant void as important empirical examples exist in many areas of micro economics. In the fully parametric setting both the model parameters and the marginal effects, the objects which are generally of primary interest, are easily obtainable. Less is known for the semiparametric index model considered here. There have been some important developments in the literature including Chesher (2005), Vytlačil and Yildiz (2007), and Shaikh and Vytlačil (2011) that discuss identification of the marginal impact of a discrete endogenous variable. However, detailed estimation of marginal effects has not been addressed for the case of sample selection.

This paper develops semiparametric estimators for both the marginal effects and the index parameters underlying them. We make no distributional assumptions and allow a model structure more general than

---

<sup>1</sup>For a survey see Vella (1998).

threshold-crossing. Our primary focus is upon the marginal effects as they have not been addressed in this setting. Unlike in the parametric case, these effects cannot be directly derived from parameter estimates in a semiparametric setting because the error distributions are unknown. Moreover, the relevant distribution is difficult to estimate because the outcome equation is only observed for the selected sample. We propose to estimate the relevant distribution by focusing on those observations in an estimated high probability set where the selection probability tends to one. The framework of this approach is developed in pioneering papers of Heckman (1990) and Andrews and Schafgans (1998) for a known high probability set. This set depends on the tail behavior of index and error distributions. Therefore, in practice it is important to study the empirical tail behavior so as to find the appropriate high probability set. In this paper we characterize the high probability set as one where the probability exceeds a cutoff that approaches one as the sample size increases. We propose and establish the theoretical properties for an estimator of this cutoff that depends on empirical tail behavior. Based on the estimated high probability set, we formulate a marginal effect estimator and provide the theory for it which takes the estimation of this set into account. This data-dependent feature of the high probability set underlying the marginal effect estimator poses a number of theoretical challenges, but is essential in empirical studies.

Estimation of the marginal effects requires estimates of the index parameters and we propose a likelihood-based procedure employing a double index formulation. Identification issues are explicitly treated in Newey (2007) although that paper does not address estimation. Our index parameter estimator employs bias adjustment mechanisms similar to those developed for single index regression models in Klein and Shen (2010). We develop an estimator based on regular kernels and show that it has both desirable theoretical properties and good finite sample performance.<sup>2</sup> It is possible to develop index parameter estimators within various frameworks (see, e.g. Gallant and Nychka (1987), Klein and Spady (1993), Ichimura and Lee (1991), Lee (1995), and Klein and Vella (2009)). Most recently Escanciano, Jacho-Chavez and Lewbel (2012) propose semiparametric estimators for the index parameters based on higher order kernels. However, the estimation of marginal effect remains unaddressed.

Section 2 provides the model while Section 3 discusses estimation strategies for the marginal effects and the index parameters. Section 4 provides the assumptions and definitions. Section 5 presents the asymptotic results. Section 6 provides simulation evidence and concluding comments are offered in Section

---

<sup>2</sup>There are other alternative methods that control for the bias under regular kernels. For example, Honore and Powell (2005) employ a jackknife approach where the final estimator is a linear combination of estimators using different windows.

7. The Appendix contains an illustrative example and all proofs.

## 2 Model

The model is a semiparametric variant on the Heckman (1974, 1979) selection model where the outcome of interest is binary. More explicitly:

$$Y_1^* = I\{g(X\beta_0, \epsilon) > 0\} \tag{1}$$

$$Y_2 = I\{h(Z\pi_0) > u\}, \tag{2}$$

where  $Y_1^*$  is only observed for the subsample for which  $Y_2 = 1$ . Here  $I\{\cdot\}$  is an indicator function;  $X$  and  $Z$  are vectors of exogenous variables;  $\epsilon_i$  and  $u_i$  are error terms with a non-zero correlation;  $g(\cdot)$  and  $h(\cdot)$  are unknown functions with  $h(\cdot)$  being increasing; and  $\beta_0$  and  $\pi_0$  are unknown parameter values. When the model is additive, and the joint distribution of the errors is parametrically known, it can be estimated by maximum likelihood. However, without separability or known error distributions, the existing estimators in the literature do not apply. Our proposed estimator for index parameters also applies to the case where  $Y_2$  does not have a threshold-crossing structure. However, for the marginal effect estimator, the theory in this paper requires that the selection equation ( $Y_2$ ) has a threshold-crossing structure. Without loss of generality, we simplify the  $Y_2$ -model by replacing  $h$  with the identity function.

As in most semiparametric models the parameters are identified up to location and scale. Writing

$$X\beta_0 = b_1(X_1 + X_2\theta_{10}) + c_1 \equiv b_1V_{10} + c_1$$

$$Z\pi_0 = b_2(Z_1 + Z_2\theta_{20}) + c_2 \equiv b_2V_{20} + c_2,$$

the  $\theta'_0$ s are identified, while the  $b$ 's and  $c$ 's are not identified. We refer to  $V_{10}$  and  $V_{20}$  as indices and assume that the model satisfies the following index restrictions:

$$\Pr(Y_1 = d_1, Y_2 = d_2|X, Z) = \Pr(Y_1 = d_1, Y_2 = d_2|V_{10}, V_{20}) \tag{3}$$

$$\Pr(Y_2 = d_2|X, Z) = \Pr(Y_2 = d_2|V_{20}) \tag{4}$$

$$\Pr(Y_1 = d_1|X, Z) = \Pr(Y_1 = d_1|V_{10}). \tag{5}$$

We note that the above conditions hold if the errors are independent of  $X$  and  $Z$ . We impose this index structure, as opposed to a non-parametric one, to improve the performance of the estimators.

### 3 Estimation

#### 3.1 Marginal Effects

Our marginal effect of interest is the change in  $\Pr(Y_1 = 1|X) = \Pr(Y_1 = 1|V_{10})$  due to a change in one of the explanatory  $X$ -variables. To motivate this marginal effect, let  $Y_2$  denote whether or not an individual decides to have a diagnostic test for a particular genetic disease and let  $Y_1$  denote whether or not an individual has that disease. We would like to know how a change in one of the  $X$ -variables affects the probability of having the disease for the entire population and not just the subgroup that received the diagnostic test. In the fully parametric case (e.g. bivariate probit with selection) the probability of having the disease  $\Pr(Y_1 = 1|V_{10})$  is a known function, and the corresponding marginal effect of interest can be directly calculated once the parameters of the model are estimated.

Now consider the semiparametric case where the functional form of this probability function is not known. Under index restrictions, the probability of interest can be written as:

$$\begin{aligned} \Pr(Y_1 = 1|V_{10}) &= \Pr(Y_1 = 1|V_{10}, V_{20}) \\ &= \Pr(Y_1 = 1|Y_2 = 1, V_{10}, V_{20}) P_2 \\ &\quad + \Pr(Y_1 = 1|Y_2 = 0, V_{10}, V_{20}) (1 - P_2). \end{aligned}$$

where  $P_2 = \Pr(Y_2 = 1|V_{10}, V_{20}) = \Pr(Y_2 = 1|V_{20})$ . We can recover the first argument on the right hand side semiparametrically. That is, we can estimate the probability of having the disease given that the individual was tested and the probability that an individual elects to be tested. The question then becomes how to recover the second part:  $\Pr(Y_1 = 1|Y_2 = 0, V_{10}, V_{20}) (1 - P_2)$ . In general, this is not estimable because we do not observe  $Y_1$  (genetic disease) when  $Y_2 = 0$  (no testing). However, if  $P_2 = 1$  this second term disappears and we can estimate the marginal effect of interest based only on the first term. In an approach related to that in Heckman (1990) and Andrews and Schafgans (1998, hereafter referred to as A&S), we estimate the marginal effect by only using those observations for which the selection probability  $P_2$  is high. With  $N$  as the full sample size,  $F$  as the distribution function for the selection error  $u$ , and  $a > 0$ , the high probability

set is defined as

$$\{v_{20} : F(v_{20}) > 1 - N^{-a}\}.$$

The probability of being in this high probability set is given by  $P_h = \Pr(V_{20} > F^{-1}(1 - N^{-a})) = 1 - G(F^{-1}(1 - N^{-a}))$ , where  $G$  is the distribution function for the selection index  $V_{20}$ . For example, when the index has a standard Weibull distribution  $G = 1 - \exp(-v_{20})$ , and the error follows a Weibull distribution with thinner tail  $F = 1 - \exp(-u^c)$ ,  $c > 1$ ,  $P_h = \exp(-[-\ln(N^{-a})]^{1/c})$ . As the error tails become thinner relative to the index tail ( $c$  increases),  $P_h$  increases. This example demonstrates that the appropriate value for  $a$  depends on the thickness of index tails relative to that for the error. As these tails are unknown, we propose a data dependent value for  $a$  and establish and establish its asymptotic properties..

To describe our estimation strategy, we introduce  $\zeta_0(\bar{v}) \equiv \Pr(Y_1 = 1 | V_{10} = \bar{v})$  and write the true marginal effect as:

$$ME = \zeta_0(v_e) - \zeta_0(v_b),$$

where  $v_b$  refers to a base or initial level of the index and  $v_e$  refers to an index evaluated at a new level for the explanatory variable of interest.<sup>3</sup> Our estimator for the probability of interest is given by:

$$\hat{\zeta}(\bar{v}, \hat{a}) \equiv \frac{\sum_j \left\{ \frac{1}{N_h} Y_{1j} K[(\bar{v} - V_{1j})/h] \right\} (Y_{2j} \hat{S}_j)}{\sum_j \left\{ \frac{1}{N_h} K[(\bar{v} - V_{1j})/h] \right\} (Y_{2j} \hat{S}_j)}$$

where  $\hat{S}_j$  is a smoothed indicator of the form in A&S that is one on a high probability set.

To motivate this estimator, notice that a traditional semiparametric estimator for the probability of interest would have the following form if there is no sample selection issue:

$$\widehat{\Pr}(Y = 1 | V = \bar{v}) = \frac{\sum_j \left\{ \frac{1}{N_h} Y_{1j} K[(\bar{v} - V_{1j})/h] \right\}}{\sum_j \left\{ \frac{1}{N_h} K[(\bar{v} - V_{1j})/h] \right\}}.$$

Our estimator  $\hat{\zeta}(\bar{v}, \hat{a})$  differs from this estimated probability in two respects. First, as we only observe  $Y_1$  when  $Y_2 = 1$ , we need to have  $Y_2$  in both the numerator and the denominator so as to select observations for which  $Y_2 = 1$ . This introduces sample selection bias, and we compensate for that by adding the smooth high probability indicator  $\hat{S}_j$ .

---

<sup>3</sup>This definition can be applied to both discrete and continuous variables. For the latter, it would also be possible to define the marginal effect as a derivative.

### 3.2 Index Parameters

While our proposed marginal effect estimator is the primary focus, it depends on estimated index parameters. These are obtained by maximizing a quasi or estimated likelihood:

$$\begin{aligned}\hat{\theta} &\equiv \arg \max_{\theta} \hat{L}(\theta), \\ \hat{L}(\theta) &\equiv \sum_{i=1}^N \tau_i \sum_{d_1, d_2} Y_i(d_1, d_2) \text{Ln} \left( \hat{P}_i(d_1, d_2; \theta) \right),\end{aligned}$$

where

$$\begin{aligned}Y_i(d_1, d_2) &= \begin{cases} I\{Y_{1i} = d_1, Y_{2i} = d_2\} & \text{for } d_2 = 1 \\ I\{Y_{2i} = d_2\} & \text{for } d_2 = 0 \end{cases} \\ \hat{P}_i(d_1, d_2; \theta) &\equiv \hat{P}r(Y_i(d_1, d_2) = 1 | V_i(\theta) = v_i(\theta)).\end{aligned}$$

Here  $V_i(\theta) = (V_{1i}(\theta), V_{2i}(\theta))$  and  $\tau_i$  is a trimming function defined below to control for small density denominators.

The properties of the estimates depend on how the likelihood probabilities are estimated. We employ regular kernels and several bias-reducing mechanisms to ensure that the estimator has desirable large sample properties and also performs well in finite samples.

To motivate these mechanisms we show below that the gradient to the quasi-likelihood is a product of terms, one of which is the derivative of the probability function,  $\nabla_{\theta} \hat{P}_i(d_1, d_2; \theta_0)$ , where  $\theta_0$  denotes the true parameter value. Subject to some issues that we address below, the key to our bias reduction mechanisms is the result due to Whitney Newey (see Klein and Shen (2010), Theorem 0) that:

$$E(\nabla_{\theta} P_i(d_1, d_2; \theta_0) | V_i(\theta_0)) = 0.$$

## 4 Assumptions and Definitions

We now provide the assumptions and definitions that we employ to establish the asymptotic properties for the estimators.

**A1.** The observations are i.i.d. from the model in (1)-(2), where the matrices  $X$  and  $Z$  have full rank with

probability 1.

- A2.** The vector of true parameter values  $(\theta_{10}, \theta_{20})$  lies in the interior of a compact parameter space,  $\Theta$ .
- A3.** The indices  $V_1$  and  $V_2$  each contains a continuous exogenous variable. Further,  $V_2$  contains at least one continuous variable which does not enter  $V_1$  in any form. The model satisfies index restrictions as in (3-5).
- A4.** Let  $g(v_1, v_2|Y_1, Y_2)$  be the conditional density for the indices. Letting  $\nabla^p g$  be any of the partials or cross partials of  $g$  up to order  $p$ , with  $\nabla^0 g = g$ , assume  $g > 0$  on all fixed compact subsets of the support for the indices, and  $\nabla^p g$ ,  $\frac{\partial}{\partial \theta}(\nabla^p g)$ , and  $\frac{\partial^2}{\partial \theta \partial \theta}(\nabla^p g)$  are bounded for  $p = 0, 1, 2, 3, 4$ .
- A5.** Let  $F$  be the distribution for the selection error,  $G$  the distribution function for the selection index, and  $G_c$  be the conditional distribution of  $v_2|V_1 = \bar{v}$ . Characterize the high probability set as  $\{v_2 : P_2 \equiv F(v_2) \geq 1 - N^{-a}\}$ . Assume: **(a)** For all  $t > T$  sufficiently large,  $1 - G(t) > 1 - F(t)$  and  $1 - G_c(t) > 1 - F(t)$ . **(b)** The marginal density for the selection index  $g(v_2)$ , is decreasing in the tail. With  $g(v_u) \equiv O(N^{-\iota})$  where  $\iota$  is a small positive number, and  $H(v_u) \equiv \frac{g(v_u)}{1-G(v_u)}$  as the hazard for  $V_2$ :  $\frac{1-F(v_u)}{1-G(v_u)} < H(v_u)N^{-(a-\iota)}$ .
- A6.** Let  $g(v_2|\bar{v})$  be the density for  $V_2$  conditioned on  $V_1 = \bar{v}$ . For all  $t > T$  sufficiently large, assume that  $O(g(t)) \geq O(g(t|\bar{v}))$ .
- A7.** Assume  $\Pr(Y_1 = d_1|Y_2 = d_2, V_1 = v_1, V_2 = v_2)$  has up to four bounded derivatives with respect to  $v_1$  at  $\bar{v}$ .

The first three assumptions are standard in index models. Assumption (A4) provides required smoothness conditions for determining the order of the bias for density estimators. Similar to A&S, assumption (A5) is needed to develop the large sample distribution for the estimator of marginal effects in the outcome equation. For (A5a), as is well known in the literature (see e.g. Kahn and Tamer (2010)), tail conditions are needed to develop the large sample distribution of these types of estimators. The error and index supports can be finite provided these tail conditions hold. We note that when the error has a bounded support that is a subset of that for the index, this assumption holds. When the index support is a subset of that for the error, this assumption will not hold. Below, once we have defined the bias-variance trade-off (D3), we will illustrate the problem that results when these conditions do not hold.

Assumption (A5b) is required for trimming arguments. Let  $v_l$  be a value of the selection index such that the selection probability  $P_2(v_2) \equiv F(v_2) \geq 1 - N^{-a}$  for  $v_2 \geq v_l$ . Let  $v_u$  be a value of the selection index such that  $g(v_2) \geq N^{-\iota}$  for  $v_2 \leq v_u$ . To avoid a conflict in these conditions, we need to guarantee that  $v_l < v_u$ . This inequality will hold provided that

$$1 - F(v_u) < 1 - F(v_l) \equiv N^{-a}.$$

Dividing both sides by  $1 - G(v_u)$  and noting that  $g(v_u) \equiv N^{-\iota} = [1 - G(v_u)]H(v_u)$ , where  $H$  is the hazard function for  $V_2$ , it suffices that

$$\frac{1 - F(v_u)}{1 - G(v_u)} < \frac{N^{-a}}{1 - G(v_u)} = H(v_u)N^{-(a-\iota)}.$$

This condition holds when error and index densities have the Weibull form:

$$1 - F(u) = \exp(-u^{c_u}); \quad 1 - G(v) = \exp(-v^{c_v}),$$

where  $c_u > c_v$  so that index tails are fatter than those of the selection error.

Assumptions (A6-7) are used in Lemma 3 to derive the order of the bias in estimating marginal effect components. In addition to the above assumptions, we also need a number of definitions for densities, probability functions and estimators. These are given below.

**D1. The estimator for marginal effects.**

$$\hat{\zeta}(\bar{v}) \equiv \frac{\sum_j \frac{1}{Nh} Y_{1j} K[(\bar{v} - V_{1j})/h] Y_{2j} \hat{S}_j}{\sum_j \frac{1}{Nh} K[(\bar{v} - V_{1j})/h] Y_{2j} \hat{S}_j},$$

where  $K$  is a regular kernel with window  $h = O(N^{-.2-\epsilon})$ ,  $\epsilon$  is a small positive value, and  $\hat{S}_j$  is a smoothed trimming function on an estimated high probability set (see (D2-3)).

**D2. The  $S$ -function.** With  $b > 0$ , the  $S$ -function (adapted from A&S) is given as:

$$S(x) = \begin{cases} 0, & R1 : x \leq 0 \\ 1 - \exp \frac{-x^k}{b^k - x^k}, & R2 : 0 < x < b \\ 1, & R3 : x \geq b. \end{cases}$$



With  $\tau$  as an indicator restricting the density for  $V_2$  to be above  $O(N^{-\iota})$  where  $\iota$  is a small positive number<sup>4</sup>,

$$x \equiv \tau \left[ \text{Ln} \left( \frac{1}{1 - P_2} \right) - \text{Ln} (N^a) \right].$$

The integer  $k$  is set to insure that  $S$  is as many times differentiable as is needed at  $x = 0$ . With  $\hat{P}_{aj}$  defined below as an estimator for the selection probability, let  $\hat{S}_j \equiv S(x(\hat{a}, \hat{P}_{aj})) \equiv S(\hat{a}, \hat{P}_{aj})$ , and  $S_0 \equiv S(x(a_{0N}, P_2)) \equiv S(a_{0N}, P_2)$ .

**D3. True and estimated high probability parameters**  $a_{0N}$  and  $\hat{a}$ . With  $K_2$  a normal twicing kernel (Newey et al. (2004)),  $h_2 = O(N^{-1})$ ,

$$\begin{aligned} \hat{E}_2(\hat{S}) &\equiv \frac{1}{N} \sum_j \hat{S}_j \\ \hat{E}_2(\hat{S}^\kappa | \bar{v}) &\equiv \frac{\sum_j \hat{S}_j^\kappa K_2[(\bar{v} - V_{1j})/h_2]}{\sum_j K_2[(\bar{v} - V_{1j})/h_2]} \\ \text{where } \hat{P}_{aj} &\equiv \frac{\sum_i Y_{2i} K_2[(V_{2j} - V_{2i})/h_2]}{\sum_i K_2[(V_{2j} - V_{2i})/h_2]}. \end{aligned}$$

Then

$$\begin{aligned} a_{0N} &= \arg \min_{a \in \mathcal{A}} \left[ hN^{1-2a} \frac{[E(S_0)]^2}{E(S_0^2 | \bar{v})} - N^{-\eta} \right]^2 \\ \hat{a} &= \arg \min_{a \in \mathcal{A}} \left[ hN^{1-2a} \frac{[\hat{E}_2(\hat{S})]^2}{\hat{E}_2(\hat{S}^2 | \bar{v})} - N^{-\eta} \right]^2 \end{aligned}$$

where  $h = N^{-2-\varepsilon}$  from (D1),  $\mathcal{A} = \{a : 0 < a < .4 - \frac{\varepsilon-\eta}{2}\}$ , and  $\eta = 1/\ln(N)^\varkappa$  for  $\varkappa$  a small positive number; for  $N$  sufficiently large  $\frac{[E(S)]^2}{E(S^2 | \bar{v})} \leq 1$  and is an increasing function of  $N^{-a}$  for  $N$  sufficiently large.<sup>5</sup>

<sup>4</sup>It can be shown that trimming based on a density estimator is asymptotically equivalent to trimming on the true density.

<sup>5</sup>Note that

$$\frac{[E(S)]^2}{E(S^2 | \bar{v})} \leq 1 \Rightarrow hN^{1-2a} \leq N^{-\eta},$$

which in turn implies that  $a \in \mathcal{A}$ .

With the S-function in (D2) smoothly restricting observations to a high probability set, the moment conditions in (D3) reflect the bias-variance trade-off in estimating the marginal effect. From Lemmas 3-4:

$$O\left(\frac{\text{bias}^2}{\text{var}}\right) \leq \frac{O(hN^{1-2a}) [E(S)]^2}{E(S^2|\bar{v})}. \quad (6)$$

To maximize the rate at which the mean-squared error, the order of the bias<sup>2</sup> and the variance are usually set to be the same. However, to ensure normality here, we set the right-hand-side to go to zero at a minimal rate  $O(N^{-\eta})$ . This requirement generates the moment condition in (D3). Detailed arguments are shown in Lemmas 3-4. The appendix also includes an illustrative example where assumption (A5a) is violated and conventional asymptotic normality fails.

**D4. Unadjusted Probabilities and Densities.** Let  $\sigma_k$  be the standard deviation for  $V_k$ ,  $k = 1, 2$ , and  $\xi$  a small positive value. For the  $Y_2$ -model, let:

$$\begin{aligned} \hat{P}_2(Y_{2i} = d_2 | V_{2i} = t_2) &\equiv \hat{f}_2(t_2; d_2) / \sum_{d_2=0}^1 \hat{f}_2(t_2; d_2), \\ \hat{f}_2(t_2; d_2) &\equiv \sum_{j=1}^N \frac{Y_{2j}^{d_2} (1 - Y_{2j})^{1-d_2}}{Nh_m} K\left[\frac{t_2 - V_{2j}}{h_m}\right]. \end{aligned}$$

where  $h_m \equiv \sigma_2 N^{-r_m}$ ,  $r_m = \frac{1}{6+\xi}$ .

For the  $Y_1$ -model, conditioned on  $Y_2 = 1$ , let:

$$\begin{aligned} \hat{P}_r(Y_{1i} = d_1 | Y_{2i} = 1, V_i = t) &\equiv \hat{f}(t; d_1) / \sum_{d_1=0}^1 \hat{f}(t; d_1) \\ \hat{f}(t; d_1) &\equiv \sum_{j=1}^N \frac{Y_{1j}^{d_1} (1 - Y_{1j})^{1-d_1} Y_{2j}}{Nh_{c1}h_{c2}} K\left(\frac{t_1 - V_{1j}}{h_{c1}}\right) K\left(\frac{t_2 - V_{2j}}{h_{c2}}\right) \end{aligned}$$

where  $h_{c1} \equiv \sigma_1 h_c$ ,  $h_{c2} \equiv \sigma_2 h_c$ ,  $h_c \equiv N^{-r_c}$ ,  $r_c = \frac{1}{8+\xi}$ .

Let  $g_1$  be the marginal density for  $V_1$  and the corresponding estimate:

$$\hat{g}_1(t_1) = \sum_{j=1}^N \frac{1}{Nh} K\left(\frac{t_1 - V_{1j}}{h_1}\right) \text{ where } h_1 = \sigma_1 N^{-.2}.$$

When the value  $t_k$  is replaced by the observation  $V_{ik}$ , the above averages are taken over the  $(N - 1)$

observations for which  $j \neq i$ .<sup>6</sup>

**D5. Smooth Trimming.** Define a smooth trimming function as:

$$\tau(z, m) \equiv [1 + \exp(Ln(N)[z - m])]^{-1}.$$

**D6. Interior Index Trimming.** Let  $\hat{V}_k^U$  and  $\hat{V}_k^L$  be upper and lower sample quantiles for the indices:  $V_k \equiv V_k(\theta)$ ,  $k = 1, 2$ . Referring to (D5), define smooth interior trimming functions as:

$$\hat{\tau}_I(t_k) \equiv \tau(\hat{V}_k^L, t_k) \tau(t_k, \hat{V}_k^U).$$

**D7. Density Adjustment.** Referring to (D4), let  $\hat{q}_2$  be a lower sample quantile for  $\hat{f}_2(V_2; d_2)$ , and  $\hat{q}$  a lower sample quantile for  $\hat{f}(V; d_1, d_2)$ . Then, define adjusted estimates as:

$$\hat{f}_2^*(t_2; d_2) = \hat{f}_2(t_2; d_2) + \hat{\Delta}_2(d_2), \quad \hat{\Delta}_2(d_2) \equiv N^{-r_2/2} [1 - \hat{\tau}_I(t_2)] \hat{q}_2$$

$$\hat{f}^*(t; d_1, d_2) = \hat{f}(t; d_1, d_2) + \hat{\Delta}(d_1, d_2), \quad \hat{\Delta}(d_1, d_2) \equiv N^{-r/2} [1 - \hat{\tau}_I(t_1) \hat{\tau}_I(t_2)] \hat{q}.$$

**D8. Adjusted Semiparametric Probability Functions.** Let:

$$\begin{aligned} \hat{P}_2^*(Y_{2i} = d_2 | V_{2i} = t_2) &\equiv \hat{f}_2^*(t_2; d_2) / \sum_{d_2=0}^1 \hat{f}_2^*(t_2; d_2) \\ \hat{P}r^*(Y_{1i} = d_1 | Y_{2i} = d_2, V_i = t) &\equiv \hat{f}^*(t; d_1, d_2) / \sum_{d_1=0}^1 \hat{f}^*(t; d_1, d_2) \\ \hat{P}_i^*(d_1, d_2; \theta) &\equiv \hat{P}r^*(Y_i(d_1, d_2) = 1 | V_i(\theta) = v_i(\theta)) \\ &\equiv \hat{P}r^*(Y_{1i} = d_1 | Y_{2i} = d_2, V_i = t) * \hat{P}_2^*(Y_{2i} = d_2 | V_{2i} = t_2). \end{aligned}$$

**D9. Likelihood Trimming.** Define  $\tau_{ix}$  as an indicator that is equal to one if all of the continuous  $X$ 's are between their respective lower and upper quantiles, and define  $\tau_{iv}$  as an indicator that is equal to one if the index vector  $V_{0i}$  is between its lower and upper quantiles.

---

<sup>6</sup>It can easily be shown that all estimators with windows depending on population standard deviations are asymptotically the same as those based on sample standard deviations. For notational simplicity, we employ population standard deviations throughout.

**D10. First and Second Stage Estimators.** We define the first stage estimator as:

$$\hat{\theta} \equiv \arg \max_{\theta} \hat{L}(\theta),$$

$$\hat{L}(\theta) \equiv \sum_{i=1}^N \tau_{ix} \sum_{d_1, d_2} Y_i(d_1, d_2) \text{Ln} \left( \hat{P}_i(d_1, d_2; \theta) \right).$$

Recall that  $\tau_{ix}$  is a trimming function based on  $X$  while  $\tau_{iv}$  is based on the index vector.<sup>7</sup> In the objective function above, replace  $\hat{P}_i$  with  $\hat{P}_i^*$  as defined in (D8), replace  $\tau_{ix}$  with  $\tau_{iv}$ , and term the new objective function as  $\hat{L}^*(\theta)$ . Then, define the second stage estimator:

$$\hat{\theta}^* \equiv \arg \max_{\theta} \hat{L}^*(\theta).$$

**D11. The Adjusted Estimator.** Letting

$$\hat{\delta}_i^*(d_1, d_2; \theta) \equiv \nabla_{\theta} \hat{P}_i^*(d_1, d_2; \theta) / \hat{P}_i^*(d_1, d_2; \theta),$$

define  $\hat{P}^o(d_1, d_2; \theta)$  as an estimated semiparametric probability function. Referring to (D7-8), the components are based on optimal window parameters:  $r = r^o = 1/6$  and  $r_2 = r_2^o = 1/5$ . Then, define a gradient correction as:

$$\hat{C}(\hat{\theta}^*) \equiv \sum_{i=1}^N \tau_{iv}(\hat{\theta}^*) \sum_{d_1, d_2} \left[ \hat{P}_i^*(d_1, d_2; \hat{\theta}^*) - \hat{P}_i^o(d_1, d_2; \hat{\theta}^*) \right] \hat{\delta}_i^*(d_1, d_2; \hat{\theta}^*).$$

With  $\hat{H}(\hat{\theta}^*)$  as the estimated Hessian, the adjusted estimator is defined as:

$$\hat{\theta}^o \equiv \hat{\theta}^* - \hat{H}(\hat{\theta}^*)^{-1} \hat{C}(\hat{\theta}^*).$$

---

<sup>7</sup>Define  $\hat{\tau}_{ix}$  and  $\hat{\tau}_{iv}$  as estimated trimming functions based on sample quantiles of  $X$  and the estimated index respectively. It can be shown using Lemma 2.18 in Pakes and Pollard(1989) that the estimators based on known trimming functions are asymptotically equivalent to those based on the true ones. For expositional simplicity, we take these trimming functions as known throughout.

## 5 Asymptotic Results

### 5.1 Marginal Effects

We now provide the asymptotic results for the marginal effect estimator. We begin with a characterization theorem underlying consistency and normality.

**Theorem 1** Referring to (D1 & D4),

$$\gamma_N(\bar{v}) \equiv \frac{\sum_j \frac{1}{Nh} [Y_{1j} - \zeta_0(\bar{v})] K[(\bar{v} - V_{1j})/h] Y_{2j} S_j}{E(S|V_1 = \bar{v}) g_1(\bar{v})},$$

which depends on true selection probabilities and a known high probability set. Then:

$$C_N(\bar{v}) \left[ \hat{\zeta}(\bar{v}) - \zeta_0(\bar{v}) \right] = C_N(\bar{v}) \gamma_N(\bar{v}) + o_p(1)$$

where  $C_N(\bar{v}) \equiv \frac{\sqrt{Nh} E(S|\bar{v})}{\sqrt{E(S^2|\bar{v})}}$  satisfies  $C_N^2(\bar{v}) = O\left(\frac{1}{\text{Var}(\gamma_N(\bar{v}))}\right)$ .

The consistency and normality results now follow:

**Theorem 2 (Consistency)** Select the high probability set as in (D3) and assume that

$$NhE(S|V_1 = \bar{v}) \rightarrow \infty$$

as  $N$  increases. Then, for the estimator defined in (D1):

$$\hat{\zeta}(\bar{v}) \xrightarrow{p} \zeta_0(\bar{v}).$$

As shown in the Appendix, this result follows from Theorem 1, because the bias and variance of  $\gamma_N(\bar{v})$  both tend to zero as  $N$  increases.

**Theorem 3 (Normality).** Let

$$\begin{aligned} \hat{V} &= \widehat{\text{Var}}(\gamma_N(ve)) + \widehat{\text{Var}}(\gamma_N(vb)) \\ \widehat{\text{Var}}(\gamma_N(\bar{v})) &= \frac{\hat{\zeta}(\bar{v}) \left[ 1 - \hat{\zeta}(\bar{v}) \right] \sum_j \frac{1}{Nh} K^2[(\bar{v} - V_{1j})/h] \hat{S}_j^2}{Nh \hat{E}_2^2(\hat{S}|\bar{v}) \hat{g}_1^2(\bar{v})} \end{aligned}$$

Then

$$\frac{\widehat{ME} - ME}{\sqrt{\widehat{V}}} \xrightarrow{d} Z \sim N(0, 1).$$

To see that this result follows from Theorem 1, we need to show that the covariance between  $\gamma_N(vb)$  and  $\gamma_N(ve)$  tends to 0 as  $N$  increases and we need to establish the relevant Lindberg condition. These results are established in the Appendix.

## 5.2 Index Parameters

To provide an overview of the theoretical arguments, we note that the consistency argument is rather standard except that we need to accommodate the bias controls used in the normality arguments. Hence we start by giving an overview of the normality arguments.

Because indicators and probabilities sum to one over all possible cells, the gradient to the objective function has the form:

$$\hat{G} = \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} \hat{\varepsilon}_i(d_1, d_2) \hat{\delta}_i(d_1, d_2; \theta_0) \tau_i, \quad (7)$$

where  $\hat{\varepsilon}_i(d_1, d_2) \equiv Y_i(d_1, d_2) - \hat{P}_i(d_1, d_2; \theta_0)$ ,  $\hat{\delta}_i(d_1, d_2; \theta_0) \equiv \nabla_{\theta} \hat{P}_i(d_1, d_2; \theta_0) / \hat{P}_i(d_1, d_2; \theta_0)$ , and we have taken the trimming function as known for expositional purposes. As is standard, the key part of the normality argument is to show that the normalized gradient converges to a normal distribution. Denoting  $\varepsilon_i(d_1, d_2) \equiv Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)$ ,  $\delta_i(d_1, d_2; \theta_0) \equiv \nabla_{\theta} P_i(d_1, d_2; \theta_0) / P_i(d_1, d_2; \theta_0)$ , and suppressing the  $(d_1, d_2; \theta_0)$  notation for simplicity, for each cell we may write the normalized gradient as:

$$\frac{1}{\sqrt{N}} \left[ \sum_{i=1}^N \varepsilon_i \delta_i \tau_i + \sum_{i=1}^N \varepsilon_i (\hat{\delta}_i - \delta_i) \tau_i + \sum_{i=1}^N (\hat{\varepsilon}_i^{\circ} - \varepsilon_i) (\hat{\delta}_i - \delta_i) \tau_i + \sum_{i=1}^N (\hat{\varepsilon}_i^{\circ} - \varepsilon_i) \delta_i \tau_i \right]$$

where  $\hat{\varepsilon}_i^{\circ}(d_1, d_2) \equiv Y_i(d_1, d_2) - \hat{P}_i^{\circ}(d_1, d_2; \theta_0)$ . We establish normality by showing that every term vanishes except the first. The second term above readily vanishes from a mean-square convergence argument. For the third term, a Cauchy-Schwartz argument enables us to separate the individual components and take advantage of the known convergence rate of each. With  $\hat{P}_i^{\circ}$  based on an optimal window, the rate of convergence for  $\hat{\varepsilon}_i^{\circ} - \varepsilon_i$  is fast enough so that the third term vanishes.

For the final term, we rely on the result due to Whitney Newey mentioned above to obtain  $E(\delta_i | V_i) = 0$ . Therefore, this multiplicative gradient component can serve as a source of bias reduction. To exploit this

residual-like property of the probability gradient, denote  $V_0$  as the matrix of observations on the indices and define  $D(V_0) \equiv E[(\hat{\varepsilon}_i^o - \varepsilon_i)|X]$ . Then from an iterated expectations argument, conditioning on  $X$ :

$$EE[(\hat{\varepsilon}_i^o - \varepsilon_i)\delta_i\tau_i|X] = E[D(V_0)E(\delta_i\tau_i|V_0)].$$

If the trimming function depended on the index, this gradient component would now have zero expectation. We design a two-stage estimator where parameter estimates from the first stage are used to construct the index and then index trimming is employed in the second stage. We then show that this fourth term is equivalent to a centered U-statistic that converges in probability to zero. To achieve consistency with index trimming, we use the adjusted probabilities in (D7, D8) so that denominators are kept away from zero, while the estimated probability still goes rapidly to the truth.

The remainder of this section provides the main asymptotic results in several theorems. Each theorem will depend on a number of intermediate results, which we state and prove as Lemmas in the Appendix. Theorem 4 below provides consistency and identification results. Theorem 5 provides the normality result.

**Theorem 4 (Consistency).** Assume that each index satisfies the identifying assumptions required for single index models.<sup>8</sup> Then, under (A1-4) and (D5-11):

$$\hat{\theta} \xrightarrow{p} \theta_0, \hat{\theta}^* \xrightarrow{p} \theta_0, \hat{\theta}^o \xrightarrow{p} \theta_0.$$

**Theorem 5 (Normality).** With  $L(\theta)$  as the limiting likelihood of  $\hat{L}^*(\theta)$  defined in (D10) and with  $H$  as its Hessian matrix, define  $H_o \equiv EH(\theta_0)$ . Then, with  $\hat{\theta}^o$  as the estimator defined in (D11) and under (A1-4) and (D5-11):

$$\sqrt{N}[\hat{\theta}^o - \theta_0] \xrightarrow{d} Z \sim N(0, -H_o^{-1}).$$

## 6 Simulation Evidence

We now consider the finite sample performance of the estimator in four different models. These differ according to: i) whether or not the model is threshold-crossing; and ii) whether the continuous variables and errors follow a Normal or Weibull distribution. The first two models we consider have threshold-crossing

---

<sup>8</sup>See, for example, Ichimura (1993) or Klein and Spady (1993).

structures. The first model (TNorm design) has normal errors and is given as:

$$\begin{aligned} Y_1^* &= I \left\{ \sqrt{2} (X_1 + X_3) > \varepsilon \right\} \\ Y_2 &= I \left\{ \sqrt{2} (X_2 - X_3) > v \right\}, \end{aligned}$$

where  $Y_1^*$  is observed when  $Y_2 = 1$ . The errors and the continuous  $X$ 's ( $X_1, X_2$ ) are generated as:

$$\begin{aligned} v, X_2 &\sim N(0, 1) \\ \varepsilon &= 2v + z, \quad z \sim N(0, 1) \\ X_1 &= X_2 + 2z_1, \quad z_1 \sim N(0, 1). \end{aligned}$$

and re-scaled to each have variance 1, while  $X_3$  is a binary variable with probability .5 and support  $\{-1, 1\}$ . Notice that the indices have standard deviation 2. For the second index, this ensures that the index has fatter tails than the error, which is theoretically needed in estimating the marginal effect.

In a second model (TWeibull design), the selection error is non-normal while the model structure stays the same. The error  $v$  follows a Weibull (1,1.5) giving a right tail probability of  $\exp(-v^{1.5})$ . We set  $X_2$  to follow Weibull (1,1) so that the tail comparison condition is satisfied. As above, all the variables and errors are rescaled to have zero mean and variance one.

In the third (NTNorm design) and fourth (NTWeibull design) models, the  $Y_1^*$  equation has a non-threshold-crossing structure:

$$Y_1^* = I \left\{ X_1 + X_3 > s \left[ 1 + (X_1 + X_3)^2 / 4 \right] \varepsilon \right\}$$

where the variables are generated as in the previous models. Note that  $s$  is chosen to ensure the right-hand-side of the inequality is rescaled to have variance one as above. Similar to the first two models above, here the third and fourth models differ according to whether Normal or Weibull distributions are employed.

For all models, we set  $N = 2000$  and conduct 1000 replications. We compare the finite sample performance of our semiparametric marginal effect estimator and the bivariate probit with selection counterpart. We also compare the parameter estimates upon which these marginal effects are based. Finally, we provide results for the estimation of the high probability set. Results for the marginal effects are shown in Table



1. Notice that there are an infinite number of marginal effects because there are an infinite number of base levels and evaluation levels. Here we report the marginal effect of moving  $X_1$  from its median level to one unit above while keeping the binary variable  $X_3$  at zero. Overall, the semiparametric estimator performs well over designs with a small bias and standard deviation. In contrast, the bivariate probit counterpart does not perform well outside of the TNorm design where bivariate probit is correct. In the TNorm case, where bivariate probit is the correct specification, it does indeed have a small bias and standard deviation. However, the advantage over the semiparametric marginal effect is minimal. The RMSE of bivariate probit is .06 compared with .07 from the semiparametric counterpart. In the TWeibull case, the semiparametric method shows significant advantage in terms of the bias. The bias of the semiparametric marginal effect estimator is almost zero, while the bivariate probit counterpart has a bias of .08, which is almost 30% of the truth (.29). When we move on to the non-threshold-crossing designs, we continue to see the semiparametric estimator performing significantly better. In the NTNORM case, the semiparametric has both smaller bias (.02 vs .10) and smaller standard deviation (.05 vs .07). In the NTWeibull case, the semiparametric estimator still performs much better than the bivariate probit in terms of RMSE (.08 vs .21). Most of the advantage comes from the standard deviation (.06 vs .20).

We also provide the index parameter estimation results in Table 2. For semiparametric estimation, the parameters are identified up to location and scale, hence we report  $\text{Ratio}_{31} = \frac{\text{coef}(X_3)}{\text{coef}(X_1)}$  in the outcome equation and  $\text{Ratio}_{32} = \frac{\text{coef}(X_3)}{\text{coef}(X_2)}$  in the selection equation. Notice that for the non-threshold-crossing designs, we report the median and median absolute deviation (MAD) for the bivariate probit estimators because there were a number of replications where bivariate probit performed extremely poorly. The semiparametric estimator, however, does not have this issue, hence we report not only median and MAD but also mean, standard deviation, and RMSE. For the selection equation, over all designs, both parametric and semiparametric estimators perform quite well. Turning to the outcome equation, both estimators perform better in normal than in non-normal designs and also better in threshold-crossing than in non-threshold-crossing designs. The non-threshold-crossing model with Weibull distributions poses the most challenge for both estimators. It is noteworthy that for all other designs, the bias and the standard deviation for the semiparametric estimator are quite small. Finally, we also investigated using higher order kernels for estimating index parameters as an alternative to the bias controls implemented here.<sup>9</sup> Due to convergence problems, we found it necessary

---

<sup>9</sup>In our Monte Carlo studies, the higher order kernel we use is the twicing kernel for both index parameter estimation and estimation of the high probability set parameter.

to examine this estimator on a two-dimensional grid, which was quite time-consuming. Accordingly, we only examined 100 replications for each design (at which point the estimator seemed quite stable). For the selection equation, the RMSE's were close with the exception of the TWeibull design where the RMSE using higher order kernels was 2.5 times larger. For the outcome equation, in all designs the RMSE under higher order kernels was of the order of 3 times larger.

Lastly, we provide the estimation results for the high probability set parameters. The means of  $\hat{a}$  with standard deviations in parentheses are as follows: .31(.004), 28(.006), 31(.004), and .28(.005) for TNorm, TWeibull, NTNORM, NTWeibull respectively. While the variances for all of the estimates are quite small, it is difficult to evaluate the performance the estimator without knowing  $a_{0N}$ . Accordingly, we examined the performance of the estimator for the Weibull example given prior to (D1) with  $c_v = 1$ . It can be shown that in that example, by setting  $b = 0$  the moment condition is approximately equivalent to:

$$\left[2 + (a_{0N} \ln N)^{\frac{1}{c_u}-1}\right] a_{0N} = .8 - (\varepsilon - \eta).$$

Since  $a_{0N}$  depends on the sample size, we examined three different sample sizes:  $N = 500, 1000,$  and  $2000$ . At each of these sample sizes, we solved the above equation for  $a_{0N}$  and conducted a Monte-Carlo experiment with 100 replications to evaluate the performance of  $\hat{a}$  at the base level of the index. The results are as follows:

<i>SAMPLE SIZE</i>	$a_{0N}$	$ BIAS $	$SD$	$RMSE$
500	.279	.037	.027	.046
1000	.280	.025	.025	.035
2000	.283	.019	.009	.020

where bias, standard deviation(SD) and RMSE are standardized by the truth  $a_{0N}$ . This table shows that our  $\hat{a}$  performs very well in terms of absolute bias, standard deviation and RMSE. It also confirms that the absolute bias, standard deviation and RMSE all decline as the sample size increases. As expected,  $a_{0N}$  increases slowly with the sample size.

## 7 Conclusions

This paper studies the binary outcome model with sample selection in a semiparametric framework. As marginal effects are often of primary interest in this type of model, we propose a semiparametric marginal effect estimator. This marginal effect estimator is based on observations in a high probability set where the selection probabilities are above a cutoff. We propose an estimator for this cutoff and establish its large sample properties. Based on that, we establish the large sample properties for our marginal effect estimator, which takes into account that the cutoff and the selection probability are estimated. In a Monte-Carlo study we find that our marginal effect estimator based on the estimated high probability set performs quite well in finite samples.

This marginal effect estimator is developed under an index framework so as to achieve good performance in finite samples. Accordingly, it depends on an estimator for index parameters. In this paper, we propose an index parameter estimator based on regular kernels with bias control mechanisms and show that the estimator is consistent and asymptotically distributed as normal. While retaining these desirable large sample properties, the Monte-Carlo results show that this estimator performs very well in finite samples.

<b>Marginal Effect Estimators</b>				
	Truth		Bivariate Probit	Semiparametric
TNorm	.34	mean	.33	.31
		std	.06	.06
		RMSE	.06	.07
TWeibull	.29	mean	.37	.29
		std	.06	.06
		RMSE	.10	.06
NTNorm	.46	mean	.36	.48
		std	.07	.05
		RMSE	.13	.05
NTWeibull	.59	mean	.63	.63
		std	.20	.06
		RMSE	.21	.08

<b>Index Parameters</b>						
	Bivariate Probit				Semiparametric	
	Outcome		Selection		Outcome	Selection
	Coef(X1)	Coef(X3)	Coef(X2)	Coef(X3)	Ratio <sub>31</sub>	Ratio <sub>32</sub>
<b>TNorm</b>						
mean	.98	1.02	1.01	-1.01	.95	-1.04
std	.05	.23	.06	.04	.07	.05
RMSE	.05	.22	.06	.04	.08	.06
<b>TWeibull</b>						
mean	1.23	1.23	1.02	-1.06	.93	-1.04
std	.09	.21	.06	.04	.06	.04
RMSE	.25	.31	.06	.08	.10	.06
<b>NTNorm</b>						
mean					.96	-1.02
median	1.07	1.13	1.00	-1.00	.96	-1.01
std					.05	.04
MAD	.12	.15	.03	.03	.04	.03
RMSE					.06	.04
<b>NTWeibull</b>						
mean					.84	-1.04
median	2.30	2.52	1.03	-1.07	.84	-1.04
std					.04	.04
MAD	1.30	1.52	.05	.07	.16	.05
RMSE					.16	.05

## References

- [1] Andrews, D. and M.Schafgans (1998): "Semiparametric Estimation of the Intercept of a Sample Selection Model" *Review of Economic Studies*, 65, 497-517.
- [2] Chesher, A. (2005): "Nonparametric Identification under Discrete Variation", *Econometrica*, 73, 1525-1550.
- [3] Escanciano, J. C., D. T. Jacho-Chavez and A. Lewbel (2012): "Uniform Convergence of Weighted Sums of Non- and Semi-parametric Residuals for Estimation and Testing," working paper.
- [4] Gallant, A. and D. Nychka (1987): "Semi-nonparametric Maximum Likelihood Estimation," *Econometrica*, 15, 363-390.
- [5] Heckman, J. (1974): "Shadow Prices, Market Wages and Labor Supply," *Econometrica*, 42(4), 679-94.
- [6] Heckman, J. (1979): "Sample Selection Bias as a Specification Error," *Econometrica*, 47(1), 153-61.
- [7] Heckman, J. (1990): "Varieties of Selection Bias," *American Economic Review*, 80, 313-18.
- [8] Honore, B. E. and J. L. Powell (2005): "Pairwise Difference Estimation of Nonlinear Models." *D. W. K. Andrews and J. H. Stock, eds., Identification and Inference in Econometric Models. Essays in Honor of Thomas Rothenberg* (Cambridge: Cambridge University Press), 520–53.
- [9] Ichimura, H. (1993): "Semiparametric least squares (SLS) and weighted SLS estimation of single index models," *Journal of Econometrics*, 58, 71-120.
- [10] Ichimura, H., and L. F. Lee (1991): "Semiparametric least squares (SLS) and weighted SLS estimation of multiple index models: Single equation estimation," in *Nonparametric and Semiparametric Methods in Econometrics and Statistics*, ed. W.Barnett, J.Powell and G.Tauchen, Cambridge University Press.
- [11] Khan, S. and E.Tamer (2010): "Irregular Identification, Support Conditions, and Inverse Weight Estimation," *Econometrica*, 6, 2021-2042.
- [12] Klein, R. and C. Shen (2010): "Bias Corrections in Testing and Estimating Semiparametric, Single Index Models," *Econometric Theory*, 1683-1718.

- [13] Klein, R. and R. Spady (1993): "An Efficient Semiparametric Estimator for the Binary Response Model," *Econometrica*, 61, 387-421.
- [14] Klein, R. and F.Vella (2009): "A Semiparametric Model for Binary Response and Continuous Outcomes Under Index Heteroscedasticity," *Journal of Applied Econometrics*, 24, 735-762.
- [15] Lee, L.F (1995): "Semi-Parametric Estimation of Polychotomous and Sequential Choice Models", *Journal of Econometrics*, 65, 381-428.
- [16] Newey, W., F. Hsieh, and J. Robins (2004): "Twicing Kernels and a Small Bias Property of Semiparametric Estimators," *Econometrica*, 72, 947-962.
- [17] Newey, W, (2007): "Nonparametric continuous/discrete choice models", *International Economic Review*, 48: 1429–1439.
- [18] Pakes, A., and D.Pollard (1989): "Simulation and the Asymptotics of Optimization Estimators," *Econometrica*, 57, 1027-1058.
- [19] Shaikh, A. M. and Vytlacil, E. J. (2011): "Partial Identification in Triangular Systems of Equations With Binary Dependent Variables," *Econometrica* 79(3), 949–955.
- [20] Vella, F. (1998): "Estimating Models with Sample Selection Bias: A Survey." *The Journal of Human Resources*, 33:1, 127-169.
- [21] Vytlacil, E. and N. Yildiz (2007): "Dummy Endogenous Variables in Weakly Separable Models," *Econometrica*, 75, 757-779.

## 8 Appendix

### 8.1 Illustrative Example

Let  $F$  and  $G$  be the following Weibull distribution functions for the selection error and index respectively:

$$F(u) = 1 - \exp(-u^{c_u}); \quad G(v) = 1 - \exp(-v^{c_v}).$$

Suppose that  $c_v > c_u$  so that error tails are fatter than index tails, which violates assumption (A5a). We show below that conventional asymptotic normality does not hold in this case.

Assuming for simplicity that selection and outcome indices are independent, from (6):

$$O\left(\frac{\text{bias}^2}{\text{var}}\right) \leq \frac{O(hN^{1-2a}) [E(S)]^2}{E(S^2|\bar{v})}$$

Referring to (D2), let  $S^+$  and  $S^-$  be the functions that have binary support where  $S^+$  is one when in R2,R3, while  $S^-$  is zero when in R1,R2. We can bound the S-ratio above by:

$$\frac{[E(S^-)]^2}{E(S^+)} \leq \frac{[E(S)]^2}{E(S^2)} \leq \frac{[E(S^+)]^2}{E(S^-)}.$$

Since  $S^+$  and  $S^-$  are binary variables we have

$$\frac{[\Pr(x > b)]^2}{\Pr(x > 0)} \leq \frac{[E(S)]^2}{E(S^2)} \leq \frac{[\Pr(x > 0)]^2}{\Pr(x > b)}.$$

The probabilities that comprise the above bounds have the form:

$$\Pr(x > 0) = \Pr(P_2 > 1 - N^{-a}) = \Pr(1 - \exp(-v^{c_u}) > 1 - N^{-a}) = \exp\left(- (a \ln N)^{c_v/c_u}\right)$$

$$\Pr(x > b) = \Pr(P_2 > 1 - \exp(-b)N^{-a}) = \exp\left(- (a \ln N + b)^{c_v/c_u}\right).$$

It now follows that:

$$\frac{O(hN^{1-2a})}{\exp[-(a \ln N)^{c_v/c_u} + 2(a \ln N + b)^{c_v/c_u}]} \leq \frac{O(hN^{1-2a}) [E(S)]^2}{E(S^2|\bar{v})} \leq \frac{O(hN^{1-2a})}{\exp[-(a \ln N + b)^{c_v/c_u} + 2(a \ln N)^{c_v/c_u}]}$$

Write the lower bound as:

$$\frac{O(hN^{1-2a})}{\exp[(a \ln N)^{c_v/c_u} \left[ -1 + 2 \left( 1 + \frac{b}{a \ln N} \right)^{c_v/c_u} \right]]}.$$

If this estimator is to be consistent, then  $a \ln N$  must go to infinity. For  $N$  sufficiently large, there then exists fixed  $\varrho > 0$  such that  $\left( 1 + \frac{b}{a \ln N} \right)^{c_v/c_u} < 1 + \varrho$ , because the left hand side converges to 1. A similar argument holds for the upper bound, therefore:

$$\frac{O(hN^{1-2a})}{\exp[(a \ln N)^{c_v/c_u} (1 + 2\varrho)]} \leq \frac{O(hN^{1-2a}) [E(S)]^2}{E(S^2|\bar{v})} \leq \frac{O(hN^{1-2a})}{\exp[(a \ln N)^{c_v/c_u} (1 - \varrho)]}.$$

Notice that the lower and upper bounds have the same order. Referring to (D3), equating this order to  $O(N^{-\eta})$  and taking logs, the order of  $a_{0N}$  is given by solving:

$$-\eta \ln(N) = \ln(hN^{1-2a_{0N}}) - (a_{0N} \ln N)^{c_v/c_u} \Rightarrow a_{0N}^{c_v/c_u} (\ln N)^{(c_v/c_u)-1} = (.8 - \varepsilon + \eta - 2a_{0N}).$$

It follows that  $a_{0N}$  must converge to zero and have order  $O((\ln N)^{c_u/c_v-1})$ ,  $c_v > c_u$ . It can now be shown that the bias and variance converge slowly to zero, but at rates such that there does not exist any fixed  $\kappa > 0$  such that  $N^\kappa [\widehat{ME} - ME]$  has an asymptotic distribution.

## 8.2 Main Results

### 8.2.1 Marginal Effects

**Proof of Theorem 1.** By definition,

$$C_N(\bar{v}) \left[ \hat{\zeta}(\bar{v}) - \zeta_0(\bar{v}) \right] = C_N(\bar{v}) \frac{\sum_j \frac{1}{Nh} [Y_{1j} - \zeta_0(\bar{v})] Y_{2j} K[(\bar{v} - V_{1j})/h] \hat{S}_j}{\sum_j \frac{1}{Nh} Y_{2j} K[(\bar{v} - V_{1j})/h] \hat{S}_j}.$$

Lemma 8 enables us to replace (up to  $o_p(1)$ ) the denominator with  $E(S|V_1 = \bar{v}) g_1(\bar{v})$ , while Lemma 9 continues to show that the numerator has the desired form.

**Proof of Theorem 2.** By Theorem 1:

$$C_N(\bar{v}) \left[ \hat{\zeta}(\bar{v}) - \zeta_0(\bar{v}) \right] = C_N(\bar{v}) \gamma_N(\bar{v}) + o_p(1).$$

Lemma 3 characterizes the order of the bias of the estimator. Recalling the definition of the high probability



parameter in (D3), the bias in the estimator vanishes. From Lemma 4, the reciprocal of the estimator variance has the order:

$$NhE(S|\bar{v})^2/E(S^2|\bar{v}) > NhE(S|\bar{v})$$

which completes the proof as  $NhE(S|\bar{v})$  tends to  $\infty$  as  $N$  increases.

**Proof of Theorem 3.** By definition,  $\widehat{ME} - ME = [\hat{\zeta}(ve) - \zeta_0(ve)] - [\hat{\zeta}(vb) - \zeta_0(vb)]$ . We begin by showing that the covariance between these two components vanishes. Notice that  $\hat{\zeta}(\bar{v}) - \zeta_0(\bar{v})$  is close to  $\gamma_N(\bar{v})$  which we can write as a sample average  $\sum_j \frac{1}{N} t_j(\bar{v})$ . The covariance is then of the form  $E[t_j(ve) t_k(vb)]$ . For  $j \neq k$ , from independence and the vanishing bias of the expectation of each term, this expectation vanishes. For  $j = k$ , the kernel function ensures that this expectation also vanishes faster than  $N^{-1/2}$  as  $V_{1j}$  cannot be close to both  $ve$  and  $vb$ . Therefore, we can calculate the variance as the sum of the variances of  $\gamma_N(ve)$  and  $\gamma_N(vb)$ .

To prove normality, we next establish a Lindberg condition for  $C_N(\bar{v}) [\hat{\zeta}(\bar{v}) - \zeta_0(\bar{v})]$ . Namely, for  $\varepsilon > 0$ , we must show that the following expectation converges to 0:

$$E \left\{ \frac{(Y_{1i} - \theta_o)^2 Y_{2i} k_i^2 S_i^2 / h}{E(S^2|\bar{v})} 1_{\{(Y_{1i} - \theta_o)^2 Y_{2i} k_i^2 S_i^2 > Nh^2 E(S^2|\bar{v})\}} \right\}.$$

Since  $(Y_{1i} - \theta_o)^2 Y_{2i} k_i^2 S_i^2$  is bounded, it suffices to show that  $Nh^2 E(S^2|v) \rightarrow \infty$ . Recalling the definition of  $S^-$  in Section 8.1:

$$Nh^2 E(S^2|v) > Nh^2 E(S^-|\bar{v})$$

$$\begin{aligned} \text{where } E(S^-|\bar{v}) &= \Pr(F(V_1) > 1 - N^{-a_{0N}} / \exp(b)), \\ &= 1 - G_c(F^{-1}(1 - N^{-a_{0N}} / \exp(b))) \\ &> 1 - F(F^{-1}(1 - N^{-a_{0N}} / \exp(b))) \text{ from (A5a)}. \end{aligned}$$

Since  $h = O(N^{-.2-\varepsilon})$  and  $0 < a_{0N} < .4$ , the normality of  $C_N(\bar{v}) [\hat{\zeta}(\bar{v}) - \zeta_0(\bar{v})]$  follows.

Turning to the marginal effects, for expositional purposes, suppose  $O(\text{Var}(\gamma_N(ve))) > O(\text{Var}(\gamma_N(vb)))$ ,

then with  $V \equiv \text{Var}(\gamma_N(v_e)) + \text{Var}(\gamma_N(v_b))$

$$\begin{aligned} \frac{1}{\sqrt{V}} &= O\left(\frac{1}{\sqrt{\text{Var}(\gamma_N(v_e))}}\right) \\ &= O(C_N(v_e)). \end{aligned}$$

Therefore, the characterization results in Theorem 1 apply to yield:

$$\frac{\widehat{ME} - ME}{\sqrt{V}} = O(C_N(v_e)) \left[ \hat{\zeta}(v_e) - \zeta_0(v_e) \right] + o_p(1)$$

Now asymptotic normality follows from the above Lindberg condition. A symmetric argument holds for the case where  $O(\text{Var}(\gamma_N(v_e))) < O(\text{Var}(\gamma_N(v_b)))$ . For the case where  $O(\text{Var}(\gamma_N(v_e))) = O(\text{Var}(\gamma_N(v_b)))$ , a Lindberg condition similar to the above applies. Therefore,  $\frac{\widehat{ME} - ME}{\sqrt{V}} \xrightarrow{d} Z \sim N(0, 1)$ . Employing similar arguments as in Lemma 8, it can be shown that  $\frac{V - \hat{V}}{V} \xrightarrow{p} 0$ . Hence the theorem follows. .

### 8.2.2 Index Parameters

**Proof of Theorem 4.** We provide the proof for  $\hat{\theta}^*$ , with the arguments for the other estimators being very similar. Lemma 10 proves that we can replace the  $\hat{P}_i^*$  with  $P_i^*$  in the objective function  $\hat{L}^*(\theta)$ , and obtain  $L^*(\theta)$  satisfying:

$$\sup_{\theta} \left| \hat{L}^*(\theta) - L^*(\theta) \right| \xrightarrow{p} 0.$$

From Lemma 11, we may ignore the probability adjustments  $\hat{\Delta}'s$  and therefore replace adjusted probabilities  $P_i^*$  in  $L^*(\theta)$  with unadjusted ones  $P_i$ . With  $L(\theta)$  as the resulting objective function:

$$\sup_{\theta} |L^*(\theta) - L(\theta)| \xrightarrow{p} 0.$$

From conventional uniform convergence arguments:

$$\sup_{\theta} |L(\theta) - E[L(\theta)]| \xrightarrow{p} 0.$$

To complete the argument, we must show that  $E[L(\theta)]$  is uniquely maximized at  $\theta_0$ . From standard arguments,  $\theta_0$  is a maximum, and the only issue is one of uniqueness. With  $\theta^*$  as any potential maximizer,

it can be shown that any candidate for a maximum must give correct probabilities for all three cells:  $(Y_1 = 1, Y_2 = 1)$ ,  $(Y_1 = 0, Y_2 = 1)$ , and  $Y_2 = 0$ . It then follows that for the  $Y_2 = 0$  cell:

$$\Pr(Y_2 = 0|V_2(\theta_2^*)) = \Pr(Y_2 = 0|X) = \Pr(Y_2 = 0|V_2(\theta_{20})).$$

Under identifying conditions for single index models,  $\theta_2^* = \theta_{20}$ . For the  $(Y_1 = 1, Y_2 = 1)$  cell:

$$\begin{aligned} \Pr(Y_1 = 1|Y_2 = 1, V_1(\theta_1^*), V_2(\theta_2^*)) \Pr(Y_2 = 1|V_2(\theta_2^*)) &= \\ \Pr(Y_1 = 1|Y_2 = 1, V_1(\theta_{10}), V_2(\theta_{20})) \Pr(Y_2 = 1|V_2(\theta_{20})). & \end{aligned}$$

Since  $\theta_2^* = \theta_{20}$ :

$$\Pr(Y_1 = 1|Y_2 = 1, V_1(\theta_{10}), V_2(\theta_{20})) = \Pr(Y_1 = 1|Y_2 = 1, V_1(\theta_1^*), V_2(\theta_{20})).$$

Solving the first probability function for  $V_1(\theta_{10})$ , for some function  $\Upsilon$  we have:

$$V_1(\theta_{10}) = \Upsilon(V_1(\theta_1^*), V_2(\theta_{20})).$$

Since  $V_2$  contains a continuous variable that does not affect  $V_1$ , differentiating both sides with respect to this variable yields  $0 = \nabla_{v_2} \Upsilon$ . Therefore,  $\Upsilon$  must only be a function of the first index and is equal to  $V_1(\theta_{10})$ . Identification now follows from conditions that identify single index models.

**Proof of Theorem 5.** From a Taylor expansion, the unadjusted estimator has the form:

$$\left(\hat{\theta}^* - \theta_0\right) = -\hat{H}(\theta^+)^{-1} \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} \left[ Y_i(d_1, d_2) - \hat{P}_i^*(d_1, d_2; \theta_0) \right] \hat{\delta}_i^*(d_1, d_2; \theta_0) \tau_{iv}.$$

where  $\theta^+$  is an intermediate point. To simplify the adjustment to this estimator, referring to (D11) we will show below:

$$\Delta \equiv \hat{H}(\hat{\theta}^*)^{-1} \hat{C}(\hat{\theta}^*) - \hat{H}(\theta^+)^{-1} \hat{C}(\theta_0) = o_p(N^{-1/2}).$$

Rewriting the above expression,  $\Delta = \Delta_1 + \Delta_2$ , where:

$$\begin{aligned}\Delta_1 &\equiv \hat{H}(\theta^+)^{-1} \hat{H}(\hat{\theta}^*)^{-1} \left[ \hat{H}(\theta^+) - \hat{H}(\hat{\theta}^*) \right] \hat{C}(\hat{\theta}^*) \\ \Delta_2 &\equiv \hat{H}(\theta^+)^{-1} \left[ \hat{C}(\hat{\theta}^*) - \hat{C}(\theta_0) \right].\end{aligned}$$

To study  $\Delta_1$ , note that lemma 15 gives a convergence rate for  $\hat{\theta}^* - \theta^+$ . Then, using a Taylor series expansion on  $\left[ \hat{H}(\theta^+) - \hat{H}(\hat{\theta}^*) \right]$  and Lemmas 1, 14 and 15, it can be shown that  $\left[ \hat{H}(\theta^+) - \hat{H}(\hat{\theta}^*) \right]$  and  $\hat{C}(\hat{\theta}^*)$  converge to zero sufficiently fast that  $\Delta_1 = o_p(N^{-1/2})$ .

For  $\Delta_2$ , Taylor expanding the second component:

$$\Delta_2 = \hat{H}(\theta^+)^{-1} \nabla \hat{C}(\hat{\theta}^* - \theta_0),$$

where  $\nabla \hat{C}$  is evaluated at an intermediate point. The first component is  $O_p(1)$  from Lemma 1; the second component is  $O_p(\frac{1}{\sqrt{Nh^3}})$  from Lemma 1; and the third component is  $O_p(h^4)$ , hence we have  $\Delta_2 = o_p(N^{-1/2})$ .

From the definition of  $\hat{\theta}^o$  in (D11) and employing the result above:

$$\begin{aligned}\sqrt{N}(\hat{\theta}^o - \theta_0) &= \sqrt{N}(\hat{\theta}^o - \theta_0 + \Delta) + o_p(1) = -\hat{H}(\theta^+)^{-1} \sqrt{N}(\hat{A}^* - \hat{B}^o) + o_p(1), \\ \hat{A}^* &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)] \hat{\delta}_i^*(d_1, d_2; \theta_0) \tau_{iv}, \\ \hat{B}^o &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [\hat{P}_i^o(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0)] \hat{\delta}_i^*(d_1, d_2; \theta_0) \tau_{iv}.\end{aligned}$$

From Lemma 12:

$$\hat{A}^* = \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)] \delta_i(d_1, d_2; \theta_0) \tau_{iv} + o_p(N^{-1/2}),$$

where  $\delta_i(d_1, d_2; \theta_0)$  is the probability limit of  $\hat{\delta}_i^*(d_1, d_2; \theta_0)$ . It can also be shown that:

$$\hat{B}^o = \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [\hat{P}_i^o(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0)] \delta_i(d_1, d_2; \theta_0) \tau_{iv} + o_p(N^{-1/2}).$$

Lemma 14b shows that  $\hat{B}^o$  is of order  $o_p(N^{-1/2})$ . The theorem now follows.

### 8.3 Intermediate Lemmas

This section provides three types of lemmas: 1) basic lemmas required by all estimators, 2) lemmas required to analyze the marginal effects estimator, and finally 3) lemmas relevant for the index estimator.

#### 8.3.1 Basic Lemmas

For  $\alpha > 0$ , referring to (D4), with  $k = 1, 2$ , define:

$$\mathcal{V}_{2N} = \{v_2 : a_2(d_2) + h_m^{1-\alpha} < v_2 < b_2(d_2) - h_m^{1-\alpha}\} \quad (8)$$

$$\mathcal{V}_N = \{v : a_{ck}(d_k) + h_c^{1-\alpha} < v_k < b_{ck}(d_k) - h_c^{1-\alpha}\}. \quad (9)$$

If the conditional density for  $V_2$ ,  $g_2(v_2|Y_2 = d_2)$  has compact support, interpret  $[a_2(d_2), b_2(d_2)]$  as that support. Similarly, interpret  $[a_{ck}(d_k), b_{ck}(d_k)]$  as the support for conditional density for  $V$ ,  $g(v|Y_1 = d_1, Y_2 = d_2)$  when it is compact. When these supports are unbounded, we let  $\mathcal{V}_{2N}$  and  $\mathcal{V}_N$  approach all of  $R^1$  and  $R^2$  respectively as  $N$  increases.

We begin with two basic lemmas on uniform and pointwise convergence rates. As the proofs of these lemmas are standard in the literature, they are not provided here but are available upon request.

**Lemma 1 (Uniform Convergence).** For  $\psi$  any  $p^{th}$  differentiable function of  $\theta$ , let  $\nabla_\theta^p(\psi)$  be the  $p^{th}$  partial derivative of  $\psi$  with respect to  $\theta$ ,  $\nabla_\theta^0(\psi) \equiv \psi$ . Let  $\hat{f}_2$  and  $\hat{f}$  be the estimators in (D4) with respective probability limits  $f_2$  and  $f$ . Then, for  $\theta$  in a compact set,  $t_2 \in \mathcal{V}_{2N}$  as defined in (8),  $t \in \mathcal{V}_N$  as defined in (9), the following rates hold for  $p = 0, 1, 2$ :

$$\begin{aligned} a) & : \sup_{t_2, \theta} \left| \nabla_\theta^p(\hat{f}_2(t_2; d_2)) - \nabla_\theta^p(f_2(t_2; d_2)) \right| = O_p \left( \min \left[ h_m^2, \frac{1}{\sqrt{N}h_m^{p+1}} \right] \right) \\ b) & : \sup_{t, \theta} \left| \nabla_\theta^p(\hat{f}(t; d_1, d_2)) - \nabla_\theta^p(f(t; d_1, d_2)) \right| = O_p \left( \min \left[ h_c^2, \frac{1}{\sqrt{N}h_c^{p+2}} \right] \right). \end{aligned}$$

**Lemma 2 (Pointwise Convergence).** Using the same notation as above in Lemma 1:

$$\begin{aligned} a) & : \left| \nabla_{\theta}^p \left( \hat{f}_2(t_2; d_2) \right) - \nabla_{\theta}^p \left( f_2(t_2; d_2) \right) \right| = O_p \left( \min \left[ h_m^2, \frac{1}{\sqrt{N h_m^{2p+1}}} \right] \right) \\ b) & : \left| \nabla_{\theta}^p \left( \hat{f}(t; d_1, d_2) \right) - \nabla_{\theta}^p \left( f(t; d_1, d_2) \right) \right| = O_p \left( \min \left[ h_c^2, \frac{1}{\sqrt{N h_c^{2p+2}}} \right] \right). \end{aligned}$$

### 8.3.2 Marginal Effects Lemmas

**Lemma 3.** Under (A4,A6,A7), with  $\zeta_0(\bar{v}) \equiv \Pr(Y_1 = 1|V_1 = \bar{v})$  and  $\gamma_N \equiv \frac{\sum_j \frac{1}{N h} [Y_{1j} - \zeta_0(\bar{v})] K[(\bar{v} - V_{1j})/h] Y_{2j} S_j}{E(S|V_1 = \bar{v}) g_1(\bar{v})}$ ,

$$|E(\gamma_N)| \leq B_N = O(N^{-a} E(S)/E(S|\bar{v})).$$

**Proof.** With  $P_2 = \Pr(Y_2 = 1|V_2)$  and  $\mu_d(V_1, V_2) \equiv E[Y_1 - \zeta_0(\bar{v}) | Y_2 = d, V_1, V_2]$ , and  $\gamma_{1N}$  as the numerator of  $\gamma_N$ :

$$\begin{aligned} E(\gamma_{1N}) & = E\left(\frac{1}{h} \mu_1(V_1, V_2) K[(\bar{v} - V_1)/h] S\right) P_2 \\ & = \iint \mu_1(\bar{v} + hz, v_2) K(z) S P_2 g(\bar{v} + hz, v_2) dz dv_2. \end{aligned}$$

Using a Taylor series expansion,

$$|E(\gamma_{1N})| \leq \left| \int \mu_1(\bar{v}, v_2) P_2 S g(\bar{v}, v_2) dv_2 \right| + |RES|.$$

Note that  $\mu_1(\bar{v}, v_2) P_2 + \mu_0(\bar{v}, v_2) (1 - P_2) = E[Y_1 - \zeta_0(\bar{v}) | V_1 = \bar{v}, V_2] = 0$ , hence for the first term on the right-hand-side:

$$\begin{aligned} \left| \int \mu_1(\bar{v}, v_2) P_2 S g(\bar{v}, v_2) dv_2 \right| & = \left| \int \mu_0(\bar{v}, v_2) (1 - P_2) S g(\bar{v}, v_2) dv_2 \right| \\ & \leq O \left( N^{-a} \int S g(\bar{v}, v_2) dv_2 \right) \\ & = O \left( N^{-a} g(\bar{v}) \int S g(v_2|\bar{v}) dv_2 \right) \\ & = O(N^{-a} E(S|\bar{v})). \end{aligned}$$

The second term on the right-hand-side ( $|RES|$ ) is a residual term from the Taylor series expansion, which

is  $O(h^2 E(S))$ . Therefore, combining those two terms, the slowest rate would be  $|E(\gamma_{1N})| = O(N^{-a} E(S))$  since  $O(h^2) < O(N^{-a})$  and  $O(E(S|\bar{v})) \leq O(E(S))$  from (A6).

**Lemma 4.** For  $\gamma_N$  defined in Lemma 3 and  $a_{0N}$  defined in (D3),

$$\frac{1}{\sqrt{\text{Var}(\gamma_N)}} = O\left(\frac{\sqrt{Nh}E(S|\bar{v})}{\sqrt{E(S^2|\bar{v})}}\right).$$

**Proof.** For  $a_{0N}$  set as in (D3),  $\frac{(E(\gamma_N))^2}{\text{Var}(\gamma_N)} \rightarrow 0$ ; hence

$$\begin{aligned} \text{Var}(\gamma_N) &= O\left(\frac{E([Y_1 - \zeta_0(\bar{v})]^2 K^2 [(\bar{v} - V_1)/h] Y_2 S^2)}{Nh^2 (E(S|V_1 = \bar{v}))^2 g^2(\bar{v})}\right) \\ &= O\left(\frac{E(K^2 [(\bar{v} - V_1)/h] S^2)}{Nh^2 (E(S|V_1 = \bar{v}))^2}\right). \end{aligned}$$

Letting  $z = (V_1 - \bar{v})/h$ , the result follows from a Taylor series expansion about  $h = 0$ .

Lemma 5 below provides a result needed to obtain the convergence rate of  $\hat{a} - a_{0N}$ .

**Lemma 5.** Let

$$c_N(a) \equiv N^{-2(a-.4)-\varepsilon+\eta}$$

$$M_1(a) \equiv E(S)^2$$

$$M_2(a) \equiv E(S^2|\bar{v})$$

and recall (D3), then:

$$c_N(\hat{a}) \left( \frac{[\hat{E}_2(\hat{S})]^2}{\hat{E}_2(\hat{S}^2|\bar{v})} - \frac{M_1(\hat{a})}{M_2(\hat{a})} \right) = O_p(N^{-\delta}) \text{ where } \delta > 0.$$

**Proof.** To prove the result, we need to show that:

$$c_N(\hat{a}) \left( \frac{[\hat{E}_2(\hat{S})]^2}{\hat{E}_2(\hat{S}^2|\bar{v})} - \frac{[\hat{E}_2(\hat{S})]^2}{M_2(\hat{a})} \right) = O_p(N^{-\delta})$$

$$\text{and } c_N(\hat{a}) \left( \frac{[\hat{E}_2(\hat{S})]^2}{M_2(\hat{a})} - \frac{M_1(\hat{a})}{M_2(\hat{a})} \right) = O_p(N^{-\delta}).$$

Here, we provide the proof for the first equation, as the proof of the second is very similar.

Employing the definition of  $\hat{a}$ , we have  $c_N(\hat{a}) \frac{[\hat{E}_2(\hat{S})]^2}{\hat{E}_2(\hat{S}^2|\bar{v})} = O(1)$ . Hence the left-hand side has the same order as:

$$\begin{aligned} & \frac{M_2(\hat{a}) - \hat{E}_2(\hat{S}^2|\bar{v})}{M_2(\hat{a})} \\ = & \frac{[\hat{E}_2(S^2[\hat{a}, P_2]|\bar{v}) - \hat{E}_2(\hat{S}^2|\bar{v})]}{M_2(\hat{a})} + \frac{[M_2(\hat{a}) - \hat{E}_2(S^2[\hat{a}, P_2]|\bar{v})]}{M_2(\hat{a})} \\ \equiv & A + B. \end{aligned}$$

Beginning with term  $A$ , from a Taylor series expansion:

$$\begin{aligned} S^2(\hat{a}, \hat{P}_a) - S^2(\hat{a}, P_2) &= \sum_{k=1}^{m-1} [S^2]^{(k)}(\hat{a}, P_2) \left[ \frac{\hat{P}_a - P_2}{1 - P_2} \right]^k / k! \\ &+ [S^2]^{(m)}(\hat{a}, P_2^+) \left[ \frac{\hat{P}_a - P_2}{1 - P_2^+} \right]^m / m! \end{aligned}$$

where  $[S^2]^{(m)}$  is the  $m$ th derivative of the function  $S^2$  w.r.t  $x$ , and  $P_2^+$  is an intermediate point. Hence we have to show

$$\begin{aligned} A_k &\equiv \left| \sum_{i=1}^N \frac{1}{N} [S^2]^{(k)}(\hat{a}, P_2) \left[ \frac{\hat{P}_a - P_2}{1 - P_2} \right]^k K_2 \left( \frac{\bar{v} - v_i}{h_2} \right) / M_2(\hat{a}) h_2 \right| = O_p(N^{-\delta}) \\ A_m &\equiv \left| \sum_{i=1}^N \frac{1}{N} [S^2]^{(m)}(\hat{a}, P_2^+) \left[ \frac{\hat{P}_a - P_2}{1 - P_2^+} \right]^m K_2 \left( \frac{\bar{v} - v_i}{h_2} \right) / M_2(\hat{a}) h_2 \right| = O_p(N^{-\delta}). \end{aligned}$$

For  $A_k$ , setting  $k = 1$  for expositional purposes, the term will be bounded above by

$$2 \sup \left| \left( \frac{\hat{P}_{ai} - P_{2i}}{1 - P_{2i}} \right) S^{(1)}(\hat{a}, P_2) \right| \sum_{i=1}^N \frac{1}{N} S(\hat{a}, P_2) K_2 \left( \frac{\bar{v} - v_i}{h_2} \right) / M_2(\hat{a}) h_2$$

where  $S^{(m)}$  is the  $m^{\text{th}}$  derivative of  $S$  w.r.t  $x$ . The  $\sup \left| \left( \frac{\hat{P}_{ai} - P_{2i}}{1 - P_{2i}} \right) S^{(1)}(\hat{a}, P_2) \right| = O_p(N^{-\delta})$  because  $\hat{P}_a$  is based on higher order kernels and converges to  $P_2$  faster than  $N^{-\hat{a}}$ , which is the order of the denominator



$1 - P_2$ .<sup>10</sup> Turning our attention to the second part of the above expression, it is bounded above by:

$$\sup_a \left| \sum_{i=1}^N \frac{1}{N} S(a, P_2) K_2 \left( \frac{\bar{v} - v_i}{h_2} \right) / h_2 - E(S(a, P_2) | \bar{v}) \right| / M_2(\hat{a}) + \frac{E(S(a, P_2) | \bar{v})}{M_2(\hat{a})}.$$

For the first term, from uniform convergence, the numerator is converging to zero at a rate arbitrarily close to  $N^{-4}$ . For the denominator, referring to (D2), notice that the  $S^2$  function is bounded below by an indicator function set to be zero when  $x$  is in either  $R1$  or  $R2$ , and one when  $x$  is in  $R3$ . Hence  $M_2(a)$  is bounded below by the conditional expectation of that indicator function, which is a probability that is of order  $N^{-a}$  provided that the tail of  $v_2 | \bar{v}$  is fatter than the error tail (A5a). Therefore,  $M_2(\hat{a}) \geq O(N^{-\hat{a}}) \geq O(N^{-(4 - \frac{\varepsilon - \eta}{2})})$  (see D3). Hence it suffices to show that  $\sup_a E(S(a, P_2) | \bar{v}) / M_2(a) = O(1)$ . Referring to (D2), notice that

$$\frac{E(S(a, P_2) | \bar{v})}{M_2(a)} = \frac{c_1 \Pr(R2 | \bar{v}) + \Pr(R3 | \bar{v})}{c_2 \Pr(R2 | \bar{v}) + \Pr(R3 | \bar{v})}$$

where

$$\begin{aligned} c_1 &\equiv E \left[ 1 - \exp \frac{-x^k}{b^k - x^k} | R2, \bar{v} \right] \\ c_2 &\equiv E \left[ \left( 1 - \exp \frac{-x^k}{b^k - x^k} \right)^2 | R2, \bar{v} \right]. \end{aligned}$$

The above ratio converges to some constant irrespective of which of the regional probabilities converges faster to zero.

For the remainder term  $A_m$ , it vanishes faster than  $N^{-\delta}$  for  $m$  sufficiently large. For term  $B$ , the argument is very similar to that above.

**Lemma 6.** Referring to Lemma 5, define

$$\begin{aligned} z(a) &\equiv N^{-a} \\ R(z(a)) &\equiv \frac{M_1(a)}{M_2(a)} \\ z_0 &\equiv z(a_{0N}), \hat{z} \equiv z(\hat{a}), z^+ \equiv z(a^+). \end{aligned}$$

Then, for  $\delta > 0$ :  $|\hat{a} - a_{0N}| = o_p(N^{-\delta})$ .

---

<sup>10</sup>Note that  $\hat{a} < .4$ , and  $S^{(1)}(\hat{a}, P_2)$  restricts  $1 - P_2$  to a middle region  $R2$  where:  $\frac{N^{-\hat{a}}}{\exp(b)} < 1 - P_2 < N^{-\hat{a}}$ .

**Proof.** From Lemma 5 and a Taylor series expansion:

$$\begin{aligned}
c_N(\hat{a}) \frac{\left[ \hat{E}_2(\hat{S}) \right]^2}{\hat{E}_2(\hat{S}^2|\bar{v})} &= c_N(\hat{a}) R(z(\hat{a})) + O_p(N^{-\delta}) \\
&= c_N(a_{0N}) R(z(a_{0N})) + O_p(N^{-\delta}) - \\
&\quad Ln(N) c_N(z^+) [2R(z^+) + z^+ R'(z^+)] [\hat{a} - a_{0N}].
\end{aligned}$$

Therefore, since  $R = \frac{[E(S)]^2}{E(S^2|\bar{v})}$  in (D3) is increasing,  $z^+ R'(z^+) > 0$ ; hence

$$[\hat{a} - a_{0N}] = O_p\left(N^{-\delta}/Ln(N) c_N(z^+) R(z^+)\right).$$

Suppose  $\hat{a} < a_{0N}$  (the argument when  $\hat{a} \geq a_{0N}$  is the same), then  $\hat{a} < a^+ < a_{0N}$ , and  $z_0 < z^+ < \hat{z}$ . Since  $R$  is increasing,  $R(z_0) < R(z^+) < R(\hat{z})$ . Therefore,  $c_N(z_0) R(z_0) < c_N(z^+) R(z^+) < c_N(\hat{z}) R(\hat{z})$ . The proof now follows.

**Lemma 7. Expectations of Kernel Products.** Let  $\{\varepsilon_{1j}, \varepsilon_{2j}, \varepsilon_{3j}\}$  be i.i.d. over  $j$  with properties:

$$\begin{aligned}
a) &: E(\varepsilon_{\gamma j}) = O(h^{2p}) \\
b) &: E\left[\varepsilon_{\gamma j}^\rho\right] = \frac{1}{h^{\rho-1}}, \rho > 1.
\end{aligned}$$

Set  $h^{4p} = O(\frac{1}{Mh})$  and denote  $\bar{\varepsilon}_\gamma = \frac{1}{M} \sum_{j=1}^M \varepsilon_{\gamma j}$ ,  $\gamma = 1, 2, 3$ , then

$$\begin{aligned}
E\left\{[\bar{\varepsilon}_1]^r [\bar{\varepsilon}_2]^s [\bar{\varepsilon}_3]^t\right\} &= O\left(h^{2p4r/4}\right) O\left(h^{2p4s/4}\right) O\left(h^{2p2t/2}\right) \\
&= O\left(h^{2p(r+s+t)}\right).
\end{aligned}$$

**Proof.** From the Cauchy–Schwartz inequality

$$E\left\{[\bar{\varepsilon}_1]^r [\bar{\varepsilon}_2]^s [\bar{\varepsilon}_3]^t\right\} < \left\{E[\bar{\varepsilon}_1]^{4r}\right\}^{1/4} \left\{E[\bar{\varepsilon}_2]^{4s}\right\}^{1/4} \left\{E[\bar{\varepsilon}_3]^{2t}\right\}^{1/2}.$$

It suffices to order one of the three terms, hence we can study a general term:  $E[\bar{\varepsilon}]^q$ . This general term

has  $q$  types of terms, with  $k$ th ( $k = 1, \dots, q$ ) type:

$$\frac{1}{M^{q-k}} \underbrace{\sum \cdots \sum}_k \frac{1}{M^k} \varepsilon_{j_1}^{i_1} \cdots \varepsilon_{j_k}^{i_k}$$

where  $i_1 + \dots + i_k = q$  and  $j_1 \neq \dots \neq j_k$ .

From the i.i.d. property of the  $\varepsilon$ 's, the expectation of this term is given as:

$$\frac{1}{M^{q-k}} E \left[ \varepsilon_{j_1}^{i_1} \right] \cdots E \left[ \varepsilon_{j_k}^{i_k} \right].$$

Suppose we study a term  $E \left[ \varepsilon_{j_t}^{i_t} \right]$ , where  $1 \leq t \leq k$ . There are two types of expectations: single power of  $\varepsilon$  and multiple powers of  $\varepsilon$ . For the single power case, from property (a):

$$E \left[ \varepsilon_{j_t}^{i_t} \right] = O(h^{2p}) \text{ for } i_t = 1.$$

For the multiple power case, from property (b):

$$E \left[ \varepsilon_{j_t}^{i_t} \right] = O \left( \left[ \frac{1}{h} \right]^{(i_t-1)} \right) \text{ for } i_t > 1$$

There are different combinations of  $i_1, \dots, i_k$  for a given  $q$  and  $k$ . To order  $\frac{1}{M^{q-k}} E \left[ \varepsilon_{j_1}^{i_1} \right] \cdots E \left[ \varepsilon_{j_k}^{i_k} \right]$ , we next need to find the combination which yields the slowest convergence rate. One observation we make here is that the slowest term is the one with the least number of single power  $\varepsilon$ 's.

When  $k \leq \frac{q}{2}$ , the slowest term would have no single power of  $\varepsilon$  in it (see below for an example). Therefore, from property (b) the convergence rate will be:

$$O \left( \left( \frac{1}{Mh} \right)^{q-k} \right).$$

When  $k > \frac{q}{2}$ , the slowest term would include at least one single power  $\varepsilon$  in it, hence from property (a) and (b) the rate will be:

$$O \left( (h^{2p})^{2k-q} \right) O \left( \left( \frac{1}{Mh} \right)^{q-k} \right).$$

To illustrate our proof strategy, suppose  $q$  is an even number (the odd number case is very similar), for

example  $q = 6$ . If we denote the type by the powers of the elements, for example 1122 would mean the  $\varepsilon_{j_1}^1 \varepsilon_{j_2}^1 \varepsilon_{j_3}^2 \varepsilon_{j_4}^2$  term, then we have the following table:

$k$	<i>Slowest Type</i>	<i>Rate</i>
1	6	$O\left[\left(\frac{1}{Mh}\right)^5\right]$
2	33	$O\left[\left(\frac{1}{Mh}\right)^4\right]$
3	222	$O\left[\left(\frac{1}{Mh}\right)^3\right]$
4	1122	$O\left[(h^{2p})^2\left(\frac{1}{Mh}\right)^2\right]$
5	11112	$O\left[(h^{2p})^4\left(\frac{1}{Mh}\right)\right]$
6	111111	$O\left[(h^{2p})^6\right]$

For  $k = \frac{q}{2}$ , the slowest term would be the one with  $k$  squared terms

$$\frac{1}{M^{q-k}} E[\varepsilon_{j_1}^2] \cdots E[\varepsilon_{j_k}^2]$$

(examples of faster terms: 114 and 123 types); hence the rate would be

$$O\left(\frac{1}{M^{q-k}} \left(\frac{1}{h}\right)^k\right) = O\left(\left(\frac{1}{Mh}\right)^{q-k}\right).$$

It can be shown that this same expression holds for all smaller  $k$ . For  $k = \frac{q}{2} + 1$ , the slowest term would have two single power  $\varepsilon$ 's and the rest are squared terms, e.g.

$$\frac{1}{M^{q-k}} E[\varepsilon_{j_1}] E[\varepsilon_{j_2}] E[\varepsilon_{j_3}^2] \cdots E[\varepsilon_{j_k}^2]$$

hence the rate would be

$$\frac{1}{M^{q-k}} O\left((h^{2p})^2 \left(\frac{1}{h}\right)^{k-2}\right) = O\left((h^{2p})^{2k-q} \left(\frac{1}{Mh}\right)^{q-k}\right).$$

This same expression holds for all larger  $k$ .

We now need to find the  $k^{th}$  term with the slowest convergence rate. Set  $h$  optimally, i.e.  $h^{4p} = O\left(\frac{1}{Mh}\right)$

and substitute it in each term above, we have

$$\begin{aligned} O\left((h^{4p})^{q-k}\right) &= O\left((h^{2p})^{2q-2k}\right) \quad \text{when } k \leq \frac{q}{2}, \\ O\left((h^{2p})^{2k-q}(h^{4p})^{q-k}\right) &= O(h^{2pq}) \quad \text{when } k > \frac{q}{2}. \end{aligned}$$

hence the slowest convergence rate is  $O(h^{2pq})$ . The lemma follows.

**Lemma 8.** Define

$$\begin{aligned} \hat{E}(Y_2S(a, P_2) | \bar{v}) &\equiv \sum_j \frac{1}{Nh} K[(\bar{v} - V_{1j})/h] Y_{2j} S(a, P_{2j}) \\ M_3(a) &\equiv E\left(\sum_j \frac{1}{Nh} K[(\bar{v} - V_{1j})/h] Y_{2j} S(a, P_{2j})\right). \end{aligned}$$

Then

$$\frac{M_3(a_{0N}) - \hat{E}\left(Y_2S\left(\hat{a}, \hat{P}_a\right) | \bar{v}\right)}{M_3(a_{0N})} \xrightarrow{p} 0.$$

**Proof.** The above term can be decomposed into the following as in Lemma 5:

$$\frac{\left[\hat{E}(Y_2S[a_{0N}, P_2] | \bar{v}) - \hat{E}\left(Y_2S\left[\hat{a}, \hat{P}_a\right] | \bar{v}\right)\right]}{M_3(a_{0N})} + \frac{\left[M_3(a_{0N}) - \hat{E}(Y_2S[a_{0N}, P_2] | \bar{v})\right]}{M_3(a_{0N})}.$$

The above two terms are similar to terms  $A$  and  $B$  in Lemma 5. Employing a Taylor series argument and utilizing the result from Lemma 6, it can be shown that the first term goes to zero in probability. The second term only requires pointwise convergence instead of the uniform convergence arguments in Lemma 5.

**Lemma 9.** For notational simplicity, we denote  $C_N = C_N(\bar{v})$  as in Theorem 1. Letting  $\vartheta_j \equiv (Y_{1j} - \zeta_0(\bar{v})) Y_{2j} K[(\bar{v} - V_{1j})/h] / h$ ,

$$C_N \left[ \hat{\zeta}(\bar{v}) - \zeta_0(\bar{v}) \right] = C_N \frac{\sum_j \frac{1}{N} \vartheta_j S(a_{0N}, P_{2j})}{M_3(a_{0N})} + o_p(1).$$

**Proof.** From Lemma 8:

$$C_N \left[ \hat{\zeta}(\bar{v}) - \zeta_0(\bar{v}) \right] = C_N \frac{\sum_j \frac{1}{N} \vartheta_j S\left(\hat{a}, \hat{P}_{aj}\right)}{M_3(a_{0N})} + o_p(1).$$

It remains to be shown that

$$C_N \frac{\sum_j \frac{1}{N} \vartheta_j \left[ S(\hat{a}, \hat{P}_{aj}) - S(a, P_{2j}) \right]}{M_3(a_{0N})} \xrightarrow{p} 0.$$

With  $S^{(m)}$  as the  $m^{\text{th}}$  derivative of  $S$  w.r.t  $x$ , for  $P_{2j}^+ \in [\hat{P}_{aj}, P_{2j}]$ ,  $a^+ \in [\hat{a}, a_{0N}]$ , a Taylor expansion provides:

$$\begin{aligned} S(\hat{a}, \hat{P}_{aj}) - S(a_{0N}, P_{2j}) &= \sum_{k=1}^K \sum_{l=1}^L \frac{1}{k!l!} T_{kl} \\ T_{kl} &\equiv S^{(k+l)}(\bar{a}, \bar{P}_{2j}) [Ln(N)]^k [\hat{a} - a_{0N}]^k \left[ \frac{\hat{P}_{aj} - P_{2j}}{1 - \bar{P}_{2j}} \right]^l \end{aligned}$$

where

$$\bar{a} \equiv \begin{cases} a_{0N} & k < K \\ a^+ & k = K \end{cases} \quad \bar{P}_{2j} \equiv \begin{cases} P_{2j} & l < L \\ P_{2j}^+ & l = L \end{cases}.$$

Substituting the Taylor Series expansion, and noting that  $[Ln(N)]^k [\hat{a} - a_{0N}]^k$  is converging to zero, we now need to show that terms of the following form converge in probability to 0:

$$C_N \frac{\sum_j \frac{1}{N} \vartheta_j S^{(k+l)}(\bar{a}, \bar{P}_{2j}) \left[ \frac{\hat{P}_{aj} - P_{2j}}{1 - \bar{P}_{2j}} \right]^l}{M_3(a_{0N})}.$$

To show that, we study the expectation of the square of the above term. Squaring the above term yields two types of elements:

$$\frac{1}{N} \frac{C_N^2}{N M_3(a_{0N})^2} \sum_j \left\{ \vartheta_j S^{(k+l)}(\bar{a}, \bar{P}_{2j}) \left[ \frac{\hat{P}_{aj} - P_{2j}}{1 - \bar{P}_{2j}} \right]^l \right\}^2 \quad \text{and} \quad (10)$$

$$\frac{C_N^2}{N^2 M_3(a_{0N})^2} \sum_j \sum_j \left\{ \vartheta_i S^{(k+l)}(\bar{a}, \bar{P}_{2i}) \left[ \frac{\hat{P}_{ai} - P_{2i}}{1 - \bar{P}_{2i}} \right]^l \right\} \left\{ \vartheta_j S^{(k+l)}(\bar{a}, \bar{P}_{2j}) \left[ \frac{\hat{P}_{aj} - P_{2j}}{1 - \bar{P}_{2j}} \right]^l \right\}. \quad (11)$$

For the first type, write:

$$\left[ \frac{\hat{P}_{aj} - P_{2j}}{1 - \bar{P}_{2j}} \right]^l = \left[ \frac{\bar{\varepsilon}_{1j}}{1 - \bar{P}_{2j}} \right]^l \left[ \frac{g_j}{\hat{g}_j} \right]^l \quad \text{where } \bar{\varepsilon}_{1j} \equiv (\hat{P}_{aj} - P_{2j}) \frac{\hat{g}_j}{g_j}.$$

From a Taylor series expansion with  $\bar{\varepsilon}_{2j} \equiv [\hat{g}_j - g_j]$  :

$$\begin{aligned} \left[ \frac{g_j}{\hat{g}_j} \right]^l &= \sum_{t=0}^m T_{tj} \\ T_{tj} &= (-1)^t \frac{1}{t!} \frac{(l+t-1)!}{(l-1)!} \left[ \frac{\bar{\varepsilon}_{2j}}{g_j} \right]^t \left( \frac{\hat{g}_j}{g_j} \right)^{(-l-m)I\{t=m\}}. \end{aligned}$$

Substituting a typical term  $T_{tj}$  into (10) we must show that the following expression converges in probability to 0:

$$\frac{1}{N} \frac{C_N^2}{NM_3(a_{0N})^2} \sum_j \left\{ \vartheta_j S'^{(k+l)}(\bar{a}, \bar{P}_{2j}) \left[ \frac{\bar{\varepsilon}_{1j}}{1 - \bar{P}_{2j}} \right]^l T_{tj} \right\}^2.$$

For non-remainder terms in the above Taylor series expansion ( $k < K, l < L, t < m$ ), the expectation of the above expression has the form:

$$\begin{aligned} & \frac{C_N^2}{NM_3(a_{0N})^2} E \frac{1}{N} \sum_j \vartheta_j^2 \left[ S'^{(k+l)}(a_{0N}, P_{2j}) \right]^2 \left[ \frac{\bar{\varepsilon}_{1j}}{1 - P_{2j}} \right]^{2l} \left[ \frac{\bar{\varepsilon}_{2j}}{g_j} \right]^{2t} \\ &= \frac{C_N^2}{NM_3(a_{0N})^2} E \left\{ E[\vartheta_j^2 | X_j] \left[ S'^{(k+l)}(a_{0N}, P_{2j}) \right]^2 \left[ \frac{\bar{\varepsilon}_{1j}}{1 - P_{2j}} \right]^{2l} \left[ \frac{\bar{\varepsilon}_{2j}}{g_j} \right]^{2t} \right\} \\ &= \frac{C_N^2}{NM_3(a_{0N})^2} E \left\{ E[\vartheta_j^2 | X_j] \left[ S'^{(k+l)}(a_{0N}, P_{2j}) \right]^2 E \left( \left[ \frac{\bar{\varepsilon}_{1j}}{1 - P_{2j}} \right]^{2l} \left[ \frac{\bar{\varepsilon}_{2j}}{g_j} \right]^{2t} \middle| V_j \right) \right\}, \end{aligned}$$

where  $E[\vartheta_j^2 | X_j]$  only depends on the indices  $V_j$ . From Lemma 7, the expectation conditioned on the indices is  $o(1)$ . Therefore, since  $C_N^2/N$  is converging to 0, the result follows as:

$$\frac{1}{M_3(a_{0N})^2} E \left\{ E[\vartheta_j^2 | X_j] \left[ S'^{(k+l)}(a_{0N}, P_{2j}) \right]^2 \right\} = O(1).$$

Turning to the expectation of the cross-product terms (11), define  $\hat{P}_{aj}[i]$  by removing from  $\hat{P}_{aj}$  its dependence on  $Y_{2i}$ . Similarly, define  $\hat{P}_{ai}[j]$ . Notice that  $\hat{P}_{aj}[i]$  and  $\hat{P}_{ai}[j]$  do not depend on  $Y_{2i}$  and  $Y_{2j}$ . Finally, denote:

$$\hat{P}_{ai} \equiv \hat{P}_{ai}[j] + \lambda_j(1 - P_{2i}); \quad \hat{P}_{aj} \equiv \hat{P}_{aj}[i] + \lambda_i(1 - P_{2j});$$

then the non-remainder terms in (11) have the following form:

$$\frac{C_N^2}{M_3(a_{0N})^2} E \left\{ \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \left[ \frac{(\hat{P}_{ai}[j] - P_{2i})}{1 - P_{2i}} + \lambda_j \right]^l \right\} \left\{ \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \left[ \frac{(\hat{P}_{aj}[i] - P_{2j})}{1 - P_{2j}} + \lambda_i \right]^l \right\}.$$

Performing the binomial expansion on  $\left[ \frac{(\hat{P}_{ai}[j] - P_{2i})}{1 - P_{2i}} + \lambda_j \right]^l$  and  $\left[ \frac{(\hat{P}_{aj}[i] - P_{2j})}{1 - P_{2j}} + \lambda_i \right]^l$ , the slowest converging term has the following form:

$$\begin{aligned} & \frac{C_N^2}{M_3(a_{0N})^2} E \left\{ \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \lambda_j \lambda_i \left( \frac{\hat{P}_{ai}[j] - P_{2i}}{1 - P_{2i}} \right)^{l-1} \left( \frac{\hat{P}_{aj}[i] - P_{2j}}{1 - P_{2j}} \right)^{l-1} \right\} \\ &= \frac{C_N^2}{M_3(a_{0N})^2} E \left\{ \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \lambda_j \lambda_i E \left[ \left( \frac{\hat{P}_{ai}[j] - P_{2i}}{1 - P_{2i}} \right)^{l-1} \left( \frac{\hat{P}_{aj}[i] - P_{2j}}{1 - P_{2j}} \right)^{l-1} \middle| V_i, V_j \right] \right\}. \end{aligned}$$

Applying Lemma 7 in a similar manner as above yields:

$$E \left[ \left( \frac{\hat{P}_{ai}[j] - P_{2i}}{1 - P_{2i}} \right)^{l-1} \left( \frac{\hat{P}_{aj}[i] - P_{2j}}{1 - P_{2j}} \right)^{l-1} \middle| V_i, V_j \right] = \frac{B(V_i, V_j)}{[(1 - P_{2i})(1 - P_{2j})]^{l-1}} o(N^{-2a(l-1)})$$

where  $B(V_i, V_j)$  is a bounded function. Since  $\lambda_j \lambda_i = \frac{1}{(1 - P_{2i})(1 - P_{2j})h_2^2 N^2} Y_i Y_j K_2 ([V_{2i} - V_{2j}] / h_2)^2$ , the absolute value of (11) is bounded above by

$$\frac{C_N^2}{N^2 h_2^2 M_3(a_{0N})^2} E \left\{ \left| \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \right| Y_i Y_j K_2 \left( \frac{V_{2i} - V_{2j}}{h_2} \right)^2 \left| \frac{B(V_i, V_j)}{[(1 - P_{2i})(1 - P_{2j})]^l} \right| o(N^{-2a(l-1)}) \right\}$$

Notice that the  $S$  derivatives restricts us to the middle region where  $1 - P_2 > N^{-a} / \exp(b)$  and hence  $\left| \frac{B(V_i, V_j)}{[(1 - P_{2i})(1 - P_{2j})]^l} \right| = O(N^{2al})$ . Since  $h_2 = N^{-1}$ , we have  $C_N^2 / h_2 N = o(1)$  and  $o(N^{2a}) / h_2 N = o(1)$ . Hence it suffices to show  $E \left| \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \right| / M_3(a_{0N})^2 = O(1)$ .

Recall that  $\vartheta_j \equiv (Y_{1j} - \zeta_0(\bar{v})) Y_{2j} K[(\bar{v} - V_{1j}) / h] / h$ . Further  $s^m$  is the  $m^{\text{th}}$  derivative of  $S$  w.r.t  $x$ , and it is zero except in region  $R2$ :

$$\frac{E \left| \vartheta_i S'^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S'^{(k+l)}(a_{0N}, P_{2j}) \right|}{M_3(a_{0N})^2} = \left( \frac{c_1 \Pr(R2|\bar{v})}{c_2 \Pr(R2|\bar{v}) + c_3 \Pr(R3|\bar{v})} \right)^2$$



where

$$\begin{aligned}
c_1 &\equiv \left\{ E \left[ \vartheta S^{(k+l)}(a_{0N}, P_2) \mid R2, \bar{v} \right] \right\} \\
c_2 &\equiv E \left[ \left( 1 - \exp \frac{-x^k}{b^k - x^k} \right) \left( \frac{1}{h} Y_2 K \left[ \frac{\bar{v} - V_1}{h} \right] \right) \mid R2, \bar{v} \right] \\
c_3 &= E \left[ \left( \frac{1}{h} Y_2 K \left[ \frac{\bar{v} - V_1}{h} \right] \right) \mid R3, \bar{v} \right].
\end{aligned}$$

Then, similar to Lemma 5, it follows that

$$E \left| \vartheta_i S^{(k+l)}(a_{0N}, P_{2i}) \vartheta_j S^{(k+l)}(a_{0N}, P_{2j}) \right| / M_3(a_{0N})^2 = O(1).$$

### 8.3.3 Index Lemmas

The next lemma proves that the estimated second-stage objective function  $\hat{L}^*(\theta)$  is uniformly close to  $L^*(\theta) \equiv \frac{1}{N} \sum_i \sum_{d_1, d_2} Y_i(d_1, d_2) \text{Ln} [P_i^*(d_1, d_2; \theta)]$ .

**Lemma 10.** With  $D \equiv \frac{1}{N} \sum_i \sum_{d_1, d_2} D_i$ , where  $D_i \equiv Y_i(d_1, d_2) \text{Ln} \left[ \hat{P}_i^*(d_1, d_2; \theta) / P_i^*(d_1, d_2; \theta) \right]$ ,

$$\sup_{\theta} |D| = o_p(1).$$

**Proof.** We prove this result when indices are restricted to be smoothly in  $\mathcal{V}_N$  and its complement.

Referring to (D6), define a smoothed indicator restricting  $v_i$  to  $\mathcal{V}_N$  in (9) as:

$$\begin{aligned}
l(v_i) &\equiv \prod_k \tau[a_k(d_k) + h_{ck}^{1-\alpha}, v_{ki}] \tau[v_{ki}, b_k(d_k) - h_{ck}^{1-\alpha}], \\
D &= \frac{1}{N} \sum_i \sum_{d_1, d_2} D_i l(v_i) + \frac{1}{N} \sum_i \sum_{d_1, d_2} D_i (1 - l(v_i)).
\end{aligned}$$

For the first term, from a Taylor series expansion of the log-probability ratio about one, from Lemma 1 it

converges to zero in probability, uniformly in  $\theta$ . For the second term, from Cauchy's inequality:

$$\begin{aligned} \sup_{\theta} \left| \frac{1}{N} \sum_i \sum_{d_1, d_2} D_i (1 - l(v_i)) \right| &\leq \sup_{\theta} \frac{1}{N} \sum_i \sum_{d_1, d_2} |D_i| (1 - l(v_i)) \\ &\leq \sup_{\theta} \sqrt{\frac{1}{N} \sum_i \sum_{d_1, d_2} D_i^2} \sup_{\theta} \sqrt{\frac{1}{N} \sum_i \sum_{d_1, d_2} |(1 - l(v_i))|}. \end{aligned}$$

It can be shown that  $\inf P_i^*(d_1, d_2; \theta)$  is bounded away from 0 and  $\inf \hat{P}_i^*(d_1, d_2; \theta)$  converges to a term bounded away from zero. Therefore, the first term above is finite. The second term converges in probability to zero as  $(1 - l(v_i))$  is smoothly zero except on a set of vanishing probability.

The next lemma proves that  $L^*(\theta)$ , which is the probability limit of  $\hat{L}^*(\theta)$  defined in (D10), is uniformly close to  $L(\theta) \equiv \frac{1}{N} \sum_i \sum_{d_1, d_2} Y_i(d_1, d_2) \text{Ln}[P_i(d_1, d_2; \theta)]$ . Therefore, we may ignore the probability adjustments  $\Delta$ 's in the adjusted likelihood,  $L^*$ .

**Lemma 11.** For  $\theta$  in a compact set:

$$\sup_{\theta} |L^*(\theta) - L(\theta)| \xrightarrow{p} 0.$$

**Proof.** The proof, which is very similar to that in Lemma 10, follows by analyzing this difference separately on  $\mathcal{V}_N$  and its complement as above. The following lemma shows that the trimming and the  $\delta$ 's in the gradient component can be taken as known.

**Lemma 12.** Referring to (D10-11) and the text right below (7), for  $\tau = \tau_v$  or  $\tau_x$ , let

$$\begin{aligned} \hat{A}^* &\equiv \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)] \hat{\delta}_i^*(d_1, d_2; \theta_0) \tau \\ A &\equiv \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)] \delta_i(d_1, d_2; \theta_0) \tau \end{aligned}$$

then

$$\hat{A}^* - A = o_p(N^{-1/2}).$$

**Proof.** Klein and Shen (2010) establish this result in a semiparametric least squares context for single index models. The argument extends to double index likelihood-based models.

Using Lemma 12, Lemma 13 provides a useful convergence rate for the initial estimator in (D10).

**Lemma 13.** For  $\hat{\theta}$  defined in (D10) and with  $h = O(N^{-r})$ ,  $r = \frac{1}{8+\xi}$  :

$$\left(\hat{\theta} - \theta_0\right) = O_p(h^2).$$

**Proof.** From a Taylor series expansion:

$$\begin{aligned} \left(\hat{\theta} - \theta_0\right) &= -\hat{H}(\theta^+)^{-1} \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} \left[ Y_i(d_1, d_2) - \hat{P}_i(d_1, d_2; \theta_0) \right] \hat{\delta}_i(d_1, d_2; \theta_0) \tau_{ix} = -\hat{H}(\theta^+)^{-1} \left[ \hat{A} - \hat{B} \right], \\ \hat{A} &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} \left[ Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0) \right] \hat{\delta}_i(d_1, d_2; \theta_0) \tau_{ix}; \\ \hat{B} &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} \left[ \hat{P}_i(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0) \right] \hat{\delta}_i(d_1, d_2; \theta_0) \tau_{ix}. \end{aligned}$$

Employing the same argument as in Lemma 12,  $\hat{A} - A = o_p(N^{-1/2})$ , and since  $A = O_p(N^{-1/2})$  we have  $\hat{A} = O_p(N^{-1/2})$ . Employing Lemma 2 and Cauchy-Schwartz argument, we have  $\hat{B} = O_p(h^2)$ , which completes the proof.

To obtain a convergence rate for the second-stage estimator and to analyze the final bias-adjusted estimator, Lemma 14 shows that the gradient component responsible for the bias in the estimator vanishes in probability.

**Lemma 14.**

$$\begin{aligned} a) \quad B^* &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ \hat{P}_i^*(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0) \right] \delta_i(d_1, d_2; \theta_0) \tau_{iv} = o_p(1) \\ b) \quad B^o &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ \hat{P}_i^o(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0) \right] \delta_i(d_1, d_2; \theta_0) \tau_{iv} = o_p(1). \end{aligned}$$

**Proof.** For a), under index trimming the adjustment factors within  $\hat{P}_i^*$  vanish exponentially. Therefore:

$$B^* = B + o_p(1), B = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ \hat{P}_i(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0) \right] \delta_i(d_1, d_2; \theta_0) \tau_{iv}.$$

Denote:

$$\hat{P}_{ci} = \hat{P}r(Y_{1i} = d_1 | Y_{2i} = d_2, V_i = t), \quad P_{ci} \equiv p \lim \hat{P}_{ci}.$$

With  $d_2 = 0$  write:

$$\hat{P}_i(d_1, d_2; \theta) - P_i(d_1, d_2; \theta) = (\hat{P}_{2i} - P_{2i}).$$

Otherwise:

$$\begin{aligned} \hat{P}_i(d_1, d_2; \theta) - P_i(d_1, d_2; \theta) &= \hat{P}_{2i}\hat{P}_{ci} - P_{2i}P_{ci} \\ &= (\hat{P}_{2i} - P_{2i})(\hat{P}_{ci} - P_{ci}) + (\hat{P}_{2i} - P_{2i})P_{ci} + P_{2i}(\hat{P}_{ci} - P_{ci}). \end{aligned}$$

For the second case (the first is similar and easier), we can rewrite the  $B$  term as:

$$B = \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ (\hat{P}_{2i} - P_{2i})(\hat{P}_{ci} - P_{ci}) + (\hat{P}_{2i} - P_{2i})P_{ci} + P_{2i}(\hat{P}_{ci} - P_{ci}) \right] \delta_i(d_1, d_2; \theta_0) \tau_{iv}.$$

For the first term in  $B$ , we may employ Cauchy's inequality and Lemma 2 to show that it vanishes in probability.

The difference between the second term in  $B$  and the following U-statistic converges in probability to zero:

$$\begin{aligned} U &\equiv \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ \left( \frac{\hat{f}_2(t_2; d_2)}{\hat{g}_2(t_2; d_2)} - P_{2i} \right) P_{ci} \right] \left[ \frac{\hat{g}_2(t_2; d_2)}{g_2(t_2; d_2)} \right] \delta_i(d_1, d_2; \theta_0) \tau_{iv} \\ &= \frac{1}{\sqrt{N}} \sum_{i=1}^N \sum_{d_1, d_2} \left[ (\hat{f}_2(t_2; d_2) - \hat{g}_2(t_2; d_2) P_{2i}) P_{ci} \right] \left[ \frac{\delta_i(d_1, d_2; \theta_0) \tau_{iv}}{g_2(t_2; d_2)} \right]. \end{aligned}$$

Notice that the U-statistic vanishes in probability from standard projection arguments; hence the second term in  $B$  vanishes. The third term in  $B$  has the same structure as the second and therefore also vanishes in probability, which completes the proof for a). The proof for b) is very similar.

**Lemma 15.** Referring to (D10), for the second stage estimator:

$$\left| \hat{\theta}^* - \theta_0 \right| = O_p \left( N^{-\frac{4}{8+\xi}} \right).$$

**Proof.** From Lemma 14, the initial estimator satisfies:  $(\hat{\theta} - \theta_0) = O_p(N^{-2r})$ . For the estimator

based on index trimming, from a standard Taylor series argument with  $\tau_{iv}$  replacing  $\tau_{ix}$  :

$$\begin{aligned} (\hat{\theta}^* - \theta_0) &= -\hat{H}^*(\theta^+)^{-1} [\hat{A}^* - \hat{B}^*], \\ \hat{A}^* &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [Y_i(d_1, d_2) - P_i(d_1, d_2; \theta_0)] \hat{\delta}^*(d_1, d_2; \theta_0) \tau_{iv}; \\ \hat{B}^* &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [\hat{P}_i^*(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0)] \hat{\delta}^*(d_1, d_2; \theta_0) \tau_{iv}. \end{aligned}$$

Referring to Lemma 13, since  $A = O_p(N^{-1/2})$ ,  $\hat{A}^* = O_p(N^{-1/2})$ .

For the  $\hat{B}^*$ -term, with  $\Delta_{Bi} \equiv [\hat{\delta}^*(d_1, d_2; \theta_0) \tau_{iv} - \delta(d_1, d_2; \theta_0) \tau_{iv}]$ :

$$\begin{aligned} \hat{B}^* &= B_1^* + \hat{B}_2^*, \\ B_1^* &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [\hat{P}_i^*(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0)] \delta(d_1, d_2; \theta_0) \tau_{iv} \\ \hat{B}_2^* &= \frac{1}{N} \sum_{i=1}^N \sum_{d_1, d_2} [\hat{P}_i^*(d_1, d_2; \theta_0) - P_i(d_1, d_2; \theta_0)] \Delta_{Bi} \end{aligned}$$

By showing that  $\hat{B}_1^*$  is close in probability to a centered U-statistic, Lemma 14, part a) proves that  $B_1^* = o_p(N^{-1/2})$ . Referring to the convergence rates in Lemma 2 and with window parameters  $r = r^* = \frac{1}{8+\xi}$ , it follows from Cauchy's inequality that  $\hat{B}_2^* = O_p\left(N^{-\frac{4}{8+\xi}}\right)$ ,  $\xi > 0$ . For these window choices, from the uniform rates in Lemma 1:  $\hat{H}^*(\theta^+) = H_o + o_p(1)$ . It now follows that  $|\hat{\theta}^* - \theta_0| = O_p\left(N^{-\frac{4}{8+\xi}}\right)$ .