

Efficient Estimation of Models with Conditional Moment Restrictions Containing Unknown Functions^a

Chunrong Ai
Department of Economics
University of Florida

Xiaohong Chen
Department of Economics
London School of Economics

August 1999

Abstract

We propose an estimation method for models of conditional moment restrictions that contain finite dimensional unknown parameters (μ) as well as unknown functions (h), and that the unknowns can interact with endogenous variables. Our proposal is to approximate h with sieves and then to estimate μ and the sieve parameters jointly by applying the method of generalized minimum distance. We show that: (1) the sieve estimator of h is consistent under certain metric, and may attain the usual nonparametric convergence rates; (2) the estimator of μ is \sqrt{n} consistent and asymptotically normally distributed; (3) the estimator for the asymptotic variance of the μ estimator is consistent and easy to compute; and (4) the optimally weighted minimum distance estimator of μ attains the semiparametric efficiency bound. These results extend the work of Hansen (1982) and Newey (1990, 1993) to a semiparametric conditional moment restrictions setting.

^aWe acknowledge helpful comments from Richard Blundell, Lars Hansen, Han Hong, Bo Honore, Roger Klein, Oliver Linton, Whitney Newey, Jim Powell, Frank Vella and Seminar participants at Texas Austin, Rochester, Chicago, Duke, Rutgers, Princeton, UCL, Carol III at Madrid, and the 1999 North American Winter Meetings of the Econometric Society. Ai's research is supported by the Warrington College of Business Administration Summer Grants for 1998 and 1999 at the University of Florida.

1. Introduction

A general framework for analyzing data $(Y; X)$ is to presume that the data satisfy some conditional moment restrictions:

$$E[\frac{1}{2}(Z; \theta_0)jX] = 0, \tag{1.1}$$

where $Z \sim (Y^0; X_z^0)^0$, $X_z \mu X$, θ_0 is a vector of unknown parameters, $\frac{1}{2}(\cdot; \cdot)$ is a vector of known functions, and $E[\frac{1}{2}(Z; \theta_0)jX]$ is the conditional expectation of $\frac{1}{2}(Z; \theta_0)$ given X , where the conditional distribution of Y given X is unspecified. A common method for estimating the parameters of interest, θ_0 , is the generalized method of moments (GMM) estimation proposed by Hansen (1982). The GMM converts the conditional moment restrictions (1.1) into the unconditional moment restrictions:

$$E[\frac{1}{2}(Z; \theta_0)W] = 0; \tag{1.2}$$

where W are instruments (e.g. X and nonlinear measurable transformations of X), and estimates θ_0 by solving the sample analog of (1.2). Under some sufficient conditions, Hansen (1982) showed that the GMM estimator of θ_0 is \sqrt{n} consistent and asymptotically normally distributed, where n is the sample size. Moreover, he established that the GMM estimator can attain the efficiency bound of the unconditional moment restrictions model (1.2). Chamberlain (1987) derived the efficiency bound of the conditional moment restrictions model (1.1). He showed that the GMM estimator attains his bound if W are the optimal instruments¹. Robinson (1987, 1991) and Newey (1990, 1993) proposed to estimate the optimal instruments nonparametrically for their respective models and showed that their estimators attain Chamberlain's (1987) bound.

In this paper, we extend all those results to the case where θ_0 may contain unknown functions. Specifically, we consider the model:

$$E[\frac{1}{2}(Z; \mu_0; h_0(\cdot))jX] = 0, \tag{1.3}$$

where $\theta_0 = (\mu_0; h_0)$ are the parameters of interest, $\mu_0 \in \mathbb{R}$ is a vector of unknown parameters, and $h_0(\cdot) = (h_{01}(\cdot); \dots; h_{0q}(\cdot))^0 \in \mathbb{H}$ is a vector of unknown functions. The dot " \cdot " in h_{0j} denotes the generic argument of $h_{0j}(\cdot)$ which may vary with applications and the subscript j . In some applications, for instance, the argument of h_{0j} is a subset of Z , or some known functions of Z up to unknown μ_0 , while in other applications the argument includes other unknown functions h_{0j} . The

¹The optimal instruments are often unknown functions of X .

argument may even include unobserved variables. In general, the argument of h_{oj} can be written as $\pm_j(Z; \theta_o; U)$, where $\pm_j(\cdot; \cdot; \cdot)$ is a vector of known functions and U are unobserved variables. But for exposition we use “ \cdot ” to denote the argument.

The extension to model (1.3) is needed because economic theories seldom produce exact functional forms and misspecifications in functional forms may lead to inconsistent parameter estimates. By specifying the model partially (i.e. including h_o as part of the unknown parameters), we can at least alleviate the inconsistency problem. In fact, many studies have already investigated some particular versions of (1.3). For example, Robinson (1988) and Fan, Härdle and Mammen (1998) studied the partially additive form $E[Y \mid X_0^0 \mu_o \mid \prod_{j=1}^q h_{oj}(X_j) \mid X] = 0$, where the argument X_j of h_{oj} is a subset of conditioning variables $X = (X_0^0; \dots; X_q^0)^0$. Powell, Stock and Stoker (1989) and Ichimura (1993) considered the index form $E[Y \mid X_1^0 \mu_{o1} \mid h_o(X_2^0 \mu_{o2}) \mid X_1; X_2] = 0$, where the argument $X_2^0 \mu_{o2}$ of h_o is a known function of conditioning variables X up to unknown μ_o . One can easily modify Newey and Powell’s (1989) nonparametric instrumental variables regression model into: $E[Y_1 \mid X_1^0 \mu_o \mid h_o(Y_2) \mid X_1; X_2] = 0$, where the argument Y_2 of h_o is a subset of the endogenous variables $Y = (Y_1; Y_2)^0$. One can also write down a version of transformation model with unknown link: $E[h_o(Y) \mid X_1^0 \mu_o \mid X_1; X_2] = 0$; where h_o is monotonic with some normalization say $h_o(0) = 1$. Chamberlain (1992) considered a particular version of (1.3): $E[\frac{1}{2}(Z; \mu_o; h_{o1}(\pm_1(Z)); \cdot \cdot \cdot; h_{oq}(\pm_q(Z))) \mid X] = 0$, where the arguments $\pm_j(Z)$ of h_{oj} for $j = 1; \dots; q$ are completely known functions of Z . Our next example will be a modified version of the semiparametric Engel curve model of Blundell, Browning and Crawford (2000): $\frac{1}{2}(Z; \theta_o) = (\frac{1}{2}_1(Z; \theta_o); \dots; \frac{1}{2}_r(Z; \theta_o))$; $\frac{1}{2}_l(Z; \theta_o) = Y_l \mid h_{ol}(X_2 \mid h_{o;r+1}(X_1^0 \mu_{o1})) \mid X_1^0 \mu_{o2l}$ and $E[\frac{1}{2}_l(Z; \theta_o) \mid X_1; X_2] = 0$ for $l = 1; \dots; r$, where the argument $X_2 \mid h_{o;r+1}(X_1^0 \mu_{o1})$ of h_{ol} contains another unknown function $h_{o;r+1}$. Additional examples can be found in Powell (1994) and Horowitz (1998).

Under the assumption that model (1.3) identifies the parameters of interest θ_o ,² we propose an estimator of θ_o that has desirable large sample properties. Our approach is to apply the method of generalized minimum distance. Heuristically, if the functional forms of $h_o(\cdot)$ and the conditional distribution of Y given X ,

²This paper mainly tries to provide root-n efficient estimators for μ_o rather than pointwise estimation of h_o . In some cases such as the semiparametric IV example and the transformation with unknown link example, the identification of h_o is in certain weak metric sense rather than pointwise strong identification, see Sections 2 and 3 for examples and precise definition of identification.

$F_{Y|X}$, were known, then the functional form of the conditional expectation

$$m(x; \theta) = \int_{\mathcal{Z}} h(y; x_Z; \mu; h(\cdot)) dF_{Y|X=x}(y)$$

would be known. The minimum distance estimation of μ_0 would be the appropriate method since condition (1.3) implies:

$$\theta_0 : \inf_{\theta=(\mu;h) \in H} E \int m(X; \theta)^0 [\mathcal{S}(X)]^i m(X; \theta)^i \quad (1.4)$$

where $\mathcal{S}(X)$ is a positive definite matrix for any given X , (and " \cdot^0 " denotes the matrix transpose in this paper). Since the functional forms of $h_0(\cdot)$ and $F_{Y|X}$ (hence $m(\cdot; \cdot)$) are unknown, we estimate $m(X; \theta)$ nonparametrically with $\hat{m}(X; \theta)$ using observations on $\frac{1}{2}(Z; \theta)$ and X . The true values θ_0 would then be estimated by minimizing the sample analog of (1.4) with $m(X; \theta)$ replaced by $\hat{m}(X; \theta)$. However, the resulting estimator could be inconsistent or converging arbitrarily slowly when the parameter space H is too large. To address this problem, we follow the sieve literature (Grenander, 1981) by approximating H with a sequence of sieve spaces H_n . The sieve spaces are "computable" (often finite-dimensional) parameter spaces which increase with sample size and become dense in the original functional space H . We then estimate the true values θ_0 by minimizing the sample analog of a nonparametric version of (1.4) with h restricted to the sieve spaces H_n :

$$\hat{\theta}_n = (\hat{\beta}_n; \hat{h}_n) : \min_{\theta=(\mu;h) \in H_n} \frac{1}{n} \sum_{i=1}^n \hat{m}(X_i; \theta)^0 [\hat{\mathcal{S}}(X_i)]^i \hat{m}(X_i; \theta)^i \quad (1.5)$$

where $\hat{\mathcal{S}}(X)$ is a consistent estimator of $\mathcal{S}(X)$.

Under a set of sufficient conditions, we show that: (1) the sieve estimator \hat{h}_n is consistent under certain metric, and may attain the usual nonparametric convergence rates; and (2) the estimator $\hat{\beta}_n$ is \sqrt{n} consistent and asymptotically normally distributed. These results extend the work of Shen (1997) and Chen and Shen (1998) to the semiparametric conditional moments framework. In addition, we show that: (3) the estimator for the asymptotic covariance of $\hat{\beta}_n$ is consistent and easy to compute; and (4) when the weighting matrix $\mathcal{S}(X)$ is set to $\mathcal{S}_0(X) = \text{Var}[\frac{1}{2}(Z; \theta_0)|X]$ (assumed to be positive-definite), the estimator $\hat{\beta}_n$ (hereafter the optimally weighted minimum distance estimator) attains the semiparametric efficiency bound of model (1.3). While establishing the latter result, we generalize the efficiency bound result of Chamberlain's (1992) to the much more general model (1.3).

The literature on semiparametric efficient estimation is large. The papers most closely related to our approach are those by Severini and Wong (1992), Shen (1997), Hansen (1982) and Newey (1990, 1993). The first two sets of authors studied efficient estimation in the semiparametric maximum likelihood setting. Hansen proposed efficient estimation for the parametric unconditional moments restrictions (1.2), while Newey investigated efficient estimation of the parametric conditional moments restrictions (1.1). Our approach is an extension of all those studies to the semiparametric conditional moments restrictions (1.3). Most importantly, our paper provides the semiparametric efficient estimation of models contain unknown functions of endogenous variables, such as the semiparametric IV example and the transformation with unknown link example.

The rest of the paper is organized as follows. Section 2 formally introduces the estimator for model (1.3) and three examples as illustration: partially additive model with unknown heteroskedastic variance; semiparametric Engel curves estimation; semiparametric instrumental variables regression. Section 3 provides general results on consistency with convergence rates of $\hat{\beta}_n$. Section 4 derives the \sqrt{n} asymptotic normal distribution of $\hat{\beta}_n$. Section 5 provides a consistent estimator for the asymptotic covariance of $\hat{\beta}_n$. Section 6 investigates the efficiency property of $\hat{\beta}_n$. Section 7 presents “low-level” sufficient conditions to obtain convergence rates of $\hat{\beta}_n$ and efficient estimators of β_n for the three examples. Section 8 illustrates our estimation procedure by a Monte Carlo study and Section 9 concludes. All technical proofs are relegated to the Appendices.

2. Estimator and Examples

Throughout the paper, we assume that the sample observations $f(Y_i; X_i) : i = 1, 2, \dots, n$ are drawn independently from the distribution of $(Y; X)$ on $Y \in X$ with $X = [X_l; X_u]$ a compact subset of \mathbb{R}^s , and that the unknown distribution of $(Y; X)$ satisfies the conditional moment restriction (1.3), where $\eta : Z \in \mathbb{A} \rightarrow \mathbb{R}^r$ a given known mapping, $Z \in (Y^0; X_z^0) \in \mathbb{Z} \subset Y \in X_z, X_z \subset X$; and $\theta_0 \in (\mu_0; h_0) \in \mathbb{A} \subset \mathbb{E} \in H$, with \mathbb{E} a compact (with nonempty interior) subset of \mathbb{R}^b and $H \subset H^1 \in \mathcal{C} \in H^q$ a functional space of q -dimensional continuous functions.

To implement the proposed procedure, we need an estimator of $m(x; \theta)$. We use the kernel method described in Silverman (1986) to estimate $m(x; \theta)$ for arbitrary $(x; \theta)$. Let $K : \mathbb{R}^s \rightarrow \mathbb{R}$ denote a known symmetric function and a_n a

bandwidth satisfying $a_n \rightarrow 0$ as $n \rightarrow \infty$. We estimate the density of X , f_X , by

$$\hat{f}_{X_i} = \frac{1}{(n_i - 1)a_n^s} \sum_{j \in I: j=1}^n K \left(\frac{X_i - X_j}{a_n} \right)$$

and the conditional expectation $m(X_i; \theta)$ by

$$\hat{m}(X_i; \theta) = \frac{[(n_i - 1)a_n^s]^{-1} \sum_{j \in I: j=1}^n \frac{1}{2} (Z_j; \mu; h(\cdot)) K \left(\frac{X_i - X_j}{a_n} \right)}{\hat{f}_{X_i}}$$

here the dot " \cdot " in $h(\cdot)$ denotes the fact that h may depend on data and hence varies with observations Z_j .

As noted in previous studies (e.g. Ichimura and Lee (1991)), one drawback of the kernel estimator is that it may be biased at the boundary points. To avoid the bias, trimming is often introduced. Suppose that X_l and X_u are known. For some $0 < \alpha < 1$, define $X_n = [X_l + \alpha a_n^s; X_u - \alpha a_n^s]$. We use the indicator function $1_{X \in X_n}$ to trim the boundaries. The proposed estimator $\hat{\theta}_n = (\hat{\beta}_n; \hat{h}_n) \in \mathbb{R}^p \times \mathbb{R}^q$ is now formally defined as the solution to:

$$\min_{\theta \in \mathbb{R}^p \times \mathbb{R}^q} \frac{1}{n} \sum_{i=1}^n 1_{X_i \in X_n} \left[\hat{m}(X_i; \theta) - \hat{g}(X_i) \right]^2 \quad (2.1)$$

To compute the estimator, the boundary points X_l and X_u are unknown and must be replaced by some consistent estimates. The order statistics, for example, are such consistent estimates. Since the order statistics converge to the boundary points at the rate $\frac{1}{n}$, replacing X_l and X_u with the order statistics will not affect the asymptotic results of the proposed estimator. In addition, we must be specific about the sieve approximation H_n . The most popular sieves in econometrics are linear sieves, which can be written as $H_n = H_n^1 \subset \dots \subset H_n^q$ with H_n^j given by

$$H_n^j = \left\{ \sum_{k=1}^q \gamma_{jk} B_k^j(\cdot) \text{ for all } \gamma_{jk} = (\gamma_{j1}, \dots, \gamma_{jk_j}) \in \mathbb{R}^{j_j} \right\}$$

where $B_k^j(\cdot), k = 1, 2, \dots, j$ denotes a sequence of known base functions (e.g. Fourier series, orthogonal polynomials, splines, power series, wavelets, etc.). The unknown functions h_{0j} ; in this case, are replaced by $h_n = (h_{n1}, \dots, h_{nq})^0$. The proposed estimator is then computed by estimating the sieve coefficients $\gamma_{jk}^j; j = 1, \dots, q$ jointly with μ . Our estimation procedure and the large sample results below

permit any types of sieve approximations including neural networks and other nonlinear sieves. The sieve method is easy to implement and convenient here since h_0 may enter $\eta(Z; \theta_0)$ nonlinearly. Moreover, the sieve method allows us to impose any known structures such as additivity and shape restrictions on h_0 implied by the population restrictions (1.3)³.

Throughout this paper, we denote the (pathwise) directional derivative of $\eta(Z; \theta_0)$ with respect to θ_0 at direction $[\dot{\theta}_i; \dot{\theta}_0]$ as:

$$\dot{\eta}_{\theta_0}[Z; \dot{\theta}_i; \dot{\theta}_0] \equiv \lim_{\zeta \rightarrow 0} \frac{\eta(Z; \theta(\zeta)) - \eta(Z; \theta_0)}{\zeta}$$

where $\theta(\zeta)$ is an one-dimensional continuous path in Θ satisfying $\theta(0) = \theta_0$ and $\theta(1) = \theta_0 + \dot{\theta}$. In many applications such as Robinson's (1988) partial linear regression and Ichimura's (1993) single index model, the directional derivative $\dot{\eta}_{\theta_0}[Z; \dot{\theta}_i; \dot{\theta}_0]$ does not depend on the endogenous variables Y . Then the kernel estimation (hence the trimming) or any nonparametric estimation of the conditional expectations function $m(x; \theta) \equiv \int \eta(y; x; \mu; h(\cdot)) dF_{Y|X=x}(y)$ is not needed. One can simply perform the following sieve generalized least squares regression:

$$\min_{\theta = (\mu; h)} \frac{1}{n} \sum_{i=1}^n \eta(Z_i; \theta)^2 [\mathbf{S}(X_i)]^{-1} \eta(Z_i; \theta); \quad (2.2)$$

which is much easier to compute when s (the dimension of X) is big.

Throughout this paper, we shall consider three examples as illustration of the usefulness of our new results. The first example is a partially additive model with conditional heteroskedastic variance of unknown form. This example slightly extends the popular partially linear regression of Robinson (1988). Although this is a simple example, but as far as we know there is still no published results on semiparametric efficient estimators under unknown heteroskedastic variance⁴. The second example is semiparametric Engel curves estimation, which is a modification of the model studied by Blundell, Browning and Crawford (2000) in their empirical investigation but nevertheless is still consistent with the rational consumer's demand behavior as discussed in Blundell et al (2000). The third example

³See e.g., Newey, Powell and Vella (1999) for using the spline sieve to estimate unknown functions with additive structure, Chen and Conley (1999) for using the shape-preserving wavelet cardinal B-spline sieve to estimate a conditional isotropic covariance function.

⁴The estimators proposed by Robinson (1988) and Li (1999) are semiparametric efficient only under conditional homoskedastic variance.

is semiparametric instrumental regression, which can be regarded as an alternative to the empirical model studied by Newey, Powell and Vella (1999, page 584, equation (7.1)).

Example 2.1 (Partially additive regression with unknown conditional heteroskedastic variance):

$$Y_i = X_{0;i}^0 \mu_0 + \sum_{j=1}^q h_{oj}(X_{j;i}) + \epsilon_i;$$

$$E[\epsilon_i | X_i] = 0; \quad E[\epsilon_i^2 | X_i] = S_o(X_i);$$

where $\theta_o = (\mu_0; h_{o1}; \dots; h_{oq})$; $Z = (Y; X_Z)$ and $X_Z = X =$ non-overlap union of $X_0; X_1; \dots; X_q$: Without loss of generality, we assume that $\dim(X_0) = b$, $\dim(X_j) = s_j$ and $\dim(X) = s = b + \sum_{j=1}^q s_j$. Obviously $\dim(\mu) = b$ and $\dim(h) = q$. It is easy to see that $E[\frac{1}{2}(Z; \theta_o)jX] = 0$ where $\frac{1}{2}(Z; \theta_o) = Y - X_0^0 \mu_0 - \sum_{j=1}^q h_j(X_j)$: Since the directional derivative

$$\frac{1}{2}_{\theta_o}[Z; \theta_o | \theta_o] = \sum_{j=1}^q X_0^0 [\mu_1 - \mu_0] + \sum_{j=1}^q [h_j(X_j) - h_{oj}(X_j)]$$

does not depend on the endogenous variable Y , we may apply the simpler sieve generalized least squares regression (2.2) to estimate $(\mu_0; h_{o1}; \dots; h_{oq})$ instead of the method of generalized minimum distance (2.1).

Assume that $h_{oj} \in H^j = B_{p;1}^{m_j}(C_j)$; ($m_j > 0; 2 < p < \infty; 0 < C_j < 1$), a Besov ball which consists of functions in $L_p(\mathbb{R}^{s_j})$ with $\|h_{oj}\|_{B_{p;1}^{m_j}} = C_j$. There are several equivalent definitions of a Besov space, we follow the one given in Meyer (1990). Suppose there exists a multivariate scaling function \hat{A} which has a compact support and continuous derivatives up to order $r_j \geq m_j$. Then there also exist $2^{s_j} - 1$ associated wavelets $\tilde{A}^e; e \in \{1, \dots, 2^{s_j} - 1\}$ such that for a given integer k_0 , $\{\hat{A}_{k_0;i}; \tilde{A}_{k;i}^e; k \geq k_0; (i;l) \in \mathbb{Z}^{2s_j}; e \in \{1, \dots, 2^{s_j} - 1\}\}$ is an orthonormal basis of $L_2(\mathbb{R}^{s_j})$, where $\hat{A}_{k;i}(z) = 2^{k-2} \hat{A}(2^{k-2} z - i)$ and $\tilde{A}_{k;i}^e(z) = 2^{k-2} \tilde{A}^e(2^{k-2} z - i)$ for $i \in \mathbb{Z}^{s_j}$, we have for any $h \in L_2(\mathbb{R}^{s_j})$,

$$h = \sum_{i \in \mathbb{Z}^{s_j}} a_{k_0;i} \hat{A}_{k_0;i} + \sum_{k \geq k_0} \sum_{i \in \mathbb{Z}^{s_j}} \sum_{e \in \{1, \dots, 2^{s_j} - 1\}} c_{k;i}^e \tilde{A}_{k;i}^e$$

Now for $m_j > 0; 2 < p < \infty; h \in B_{p;1}^{m_j}(\mathbb{R}^{s_j})$ if and only if

$$\|h\|_{B_{p;1}^{m_j}} = \|a_0\|_p + \sup_{k \geq 0} 2^{k(m_j + s_j(1-2^{-1-p}))} \|c_k\|_p < \infty$$

where

$$jja_{0:i}jj_p^p = \prod_{i \in Z^{S_j}} ja_{0:i}j^p; jjc_{k:i}jj_p^p = \prod_{i \in Z^{S_j} \setminus e} \prod_{i \in 2^{S_j} \setminus i} jc_{k:i}j^p$$

We consider the wavelet sieve space:

$$H_n^j = \left\{ h \in L_2(\mathbb{R}^{S_j}) : h = \sum_{k=0}^{\infty} \sum_{i \in Z^{S_j} \setminus e} \sum_{i \in 2^{S_j} \setminus i} c_{k:i}^e \tilde{A}_{k:i}^e : jjhjj_{B_{p;1}^{m_j}} \leq C_j \right\} \quad (2.3)$$

In (2.2), we let \mathbf{I} be the identity matrix and $H_n^j, j = 1, \dots, q$, be the wavelet sieve space (2.3). Then by Corollary 3.2 in Section 3, the root-mean square convergence rate for the wavelet sieve estimator of $h_{oj}, j = 1, \dots, q$, will be $O_p(n^{-1/(m_j + s_j)})$, and by Theorem 4.1, the corresponding estimator of μ_0 will be \sqrt{n} asymptotic normal. Moreover when $\hat{S}(X)$ in (2.2) is a consistent estimator of $S_0(X)$ as provided in Section 6, then by Theorems 6.1 and 6.2 in Section 6, the corresponding estimator of μ_0 will be \sqrt{n} efficient with asymptotic variance V_0^{-1} :

$$V_0 = \inf_{w_j} E \left[\sum_{j=1}^q [X_0^0] w_j(X_j) \right]^0 S_0(X)^{-1} \left[\sum_{j=1}^q [X_0^0] w_j(X_j) \right]; \quad (2.4)$$

where $w = (w_1, \dots, w_q)$; w_j measurable function of X_j , and $w_j \in B_{p;1}^{m_j}(C_j)$ for all $j = 1, \dots, q$.

Example 2.2 (Partially additive IV regression with unknown conditional heteroskedastic variance, see Chamberlain (1992, model (3.5), page 579)):

$$Y_{1i} = Y_{2i}^0 \mu_0 + \sum_{j=1}^q h_{oj}(X_{j;i}) + \epsilon_i;$$

$$E[\epsilon_i | X_i] = 0; \quad E[\epsilon_i^2 | X_i] = S_0(X_i);$$

where $\mathbb{R}_0 = (\mu_0; h_{o1}; \dots; h_{oq})$; $Z = (Y; X_z)$; $Y = (Y_1; Y_2^0)$; $X_z =$ non-overlap union of X_1, \dots, X_q ; and $X = (X_z^0; X_0^0)$. Without loss of generality, we assume that $\dim(X_0) = \dim(Y_2) = b$, and the rest of the dimensions and the assumptions on h_{oj} are the same as those in Example 2.1. It is easy to see that $E[\frac{1}{2}(Z; \mathbb{R}_0) | X] = 0$ where $\frac{1}{2}(Z; \mathbb{R}_0) = Y_1 - Y_2^0 \mu_0 - \sum_{j=1}^q h_j(X_j)$: Since the directional derivative

$$\frac{1}{2}_{\mathbb{R}_0}[Z; \mathbb{R}_0; \mathbb{R}_0] = Y_2^0 [\mu_1 - \mu_0] - \sum_{j=1}^q [h_j(X_j) - h_{oj}(X_j)]$$

depend on the endogenous variable Y_2 , we apply the generalized minimum distance (2.1) to estimate $(\mu_0; h_{01}; \dots; h_{0q})$. By Lemma A.1 in Appendix, we have $\mathfrak{m}(X; \mathbb{R}) \rightarrow m(X; \mathbb{R}) = \int \mathbb{E}[Y_1|X] g + \mu^0 \int \mathbb{E}[Y_2|X] g = o_p(n^{-1/4})$ uniformly over $X \in \mathcal{X}_n; \mathbb{R} \in \mathcal{A}_n$. In (2.1), we let \mathfrak{S} be the identity matrix and H_n^j , $j = 1; \dots; q$; be the wavelet sieve space (2.3). Then by Theorem 3.1 in Section 3, the root-mean square convergence rate for the wavelet sieve estimator of h_{0j} , $j = 1; \dots; q$; will be $o_p(n^{-1/4})$, and by Theorem 4.1, the corresponding estimator of μ_0 will be \sqrt{n} asymptotic normal. Moreover when $\mathfrak{S}(X)$ in (2.1) is a consistent estimator of $S_0(X)$ as provided in Section 6, then by Theorems 6.1 and 6.2 in Section 6, the corresponding estimator of μ_0 will be \sqrt{n} efficient with asymptotic variance V_0^{-1} :

$$V_0 = \inf_{w_j(\cdot)} \mathbb{E} \left[\sum_{j=1}^q w_j(X_j) g^0 [S_0(X)]^{-1} \int \mathbb{E}[Y_2|X] w_j(X_j) g \right]^2 \quad (2.5)$$

where $w = (w_1; \dots; w_q)$ with w_j measurable function of X_j , and $w_j \in B_{p,1}^{m_j}(C_j) \cap \mathcal{H}_{0j}$ for all $j = 1; \dots; q$.

Example 2.3 (Semiparametric sample selection): for $l = 1; 2$

does not depend on the endogenous variable Y , we may apply the simpler sieve generalized least squares regression (2.2) to estimate $(\mu_{01}; \mu_{021}; \mu_{022}; h_{01}; h_{02})$ instead of the method of generalized minimum distance (2.1).

Example 2.4 (Semiparametric IV regression):

$$\begin{aligned} Y_{1i} &= X_{1i}^0 \mu_0 + h_0(Y_{2i}) + \varepsilon_i \\ \mathbb{E}[\varepsilon_i | X_{1i}; X_{2i}] &= 0; \quad \mathbb{E}[\varepsilon_i^2 | X_{1i}; X_{2i}] = S_0(X_{1i}; X_{2i}); \end{aligned}$$

where $\mathbb{R}_0 = (\mu_0; h_0)$; $Z = (Y; X_Z)$; $Y = (Y_1; Y_2)^0$; $X_Z = X_1$ and $X = (X_1^0; X_2)^0$. Without loss of generality, we assume that $\dim(X_1) = b$, $\dim(X_2) = 1$ and $\dim(X) = s = b + 1$, also $\dim(Y_1) = \dim(Y_2) = 1$, $\dim(\mu) = b$ and $\dim(h) = 1$. It is easy to see that $\mathbb{E}[\frac{\partial}{\partial \mu} (Z; \mathbb{R}_0) | X] = 0$ where $\frac{\partial}{\partial \mu} (Z; \mathbb{R}_0) = Y_1 + X_1^0 \mu + h(Y_2)$. Since the directional derivative

$$\frac{\partial}{\partial \mu} (Z; \mathbb{R}_0) | \mathbb{R}_0 = X_1^0 [\mu - \mu_0] + [h(Y_2) - h_0(Y_2)]$$

depends on the endogenous variable Y_2 , we apply the method of generalized minimum distance (2.1) to estimate $(\mu_0; h_0)$. By Lemma A.1 in Appendix, we have $\mathfrak{m}(X; \mathbb{R}) \rightarrow m(X; \mathbb{R}) = \int \mathbb{E}[Y_1|X] g + \int \mathbb{E}[h(Y_2)|X] g - \int \mathbb{E}[h(Y_2)jX] g = o_p(n^{-1/4})$ uniformly over $X \in \mathcal{X}_n; \mathbb{R} \in \mathcal{A}_n$. In (2.1), we let \mathfrak{S} be the identity

matrix and H_n be the sieve. Then by Theorem 3.1 in Section 3, the convergence rate in certain metric (to be specified later) for the sieve estimator of h_0 will be $o_p(n^{-1/4})$, and by Theorem 4.1, the corresponding estimator of μ_0 will be $\sqrt{p/n}$ asymptotic normal. Moreover when $\hat{S}(X)$ in (2.1) is a consistent estimator of $S_0(X)$ as provided in Section 6, then by Theorems 6.1 and 6.2 in Section 6, the corresponding estimator of μ_0 will be $\sqrt{p/n}$ efficient with asymptotic variance V_0 :

$$V_0 = \inf_w E \int f(X) \left[E[w(Y_2) | X] g[S_0(X)] \right]^2 f(X) \left[E[w(Y_2) | X] g \right]^2 \quad (2.6)$$

where w is a measurable function of Y_2 , and $w \in B_{p,1}^m(C)$.

3. Consistency

In this section, we first present a set of sufficient conditions, then discuss those conditions, and finally state the theorem on consistency and rates of convergence of $\hat{\theta}_n$. We begin by introducing notations which will aid presentation of sufficient conditions. Let $\|\cdot\|_2$ denote a pseudo-norm on A , and $N(A_n; \epsilon)$ the minimum bracket-number of radius ϵ balls that cover the sieve space A_n . In this paper, we require the pseudo-norm $\|\cdot\|_2$ to be locally equivalent to $K(\mathbb{R}_0; \mathbb{R})$, where

$$K(\mathbb{R}_0; \mathbb{R}) = E \int f(X) \int X_n g(m(X; \mathbb{R})) \hat{S}(X) \int m(X; \mathbb{R}) g = 2:$$

That is, there exist two constants $c_1, c_2 > 0$ such that $c_1 \frac{1}{K(\mathbb{R}_0; \mathbb{R})} \leq \|\cdot\|_2 \leq c_2 \frac{1}{K(\mathbb{R}_0; \mathbb{R})}$ hold for any \mathbb{R} in a neighborhood of \mathbb{R}_0 . For example, when $\frac{1}{2}(Z; \mathbb{R})$ is Lipschitz continuous with respect to \mathbb{R} , we can choose the pseudo-metric $\|\cdot\|_2$ to be $E[(\mu - \mu_0) \int \mu(Z)(\mu - \mu_0)] + \sum_{j=1}^q E[\int h_j(Z) f h_j(\mathbb{R}) - h_{0j}(\mathbb{R}) g^2]$ for some known weighting functions $\int \mu(Z)$ and $\int h_j(Z)$. To avoid confusion, all vectors refer to column vectors unless stated otherwise.

An estimator $\hat{\theta}_n$ is said to be consistent for \mathbb{R}_0 if $\|\hat{\theta}_n - \mathbb{R}_0\|_2 \rightarrow 0$ in probability, denoted as $\|\hat{\theta}_n - \mathbb{R}_0\|_2 = o_p(1)$. Let $f_{\pm n}$ denote a positive sequence that decreases to zero as $n \rightarrow \infty$. $\hat{\theta}_n$ is said to be consistent for \mathbb{R}_0 at a rate (strictly) faster than $f_{\pm n}$ if $\|\hat{\theta}_n - \mathbb{R}_0\|_2 = o_p(f_{\pm n})$. Our goal in this section is to show that $\hat{\theta}_n$ is consistent at a rate faster than $n^{-1/4}$. To attain this goal, we need conditions to ensure that H_n converges to H , $m(X; \mathbb{R})$ converges to $m(X; \mathbb{R}_0)$ uniformly over X and \mathbb{R} , and $\hat{S}(X)$ converges to $S(X)$ uniformly over X at rates faster than $n^{-1/4}$. We also need conditions on the local curvature of $\frac{1}{2}(\mathbb{R}; \mathbb{R})$. Below we present such conditions.

Assumption 1. The data $f(Y_i; X_i) : i = 1; 2; \dots; n$ are i.i.d..

Assumption 2. $E[\frac{1}{2}(Z; \theta)jX] = 0$ holds for any $\theta \in A$ $\|k_{\theta}\|_2 = 0$.

Assumption 3. (i) $\hat{S}(X) = S(X) + o_p(n^{-1/4})$ uniformly over $X \in X_n$; and (ii) there exist some positive constants c_1 and c_2 such that $c_1 \min_{\lambda} (S(X)) \leq \max_{\lambda} (S(X)) \leq c_2$ for all $X \in X$, where $\lambda(S(X))$ denote eigenvalues of $S(X)$.

Assumption 4. (i) The density f_X is strictly positive on the compact support X ; and (ii) f_X has up to J -th derivatives, where $J \geq 2$.

Assumption 5. $K(x)$: (i) is a bounded function with bounded support, and Lipschitz continuous; (ii) is symmetric around origin, and integrates to one; and (iii) is of order $J \geq 2$, i.e.

$$\int_{\mathbb{R}^d} (x^1)^{j_1} \dots (x^s)^{j_s} K(x) dx = 0 \text{ for } 1 \leq j_1 + \dots + j_s \leq J-1;$$

$$\int_{\mathbb{R}^d} (x^1)^{j_1} \dots (x^s)^{j_s} K(x) dx \neq 0 \text{ for } j_1 + \dots + j_s = J;$$

where x^i denotes the i -th element of x .

Assumption 6. $E[\frac{1}{2}(Z; \theta)jX]$ has up to J -th (total) derivatives with respect to X ; and the J -th derivative (w.r.t. X) is uniformly bounded in $(X; \theta) \in X \times A_n$.

Assumption 7. (i) There exists a measurable function $\hat{A}_1 : Z \rightarrow [0; 1)$ such that $E[\hat{A}_1(Z)^p] < 1$ for some integer $p \geq 2$ and $\frac{1}{2}(Z; \theta)j \hat{A}_1(Z)$ for all $\theta \in A_n$; and (ii) there exist a sequence of positive numbers c_n and a measurable function $\hat{A}_2 : Z \rightarrow [0; 1)$ such that $E[\hat{A}_2(Z)jX] < 1$ and

$$\frac{1}{2}(Z; \theta_1)j \frac{1}{2}(Z; \theta_2)j \leq c_n \hat{A}_2(Z) \|k_{\theta_1} - k_{\theta_2}\|_2 \text{ for all } \theta_1, \theta_2 \in A_n;$$

Assumption 8. (i) $\frac{c_n \log N(A_n; c_n^{-1} n^{1/4}) E a_n^{s(s+1)} n^{s-4} E \max_{1 \leq i \leq p} E a_n^{i^2(s+1)=p}}{n^{p-1}} \rightarrow 0$ as $n \rightarrow \infty$,

and (ii) $n a_n^{4J} \rightarrow 0$ as $n \rightarrow \infty$.

Assumption 9. For any $\theta \in A$, there exists $\theta_n \in A_n$ such that $\|k_{\theta} - k_{\theta_n}\|_2 = o(n^{-1/4})$ as $n \rightarrow \infty$.

Assumption 10. $\alpha_n = o(n^{-1/4})$, where $\alpha_n \in (0; 1]$ denote the solution to:

$$\inf_{\alpha \in (0; 1]} \frac{1}{2} \int_{\mathbb{R}^d} \frac{r^2}{\log N(A_n; \alpha^{-1} n^{1/4})} \frac{1}{\alpha^{p-1}} d\mu \leq O(\alpha^{-p/n});$$

In the following, we use $\mathcal{B}_o(\delta) = \{Z \in \mathbb{R}^n : \|Z - \theta_o\| \leq \delta\}$ to denote a (radius δ) neighborhood of θ_o , and $U_n = \{Z \in \mathbb{R}^n : \|Z - \theta_o\| \leq o(n^{-1/8})\}$ the unit circle.

Assumption 11. (i) $E f_{\mathcal{B}_o}^j[Z; \hat{A}]X$ has up to J -th (total) derivatives w.r.t X ; and the J -th derivative (w.r.t X) is uniformly bounded over $X \in \mathbb{R}^n$, and $\hat{A} \in U_n$; (ii) there exists a measurable function $\hat{A}_3(Z)$ such that $E[\hat{A}_3(Z)^p]$ is finite and $j_{\mathcal{B}_o}^j[Z; \hat{A}] - \hat{A}_3(Z)$ holds for any $Z \in \mathbb{R}^n$ and $\hat{A} \in U_n$; and (iii) there exist a sequence of positive numbers c_n and a measurable function $\hat{A}_4(Z)$ such that $E[\hat{A}_4(Z)X]$ is finite and

$$j_{\mathcal{B}_o}^j[Z; \hat{A}_1] - j_{\mathcal{B}_o}^j[Z; \hat{A}_2] \leq c_n \hat{A}_4(Z) \|\hat{A}_1 - \hat{A}_2\|_2$$

holds for any $Z \in \mathbb{R}^n$ and $\hat{A}_1, \hat{A}_2 \in U_n$.

Define the linear approximation error by

$$r_{\mathcal{B}_o}^j[Z; \theta_o] = j_{\mathcal{B}_o}^j[Z; \theta_o] - j_{\mathcal{B}_o}^j[Z; \theta_o]$$

Assumption 12. (i) there exists a measurable function $\hat{A}_5(Z)$ such that $E[\hat{A}_5(Z)^p]X$ exists and

$$j_{\mathcal{B}_o}^j[Z; \theta_o] - \hat{A}_5(Z) \|\theta_o\|_2^2$$

holds for any $Z \in \mathbb{R}^n$ and $\theta_o \in \mathcal{B}_o(o(n^{-1/8}))$;

(ii) there exist a sequence of positive numbers c_n and a measurable function $\hat{A}_6(Z)$ such that $E[\hat{A}_6(Z)X]$ is finite and

$$j_{\mathcal{B}_o}^j[Z; \theta_1] - j_{\mathcal{B}_o}^j[Z; \theta_2] \leq c_n \hat{A}_6(Z) \|\theta_1 - \theta_2\|_2$$

holds for any $Z \in \mathbb{R}^n$ and $\theta_1, \theta_2 \in \mathcal{B}_o(o(n^{-1/8}))$.

Theorem 3.1. Suppose Assumptions 1-12 are satisfied. If $\hat{\theta}_n$ is the solution to (2.1), then $\|\hat{\theta}_n - \theta_o\|_2 = o_p(n^{-1/4})$:

Corollary 3.2. Suppose Assumptions 1-3, 7(ii), 9-10 are satisfied, and that the directional derivative $j_{\mathcal{B}_o}^j[Z; \theta_o]$ does not depend on the variables Y . If $\hat{\theta}_n$ is the solution to (2.2), then $\|\hat{\theta}_n - \theta_o\|_2 = o_p(n^{-1/4})$:

Discussions: (1) Assumptions 1, 2 and 3 are basic regularity conditions. Assumption 1 rules out dependent data. This condition, though restrictive, is not critical for our results on consistency of $\hat{\theta}$ and $\hat{\theta}_n$ asymptotic normality of $\hat{\theta}$.

These results still can be proved for weakly dependent data using the approach developed here and those developed in Chen and Shen (1998). The independence assumption is, however, critical for the semiparametric efficiency result. Assumption 2 is a global identification condition, and must be satisfied for our results to hold. Verification of this condition requires the knowledge of structures of specific models. Note that this condition only identifies μ_0 up to the equivalence class with respect to the norm $\|\cdot\|_2$. In Section 4, we shall impose additional condition to identify μ_0 uniquely. Assumption 3(i) requires that the estimator of the weighting matrix converges to some positive definite matrix uniformly over the regressors at a rate faster than $n^{-1/4}$. This assumption is not restrictive and can be satisfied by many estimators. For example, the identity weighting matrix satisfies this condition. The kernel estimator proposed in Section 6 also satisfies this condition. Assumption 3(ii) requires the weighting matrix to be bounded from above and below. Conditions of this sort are common in the weighted Least Squares regression literature.

(2) Assumptions 4-8, 11 and 12 are conditions related to the kernel estimation. They (except for 7(ii)) should be dropped whenever the initial estimation of $m(X; \theta)$ is not needed. Assumptions 4 - 8 imply that the kernel estimator $\hat{m}(X; \theta)$ converges to $m(X; \theta)$ at a rate faster than $n^{-1/4}$ uniformly over $X \in X_n; \theta \in \Theta_n$ (see Corollary A.2 in the Appendix A). These conditions together with Assumption 3(i) imply that the kernel-estimated sample criterion function in (2.1) converges to $\frac{1}{n} \sum_{i=1}^n 1_{\{X_i \in X_n\}} [S(X_i)]^{-1} m(X_i; \theta)$ uniformly over θ at a rate faster than $n^{-1/4}$. Under additional Assumptions 11 and 12, the convergence rate at a local neighborhood of true θ_0 is even faster, (see the proof of Theorem 3.1 in the Appendix B for details).

Assumption 4(i) rules out unbounded regressors and zero density values. This condition can be dropped when trimming is used to control for small density values of f_X , see e.g. Andrews (1995). Assumption 4(ii) is a smoothness condition. It is not critical for the consistency of $\hat{\theta}_n$ but critical for the \sqrt{n} consistency of $\hat{\beta}_n$. Nevertheless, it is familiar in the semiparametric literature. Assumption 5 also is familiar in the semiparametric econometric literature. It is satisfied by many commonly used kernel functions. Assumption 6 requires $m(\cdot)$ to be sufficiently smooth with respect to X . Conditions of this sort are familiar in the semiparametric literature, (e.g. Powell et al. (1989), Horowitz (1996)). Again, this condition is not critical for the consistency of $\hat{\theta}_n$ but critical for the \sqrt{n} consistency of $\hat{\beta}_n$. Assumption 7(i) is an envelope condition which resembles the envelope condition commonly imposed in the parametric literature. The envelope condition plays

an important role in the parametric asymptotics literature and it is expected to play a similar role in this paper. Assumption 7(ii) is a Lipschitz condition which requires that the derivative satisfy an envelope condition. Assumption 8 imposes restrictions on the rates at which the kernel bandwidth $a_n \rightarrow 0$. Assumption 8(i) controls the variance component of the kernel estimator and Assumption 8(ii) controls the bias component, both converging at a rate faster than $n^{-1/4}$ uniformly over $X \times X_n$. Assumptions 11 and 12 are conditions on the local curvature of $\frac{1}{2}(\psi; \psi)$. Assumption 11 implies that, on the unit circle centered at θ_0 , the directional derivative function is smooth, satisfies an envelope condition and is Lipschitz continuous. While Assumption 12 requires that, in a neighborhood of θ_0 , the linear approximation error is sufficiently smooth, bounded by a quadratic function of $\theta - \theta_0$, and Lipschitz continuous.

(3) Assumptions 7, 9 and 10 are conditions related to the sieve estimation of θ_0 . Assumption 9 requires that the sieve approximation error must converge to zero faster than $n^{-1/4}$, while Assumption 10 requires that the size of the sieves A_n should not be too large. To verify these conditions, we must compute the sieve approximation error $k_{\theta_0}(\cdot) - \inf_{\theta \in A_n} k_{\theta}(\cdot)$ and the bracketing number $N_{[]}(\cdot, \|\cdot\|_2, A_n, \epsilon)$. While both depend on the smoothness and dimensionality of $A_n \subset \mathcal{E} \subset H_n$, $k_{\theta_0}(\cdot) - \inf_{\theta \in A_n} k_{\theta}(\cdot)$ also depends on the structure of $A_n \subset \mathcal{E} \subset H$. Let $\|\cdot\|_2$ denote the Euclidean norm on \mathcal{E} and $d(h; h_0)$ a distance metric on H such that the product metric $\|\cdot\|_2 \times d(h; h_0)$ is locally equivalent to $k_{\theta_0}(\cdot) - \inf_{\theta \in A_n} k_{\theta}(\cdot)$ on A . Let $N(H_n; \epsilon)$ denote the minimum bracketing number of radius ϵ balls (in d distance) covering the sieve space H_n . Then obviously $k_{\theta_0}(\cdot) - \inf_{\theta \in A_n} k_{\theta}(\cdot) \leq d(h; \inf_{\theta \in A_n} h)$; the approximation error rate between $h \in H$ and its corresponding sieve $\inf_{\theta \in A_n} h \in H_n$. And $N(A_n; \epsilon) \leq N(H_n; \epsilon) \leq \epsilon^{-b}$; where ϵ^{-b} is the minimum number of radius ϵ balls (in Euclidean distance) covering \mathcal{E} . For many sieves, calculations of $d(h; \inf_{\theta \in A_n} h)$ and $N(H_n; \epsilon)$ can be found in DeVore and Lorentz (1993) and Van der Vaart and Wellner (1996). Such calculations can also be found in econometrics applications such as Fenton and Gallant (1996) for Hermite polynomials, Newey (1997) for spline and power series, Chen, Hansen and Scheinkman (1997) for shape-preserving wavelet B-spline with unbounded support, and Chen and White (1999) for multivariate neural networks.

Theorem 3.1 shows that the proposed estimator is consistent at a rate faster than $n^{-1/4}$. This result is stronger than the simple consistency $\| \hat{\theta}_n - \theta_0 \|_2 = o_p(1)$. The stronger result is needed for obtaining the L^2 consistency of $\hat{\beta}_n$, to which we now turn to.

4. Asymptotic Normality

The standard approach for deriving the asymptotic distribution is to assume that $\Theta_n = (\beta_n; \mathbf{h}_n)$ is an interior solution and then to apply Taylor expansion. This approach may not work well here for two reasons. First, while it is reasonable to expect β_n to be an interior solution, \mathbf{h}_n often is not and hence the first order condition may not hold. Second, even if the first order condition is satisfied, we need to compute the projection of the derivatives with respect to μ onto the derivatives with respect to h . Because h is infinite dimensional, computing such projection is often complicated and may not have a close-form solution. An alternative approach is to apply the Riesz representation theorem. The Riesz representation theorem links β_n to the derivatives of the sample criterion function without relying on the first order condition and hence does not have the problems mentioned above. This approach was used in Shen (1997) and Chen and Shen (1998) and is adopted in this paper.

To derive the asymptotic distribution of β_n , it suffices to derive the asymptotic distribution of $f(\Theta_n) = \beta_n$ for any fixed non-zero $\beta \in \mathbb{R}^b$. Let

$$m_{\beta_0}[X; \Phi] = E[\beta_0[Z; \Phi]X]$$

denote the conditional expectation of the directional derivative. Define a Hilbert-norm on A as

$$\|h\|_{\beta_0}^2 = E[m_{\beta_0}[X; \beta_0]^\top \Sigma(X)^{-1} m_{\beta_0}[X; \beta_0] h]$$

Let \bar{A} denote $E \in \bar{H}$ with \bar{H} the linear completion of H under the norm $\|h\|_{\beta_0}$, and let $\langle h, i \rangle$ denote the inner product induced by the norm $\|h\|_{\beta_0}$ on \bar{A} . Since $f(\beta)$ is a linear functional on \bar{A} , it is bounded (i.e. continuous) if and only if

$$\sup_{\|h\|_{\beta_0} > 0, h \in \bar{A}} \frac{|f(\beta) - f(\beta_0)|}{\|h\|_{\beta_0}} < \infty :$$

The Riesz representation theory states that there exists a representer $v \in \bar{A}$ satisfying $\|v\|_{\beta_0} = \sup_{\|h\|_{\beta_0} > 0, h \in \bar{A}} \frac{|f(\beta) - f(\beta_0)|}{\|h\|_{\beta_0}}$ and

$$f(\beta) = f(\beta_0) + \langle h, v \rangle_{\beta_0} :$$

Hence, $f(\Theta_n) - f(\beta_0) = \langle \mathbf{h}_n, v \rangle_{\beta_0}$. Our aim in this section is to compute the representer v and then link $f(\Theta_n) - f(\beta_0)$ to the pathwise directional derivatives of the sample criterion function.

First, we compute the representor v^α . Let $W = H \times h_0$ and let \overline{W} denote the linear completion of W under the Hilbert norm. For any $h \in \overline{H}$, there exist $w^j(t) \in \overline{W}$ for $j = 1; \dots; b$ such that $h \times h_0 = \sum_{j=1}^b (w^j(t); t) (\mu_j - \mu_0) = \sum_{j=1}^b w^j(t) (\mu_j - \mu_0)$. Let $\frac{1}{2}_h[Z; \Phi h]$ denote the directional derivative with respect to h at direction Φh . Define $\frac{1}{2}_{h_0}[Z; w(t)] = (\frac{1}{2}_{h_0}[Z; w^1(t)]; \dots; \frac{1}{2}_{h_0}[Z; w^b(t)])$ and $D_w(X)$ the $r \times b$ matrix valued function with the j th column $D_{w^j}(X)$; $j = 1; \dots; b$ defined as $D_{w^j}(X) = E \frac{1}{2}_{\mu_j}(Z; \otimes_0) \times \frac{1}{2}_{h_0}[Z; w^j(t)] X$. Then

$$D_w(X) = E \frac{1}{2}_{\mu^0}(Z; \otimes_0) \times \frac{1}{2}_{h_0}[Z; w(t)] X \quad (4.1)$$

where $\frac{1}{2}_{\mu^0}(Z; \otimes)$ is the ordinary derivative with respect to $\mu^0 = (\mu_1; \dots; \mu_b)$. By definition we have $m_{\otimes_0}[X; \otimes_0] = D_w(X)(\mu_j - \mu_0)$, and

$$k_{\otimes_0}^2 = (\mu_j - \mu_0)^0 E \frac{1}{2}_{h_0}[Z; w(t)] X^i D_w(X)^i (\mu_j - \mu_0)$$

Let $w^\alpha(t) = (w^{\alpha 1}(t); \dots; w^{\alpha b}(t))$ be the solution to

$$\inf_{w^j(t) \in \overline{W}; j=1; \dots; b} E [D_w(X)^0 S(X)^i D_w(X)^i] \quad (4.2)$$

where "inf" is in positive semidefinite matrix sense. Simple calculation yields

$$\begin{aligned} kv^\alpha k^2 &= \sup_{\mu_i \mu_0 \in \mathcal{O}} \frac{(\mu_j - \mu_0)^0 (\mu_j - \mu_0)}{(\mu_j - \mu_0)^0 [\inf_w E [D_w(X)^0 S(X)^i D_w(X)^i]] (\mu_j - \mu_0)} \\ &= \frac{1}{\inf_w E [D_w^\alpha(X)^0 S(X)^i D_w^\alpha(X)^i]} \end{aligned} \quad (4.3)$$

Clearly, $kv^\alpha k < 1$ if and only if $E [D_w^\alpha(X)^0 S(X)^i D_w^\alpha(X)^i]$ is positive-definite. The Riesz representor $v^\alpha = (v_\mu^\alpha; v_h^\alpha)$ is now given by

$$\begin{aligned} v_\mu^\alpha &= E \frac{1}{2}_{h_0}[Z; w^\alpha(t)] X^i D_w^\alpha(X)^i S(X)^i D_w^\alpha(X)^i \\ v_h^\alpha &= \sum_{j=1}^b w^{\alpha j}(t) E \frac{1}{2}_{\mu_j}(Z; \otimes_0) \times \frac{1}{2}_{h_0}[Z; w^j(t)] X^i D_w^\alpha(X)^i S(X)^i D_w^\alpha(X)^i \end{aligned}$$

Hence, $m_{\otimes_0}[X_i; v^\alpha] = D_w^\alpha(X) (E [D_w^\alpha(X)^0 S(X)^i D_w^\alpha(X)^i])^{-1}$

Next, we link $f(\mathcal{O}_n) \times f(\otimes_0)$ to the pathwise directional derivatives of the sample criterion function. After some calculation in the Appendix C, we obtain:

$$P_n^-(f(\mathcal{O}_n) \times f(\otimes_0)) = \sum_{i=1}^n \frac{1}{n} X^i m_{\otimes_0}[X_i; v^\alpha]^0 [S(X_i)]^i \frac{1}{2}(Z_i; \otimes_0) + o_p(1) \quad (4.4)$$

Applying the Lindeberg-Levy central limit theorem proves that $\sqrt{n}(\hat{\beta}_n - \beta_0)$ is asymptotically normally distributed with mean zero and variance V^{-1} , where

$$V = \frac{1}{n} \left(E[D_{w^a}(X)^0 [S(X)]^i {}^1 D_{w^a}(X) g^a] \right)^2 + \frac{1}{n} \left(E[D_{w^a}(X)^0 [S(X)]^i S_0(X) [S(X)]^i {}^1 D_{w^a}(X) g^a] \right)^2 + \frac{1}{n} \left(E[D_{w^a}(X)^0 [S(X)]^i {}^1 D_{w^a}(X) g^a] \right)^2 \quad (4.5)$$

and $S_0(X) = \text{Var}[\frac{1}{2}(Z; \beta_0) | X]$. Hence, $\hat{\beta}_n$ is \sqrt{n} consistent and asymptotically normally distributed with covariance V^{-1} .

We now present a set of conditions in addition to those in Section 3 to ensure validity of equation (4.4).

Assumption 13. There exists a measurable function $\hat{A}_7(Z)$ such that $E[\hat{A}_7(Z) | X]$ is finite and

$$\hat{A}_7(Z) = k_{\beta_1} | \beta_1 - \beta_0 | + k_{\beta_2} | \beta_2 - \beta_0 |$$

holds for any $Z \in \mathcal{Z}$ and $\beta_1, \beta_2 \in \mathcal{B}_0(o(n^{1/4}))$:

Assumption 14. (i) $E[D_{w^a}(X)^0 S(X)^i {}^1 D_{w^a}(X)]$ is finite and positive-definite; and (ii) $S_0(X)$ is finite positive definite uniformly over $X \in \mathcal{X}$.

Assumption 13 is similar to but not implied by Assumption 12. It requires the moment functions $\frac{1}{2}(Z; \beta)$ to have an extra curvature in the neighborhood of β_0 . Assumption 14 is a local identification condition of μ_0 . Assumptions 2, 3 and 14 together uniquely identify μ_0 . Moreover, they ensure that μ_0 is \sqrt{n} consistently estimated as we show in the next theorem.

Theorem 4.1. Suppose Assumptions 1-14 hold. Then $\sqrt{n}(\hat{\beta}_n - \beta_0)$ is asymptotically normally distributed with mean zero and covariance matrix V^{-1} .

Theorem 4.1 shows that the proposed estimator for the parametric component (μ_0) attains the parametric rates. It provides an asymptotic distribution for the parametric component, which is needed for statistical inference. The asymptotic covariance has three terms. If we can interpret $D_{w^a}(X)$ as regressors in a system of equations with $\frac{1}{2}\mu_0(Z; \beta_0) + \frac{1}{2}w^a$ as dependent variables, then V^{-1} has an interpretation of the covariance of the weighted minimum distance estimator. The three terms in V correspond to the three terms in the covariance of the (incorrectly) weighted LS estimator in the presence of heteroskedasticity.

As we noted in the introduction, the nonparametric component (h_0) can be concentrated out analytically in some applications such as the partial linear regression (Robinson, 1988) and the single index model (Ichimura, 1993). Although the proposed procedure and Theorems 3.1 and 4.1 are applicable to those applications, it may be worth taking advantage of the concentrated out nonparametric component. Let $h^\mu(\mathfrak{t}; \mu)$ denote the concentrated out nonparametric component, which satisfies $h^\mu(\mathfrak{t}; \mu_0) = h_0(\mathfrak{t})$, and often depends on some unknown functions. Suppose that, for each fixed μ , $h^\mu(\mathfrak{t}; \mu)$ is estimated by some nonparametric estimator $\hat{h}^\mu(\mathfrak{t}; \mu)$ using kernel or sieve methods. The estimator, $\hat{\beta}_n$; for the parametric component is then obtained as the solution to

$$\min_{\mu \in \mathcal{E}} \sum_{i=1}^n 1(X_i \in X_n) \mathfrak{m}(X_i; \mu; \hat{h}^\mu(\mathfrak{t}; \mu))^\top [\hat{S}(X_i)]^{-1} \mathfrak{m}(X_i; \mu; \hat{h}^\mu(\mathfrak{t}; \mu));$$

Suppose that $\hat{h}^\mu(\mathfrak{t}; \mu)$ is differentiable with respect to μ and that $\hat{h}^\mu(\mathfrak{t}; \mu)$ and its derivatives with respect to μ converge to the true values in probability uniformly at rates faster than $n^{-1/4}$. Then Theorem 4.1 can also be proved using the techniques developed in Andrews (1994), Newey (1994) or Pakes and Olley (1995).

5. Covariance Estimator

To conduct statistical inference on the parametric component, a consistent covariance estimator is needed. In this section, we provide such an estimator. We first present consistent estimators of $S(X)$; $S_0(X)$; and $D_{w^\mu}(X)$ respectively, and then present a consistent estimator of V . Consistent estimator of $S(X)$ already exists. To estimate $S_0(X)$, we use the kernel method:

$$\hat{S}_0(X_i) = \frac{[(n_i - 1)a_n^s]^{-1} \sum_{j \in \mathcal{I}; j=1}^n \frac{1}{2} (Z_j; \mathfrak{D}_n) \frac{1}{2} (Z_j; \mathfrak{D}_n)^\top K \left(\frac{X_i - X_j}{a_n} \right)}{\hat{f}_{X_i}};$$

To estimate $D_{w^\mu}(X)$; we estimate $w^\mu(\mathfrak{t})$. Recall that $D_{w^k}(X)$ and $w^k(\mathfrak{t})$ are the k -th columns of $D_w(X)$ and $w(\mathfrak{t})$ respectively, and that $w^k(\mathfrak{t})$; $k = 1; \dots; b$ is the solution to:

$$\inf_{w^k(\mathfrak{t}) \in \mathcal{W}} \mathbb{E}^h [D_{w^k}(X)^\top [S(X)]^{-1} D_{w^k}(X)];$$

Hence, $w^k(\mathfrak{t})$ can be estimated by applying a nonparametric least squares regres-

sion. Specifically, let μ_k denote the k -th element of μ and let

$$D_{w^k}(X_i) = \frac{\frac{1}{n} \sum_{j \in \mathcal{I}_n} \frac{\partial \frac{1}{2}(Z_j; \mathbf{b}_n)}{\partial \mu_k} \frac{1}{h_n} [Z_j; w^k(\mathbf{t})] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n}}{f_{X_i}}$$

denote the kernel estimator of $D_{w^k}(X_i)$: Let \overline{W}_n denote the linear completion of \mathcal{I}_n , then \overline{W}_n is a sieve approximation to \overline{W} . We now estimate $w^k(\mathbf{t})$ by $\hat{w}^k(\mathbf{t})$, the solution to

$$\min_{w^k(\mathbf{t}) \in \overline{W}_n} \sum_{i=1}^n 1(X_i \in \mathcal{I}_n) \mathbf{D}_{w^k}(X_i) [\hat{w}^k(X_i)]^2 \mathbf{D}_{w^k}(X_i);$$

Let \hat{W} denote the matrix whose k -th column is $\hat{w}^k(\mathbf{t})$. Then, $D_{w^k}(X)$ is estimated by $\hat{D}_{w^k}(X)$ and V is estimated by

$$\hat{V} = \frac{1}{n} \sum_{i=1}^n \frac{1}{f_{X_i}} \frac{1}{h_n} \sum_{j \in \mathcal{I}_n} \frac{\partial \frac{1}{2}(Z_j; \mathbf{b}_n)}{\partial \mu_k} \frac{1}{h_n} [Z_j; w^k(\mathbf{t})] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n} \frac{1}{f_{X_i}} \frac{1}{h_n} \sum_{j \in \mathcal{I}_n} \frac{\partial \frac{1}{2}(Z_j; \mathbf{b}_n)}{\partial \mu_k} \frac{1}{h_n} [Z_j; w^k(\mathbf{t})] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n}$$

Assumption 15. (i) There exists a measurable function $\hat{A}_8(Z)$ such that $E[\hat{A}_8(Z)jX]$ exists and is finite, and $j \frac{\partial \frac{1}{2}(Z; \mathbf{b}_n)}{\partial \mu_k} \frac{1}{h_n} [Z; w^k(\mathbf{t})] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n} \hat{A}_8(Z) \leq k \leq b$; $Z \in \mathcal{Z}$, $\mathbf{b}_n \in \mathcal{B}_0(o(n^{1/4}))$; and (ii) there exists a measurable function $\hat{A}_9(Z)$ such that $E[\hat{A}_9(Z)jX]$ exists and is finite, and $j \frac{1}{h_n} [Z; v(\mathbf{t})] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n} \hat{A}_9(Z) \leq k \leq b$; $Z \in \mathcal{Z}$, $v \in \overline{W}_n$ with $\|v\|_2 = 1$, $\mathbf{b}_n \in \mathcal{B}_0(o(n^{1/4}))$.

Assumption 16. (i) there exists a measurable function $\hat{A}_{10}(Z)$ such that $E[\hat{A}_{10}(Z)^p]$ is finite and $j \frac{1}{h_n} [Z; \mathbf{b}_n] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n} \hat{A}_{10}(Z) \leq k \leq b$ and $\mathbf{b}_n \in \mathcal{B}_0(o(n^{1/4}))$; (ii) there exists a measurable function $\hat{A}_{11}(Z)$ such that $E[\hat{A}_{11}(Z)jX]$ is finite and $j \frac{1}{h_n} [Z; \mathbf{b}_1] \mathbf{1}_{\mathcal{I}_n} \frac{X_{ij} X_j}{a_n} \hat{A}_{11}(Z) \leq k \leq b$ and $\mathbf{b}_1 \in \mathcal{B}_0(o(n^{1/4}))$.

Theorem 5.1. Under Assumptions 1-16, we have $\hat{V} = V + o_p(1)$:

It is worth noting that, for linear sieves, computing $\hat{w}^k(\mathbf{t})$ and hence the covariance estimator does not require nonlinear optimizations. A simple pooled and weighted LS regression of the derivative of $\frac{1}{2}$ with respect to μ_k on the derivatives of $\frac{1}{2}$ with respect to h is all we need. Thus, our covariance estimator is easy to compute.

6. Efficiency

Having established the asymptotic distribution of the parametric component, we now investigate the efficiency of $\hat{\beta}_n$. Clearly, the efficiency depends on the choice of the weighting matrix $S(X)$. It is straightforward to show that the optimal weighting matrix in the sense of minimizing the asymptotic variance of the estimator $\hat{\beta}_n$ is $S(X) = S_0(X) \sim \text{Var}(\frac{1}{2}(Z; \otimes_0)jX)$. Thus, when $S(X)$ is set to $S_0(X)$; the corresponding estimator is the best in the class of all weighted minimum distance estimators and hence called the optimally weighted minimum distance estimator. The asymptotic covariance of the optimally weighted minimum distance estimator is V_0^{-1} , where

$$V_0 = E \left[\frac{h}{D_{w_0}(X)} \left[S_0(X) \right]^{-1} \frac{i}{D_{w_0}(X)} \right] \quad (6.1)$$

$$= \inf_{w(\cdot)} E \left[\frac{h}{D_w(X)} \left[S_0(X) \right]^{-1} \frac{i}{D_w(X)} \right] : \quad (6.2)$$

The question now is whether the optimally weighted sieve minimum distance estimator is also the best in the class of all regular and \sqrt{n} consistent estimators for model (1.3). To address this question, one would normally have to compute the semiparametric efficiency bound of model (1.3) and then compare the semiparametric efficiency bound with the asymptotic covariance matrix of the optimally weighted sieve minimum distance estimator. Unfortunately, in our general setup (1.3), there is no analytical expression of the semiparametric efficiency bound, see Chamberlain (1992) for an example. Fortunately, we do not need the explicit expression of the bound and are still able to present the following efficiency result:

Theorem 6.1. Suppose Assumptions 1-14 hold with $S = S_0$. Then the optimally weighted sieve minimum distance estimator for the parametric component in model (1.3) attains the semiparametric efficiency bound V_0^{-1} .

Theorem 6.1 shows that the optimally weighted sieve minimum distance estimator is efficient. It also shows that the V_0^{-1} defined in (6.1) characterizes the semiparametric efficiency bound of model (1.3). While the former result improves upon the existing literature on semiparametric efficient estimation, the latter result extends Chamberlain's (1992) which computes the semiparametric efficiency bound for a special case of (1.3) where $h_0 = (h_{01}(\pm_1(Z)); \dots; h_{0q}(\pm_q(Z)))^0$, $\pm_j(Z)$, $j = 1; \dots; q$ are completely known.

It is known that the generalized sieve minimum distance estimator is efficient for a parametric regression model when the variance of the error term is used as the weight. For a more general parametric conditional moment restrictions model (1.1), the GMM estimator is efficient when the moment functions are weighted by the covariance matrix of the moment functions (Newey 1990, 1993). For a semiparametric conditional moment restrictions setting (1.3), one would expect the same efficiency result to hold when the covariance of the moment functions is used as the weighting matrix. Theorem 6.2 simply confirms this intuition.

To obtain the efficient estimator, we need the weighting matrix $S_0(X)$, which can be estimated by the kernel estimator $\hat{S}_0(X)$ if an initial consistent estimator $\hat{\theta}_n$ is available. The following three-step procedure describes how to obtain the efficient estimator:

Step 1. Obtain an initial consistent estimator $\hat{\theta}_n$ as the solution to:

$$\min_{\theta \in \Theta_n} \frac{1}{n} \sum_{i=1}^n 1' f(X_i; \theta) \Sigma_n^{-1} [f(X_i; \theta)] \Sigma_n^{-1} f(X_i; \theta):$$

Step 2. Obtain a consistent estimator of the optimal weighting matrix by

$$\hat{S}_0(X_i) = \frac{[(n-1)a_n]^{-1} \sum_{j \in I; j=1}^n \frac{1}{2} (Z_j; \hat{\theta}_n) \frac{1}{2} (Z_j; \hat{\theta}_n)' K \left(\frac{X_i - X_j}{a_n} \right)}{f(X_i)}:$$

Step 3. Obtain the optimally weighted estimator $\hat{\theta}_n = (\hat{\beta}_n; \hat{\eta}_n)$ by solving

$$\min_{(\mu, h) \in \Theta_n} \frac{1}{n} \sum_{i=1}^n 1' f(X_i; \theta) \Sigma_n^{-1} \hat{S}_0(X_i)^{-1} f(X_i; \theta):$$

From Theorem 3.1, the initial estimator in Step 1, $\hat{\theta}_n$, converges to θ_0 at a rate faster than $n^{-1/4}$. To apply Theorem 4.1 to $\frac{1}{n} \sum_{i=1}^n \beta_n - \mu_0$, we need to show the estimated weighting matrix in Step 2 converges in probability to $S_0(X_i)$ uniformly over $X_i \in X_n$ at a rate faster than $n^{-1/4}$. This uniform convergence can be proved with the following additional condition.

Assumption 17. $E \left[\frac{1}{2} (Z; \theta) \frac{1}{2} (Z; \theta)' \mid X \right]$ has up to J -th derivatives with respect to X and the J -th derivative is uniformly bounded in $X \in X$ and $\theta \in \Theta_n(o(n^{-1/4}))$.

Theorem 6.2. Under Assumptions 1, 2, 4 - 12, 14(ii), 16 and 17, we obtain: (i) $\hat{S}_0(X_i) = S_0(X_i) + o_p(n^{-1/4})$ uniformly over $X_i \in X_n$; (ii) If in addition, V_0 is positive-definite and Assumption 13 hold, then: $\sqrt{n}(\hat{\beta}_n - \mu_0) \Rightarrow N(0; V_0^{-1})$.

It is worth noting that Step 2 can be repeated using $\hat{\theta}_n$ to re-estimate the weighting matrix and Step 3 is then repeated to obtain a new estimator. The new estimator for μ_0 is asymptotically as efficient as $\hat{\beta}_n$. Thus, there is no efficiency gain by repeating Step 2 and Step 3. This remains true even if Step 2 and Step 3 are repeated many times. It is possible, however, the new estimator have better finite sample properties.

7. Applications

In this section we provide relative “low-level” sufficient conditions to obtain efficient estimators for the three examples in Section 2.

Example 2.1 Continued: We estimate $\theta_0 = (\mu_0; h_{01}; \dots; h_{0q})$ using the simpler sieve generalized least squares regression (2.2). Consider the distance:

$$\begin{aligned} K(\theta_0; \theta) &= \frac{1}{2} E \left(\frac{f(X_0; \mu) + \sum_{j=1}^q [h_j(X_j) - h_{0j}(X_j)] g^0(X) i^{1/2}}{f(X_0; \mu) + \sum_{j=1}^q [h_j(X_j) - h_{0j}(X_j)] g} \right)^2 \\ &= \frac{1}{2} \|\theta - \theta_0\|_2^2 = \frac{1}{2} \|\theta - \theta_0\|_2^2 \end{aligned}$$

Suppose V_0 given by (2.4) is finite positive definite. Then the three-step procedure in Section 6 leads to \sqrt{n} -efficient estimator of μ_0 with asymptotic variance V_0^{-1} .

Example 2.2 Continued: we apply the generalized minimum distance (2.1) to estimate $(\mu_0; h_{01}; \dots; h_{0q})$. Consider the distance:

$$\begin{aligned} K(\theta_0; \theta) &= \frac{1}{2} E \left(\frac{f(\mu; \mu_0) E[Y_2 | X] + \sum_{j=1}^q [h_j(X_j) - h_{0j}(X_j)] g^0(X) i^{1/2}}{f(\mu; \mu_0) E[Y_2 | X] + \sum_{j=1}^q [h_j(X_j) - h_{0j}(X_j)] g} \right)^2 \\ &= \frac{1}{2} \|\theta - \theta_0\|_2^2 = \frac{1}{2} \|\theta - \theta_0\|_2^2 \end{aligned}$$

Suppose V_0 given by (2.5) is finite positive definite. Then the three-step procedure in Section 6 leads to \sqrt{n} -efficient estimator of μ_0 with asymptotic variance V_0^{-1} .

Example 2.3 (Semiparametric sample selection): for $l = 1; 2$

does not depend on the endogenous variable Y , we may apply the simpler sieve generalized least squares regression (2.2) to estimate $(\mu_{01}; \mu_{021}; \mu_{022}; h_{01}; h_{02})$ instead of the method of generalized minimum distance (2.1).

Example 2.4 Continued: we apply the method of generalized minimum distance (2.1) to estimate $(\mu_0; h_0)$. Consider the distance:

$$K(\theta_0; \theta) = \frac{1}{2} \|\theta - \theta_0\|_2^2 = \frac{1}{2} \|\theta - \theta_0\|_2^2$$

$$= E[\text{ff}X_1^0(\mu; \mu_0) + E[h(Y_2) | h_0(Y_2) | X]g^0]S^{i-1}fX_1^0(\mu; \mu_0) + E[h(Y_2) | h_0(Y_2) | X]gg$$

For identification, we assume that there is no linear combination of X_1 enters additively in $E[h_0(Y_2) | X_1; X_2]$. We also assume that Assumptions 1-5 are satisfied.

Suppose V_0 given by (2.6) is finite positive definite. Then the three-step procedure in Section 6 leads to \sqrt{n} -efficient estimator of μ_0 with asymptotic variance V_0^{-1} .

8. Conclusion

In this paper, we present a consistent estimator with rate for the general semiparametric conditional moments restrictions model (1.3). We derive the \sqrt{n} asymptotic normality of the estimator for the parametric component and provide a consistent estimator of its asymptotic covariance matrix. We show that the optimally weighted sieve minimum distance estimator of the parametric component can attain the semiparametric efficiency bound of model (1.3). We finally provide a three-step procedure to obtain the efficient estimator.

Our results can be extended in several directions. First, model (1.3) assumes all moment restrictions hold conditional on the same set of regressors. This assumption may rule out many interesting applications. For instance, consider a partially linear regression model with endogenous regressors: $Y = X_1\mu_0 + h_{01}(X_2) + u_1$ and $X_2 = h_{02}(X_3) + u_2$, where u_1 is correlated with u_2 , $E[u_1u_2 | X_1; X_3] = E[u_1u_2]$ and $E[u_2 | X_3] = 0$. Let $\eta_1(Z; \mu_0; h_0) = Y | X_1\mu_0 | h_{01}(X_2) | h_{03}(X_2 | h_{02}(X_3))$ and $\eta_2(Z; \mu_0; h_0) = X_2 | h_{02}(X_3)$. Then we have $E[\eta_1(Z; \mu_0; h_0) | X_1; X_2; X_3] = 0$ and $E[\eta_2(Z; \mu_0; h_0) | X_3] = 0$.

Our next example is a modified version of the sample selection model considered by Das, Newey and Vella (1999):

$$Y_2 = Y_1 E[h_{02}(X_2) + X_3^0\mu_0 + U]; Y_1 = 1fg(X_1) \text{ , eg}$$

where g and h_{02} are unknown functions. Let $X = (X_1; X_2; X_3)$ and suppose $E[Y_1 | X] = h_{01}(X_1)$ and $E[U | Y_1 = 1; X] = h_{03}(h_{01}(X_1))$ where both h_{01} and h_{03} are unknown functions. It is easy to show that equation (1.3) is satisfied with

$$\eta(Z; \mu_0) = (Y_1 | h_{01}(X_1); Y_2 | h_{01}(X_1)[h_{02}(X_2) + X_3^0\mu_0 + h_{03}(h_{01}(X_1))])^0;$$

where $h_0(\mathbb{C}) = (h_{01}(\mathbb{C}); h_{02}(\mathbb{C}); h_{03}(\mathbb{C}))^0$. In this example, the argument of h_{01} is X_1 , the argument of h_{02} is X_2 , and the argument of h_{03} is the unknown function h_{01} . This example shows that different moment restrictions may hold conditional on

different set of regressors. A related problem is that some conditioning variables have index structure such as those in Ai (1997) and Andrews (1995). Our procedure and results need to be extended to models of this sort. Second, although our results are sufficient for statistical inference on the parametric component, they are not sufficient for statistical inference on the nonparametric component. In some applications, we may be interested in testing the functional forms of $h_0(\cdot)$ (e.g. completely unspecified $h_0(\cdot)$ against the index form or the index form against a parametric specification of $h_0(\cdot)$). We may also be interested in testing the joint restrictions of both the parametric component and the nonparametric component. We are currently undertaking a project to develop more efficient tests for all sorts of restrictions of model (1.3). Third, as we noted before, although our results on convergence rates and \sqrt{n} asymptotic normality can be easily extended to weakly dependent time series data, the problem of semiparametric efficiency bound with time series data is non-trivial. Since most dynamic models of asset pricing fit into the framework (1.3) with stationary Markov data, see e.g. Hansen and Richard (1987), we plan to address the efficiency issue in a separate paper. Lastly but not the least importantly, we need to address problems such as finite sample performance, higher-order analysis and choice of smoothing parameters in subsequent research.

References

- [1] Ai, C. (1997): "A Semiparametric Maximum Likelihood Estimator", *Econometrica*, 65, 933-964.
- [2] Andrews, D. (1994): "Asymptotics for Semi-parametric Econometric Models via Stochastic Equicontinuity", *Econometrica*, 62, 43-72.
- [3] Andrews, D. (1995): "Nonparametric kernel estimation for semiparametric models", *Econometric Theory*, 11, 560-596.
- [4] Blundell, R., M. Browning and I. Crawford (2000): "Nonparametric Engel Curves and Revealed Preference", UCL Working Paper.
- [5] Chamberlain, G. (1987): "Asymptotic Efficiency in Estimation with Conditional Moment Restrictions," *Journal of Econometrics*, 34, 305-334.
- [6] Chamberlain, G. (1992): "Efficiency Bounds for Semiparametric Regression," *Econometrica*, 60, 567-596.
- [7] Chen, X. and T. Conley (1999): "A New Semiparametric Spatial Model for Panel Time Series." Northwestern Working Paper.
- [8] Chen, X. and X. Shen (1998): "Sieve Extremum Estimates for Weakly Dependent Data," *Econometrica*, 66, 289-314.
- [9] Chen, X. and H. White (1999): "Improved Rates and Asymptotic Normality for Nonparametric Neural Network Estimators", *IEEE Tran. Information Theory*, 45, 682-691.
- [10] Chen, X., L.P. Hansen, and J. Scheinkman (1997): "Shape-preserving Estimation of Diffusions", manuscript, University of Chicago, Dept. of Economics.
- [11] Das, M., W.K. Newey and F. Vella (1999): "Non-parametric Estimation of the Sample Selection Model", manuscript, MIT Dept. of Economics.
- [12] DeVore, R.A. and G. Lorentz (1993): *Constructive Approximation*. New York: Springer-Verlag.

- [13] Fan, J., W. Härdle and E. Mammen (1998): "Direct Estimation of Low Dimensional Components in Additive Models," *Annals of Statistics* 26, 943-971.
- [14] Fenton, V. and A.R. Gallant (1996): "Convergence Rate of SNP Density Estimators," *Econometrica* 64, 719-727.
- [15] Grenander, U. (1981): *Abstract Inference*, New York: Wiley Series.
- [16] Hansen, L.P. (1982): "Large Sample Properties of Generalized Method of Moments Estimators," *Econometrica*, 50, 1029-1054.
- [17] Hansen, L.P. and S.F. Richard (1987): "The Role of Conditioning Information in Deducing Testable Restrictions Implied by Dynamic Asset Pricing Models", *Econometrica*, 55, 587-613.
- [18] Horowitz, J. (1996): "Semiparametric Estimation of a Regression Model with an Unknown Transformation of the Dependent Variable," *Econometrica* 64:1, 103-137.
- [19] Horowitz, J. (1998): *Semiparametric Methods in Econometrics*, Springer-Verlag.
- [20] Ichimura, H. and L.F. Lee (1991): "Semiparametric Least Squares Estimation of Multiple Models: Single Equation Estimation," in *Nonparametric and Semiparametric Models in Econometrics and Statistics*, ed. by W. Barnett, J. Powell, and G. Tauchen. Cambridge: Cambridge University Press.
- [21] Ichimura, H. (1993): "Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single Index Models," *Journal of Econometrics* 58, 71-120.
- [22] Linton, O. (2000): "Edgeworth Approximations for Semiparametric Instrumental Variable Estimators and Test Statistics", LSE Working Paper.
- [23] Newey, W.K. (1990): "Efficient Instrumental Variables Estimation of Non-linear Models," *Econometrica*, 58, 809-837.
- [24] Newey, W.K. (1993): "Efficient Estimation of Models with Conditional Moment Restrictions, in *Handbook of Statistics*, Vol. 11, G.S. Maddala, C.R. Rao, and H.D. Vinod, eds., Amsterdam: North-Holland.

- [25] Newey, W.K. (1994): "The Asymptotic Variance of Semiparametric Estimators", *Econometrica* 62, 1349-1382.
- [26] Newey, W.K. (1997): "Convergence Rates and Asymptotic Normality for Series Estimators", *Journal of Econometrics* 79, 147-168.
- [27] Newey, W.K. and J.L Powell (1989): "Nonparametric Instrumental Variables Estimation", MIT Working Paper.
- [28] Newey, W.K., J.L. Powell, and F. Vella (1999): "Nonparametric Estimation of Triangular Simultaneous Equations Models", *Econometrica* 67, 565-603.
- [29] Pakes, A. and S. Olley (1995): "A Limit Theorem for A Smooth Class of Semiparametric Estimators", *Journal of Econometrics* 65, 295-332.
- [30] Powell, J., J. Stock, and T. Stoker (1989): "Semiparametric Estimation of Index Coefficients", *Econometrica* 57, 1403-1430.
- [31] Powell, J. (1994): Estimation of Semiparametric Models, in *Handbook of Econometrics*, ed. by R. Engle and D. McFadden.
- [32] Robinson, P. (1987): "Asymptotically Efficient Estimation in the Presence of Heteroskedasticity of unknown Form," *Econometrica* 55, 875-892.
- [33] Robinson, P. (1988): "Root-N-Consistent Semiparametric Regression," *Econometrica* 56, 931-954.
- [34] Robinson, P. (1991): "Best Nonlinear Three-Stage Least Squares Estimation of Certain Econometric Models," *Econometrica* 59, 755-786.
- [35] Severini, T. and H.W. Wong (1992): "Profile Likelihood and Conditionally Parametric Models," *The Annals of Statistics*, 20(4), 1768-1802.
- [36] Shen, X. (1997): "On Methods of Sieves and Penalization," *The Annals of Statistics* 25(6), 2555-2591.
- [37] Silverman, B. (1986): *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall.
- [38] Van der Vaart, A. and J. Wellner (1996): *Weak Convergence and Empirical Processes: with Applications to Statistics*. New York: Springer

Mathematical Appendices

This appendix consists of several parts. In Appendix A, we present some uniform convergence results for the kernel estimators. In Appendix B, we prove the consistency of the proposed estimator $\hat{\theta}_n$ and compute the convergence rates. In Appendix C, we prove several Lemmas which are then used to establish the asymptotic distribution of the estimator $\hat{\beta}_n$. In Appendix D, we prove the consistency of the variance estimator $\hat{\Psi}$ and verify the efficiency property of the optimally weighted LS estimator.

Appendix A: Some Uniform Convergence Results

Let $G(Z; \theta)$ denote a generic measurable function of the data $Z \in \mathcal{Z}$ and the parameters $\theta \in \mathcal{A}$. Recall that $f_X(x)$ is the joint density function of the s -dimensional random variables X , $K(x)$ the kernel function, and a_n the bandwidth. Define

$$g(x; \theta) = f_X(x) E[G(Z; \theta) | X = x]$$

Denote the kernel estimator of $g(X_i; \theta)$ by

$$\hat{g}(X_i; \theta) = \frac{1}{(n_i - 1)a_n^s} \sum_{j \in i, j=1}^n G(Z_j; \theta) K\left(\frac{X_i - X_j}{a_n}\right)$$

Also, recall that X_n is the trimmed support of X ; A_n the sieve approximation of A , $\|\cdot\|_2$ the norm on A , and $N(A_n; \epsilon)$ the minimum number of mutually exclusive partitions of $A_n = \bigcup_{k=1}^{N(A_n; \epsilon)} A_{kn}$ such that for any $\theta_1, \theta_2 \in A_{kn}$, we have $\|\theta_1 - \theta_2\|_2 \leq \epsilon$.

Lemma A.1: Suppose that the data $\{Z_i; i = 1, 2, \dots, n\}$ satisfy Assumption 1, the density of X satisfies Assumption 4, and the kernel function satisfies Assumption 5. In addition, suppose the followings hold:

- A.1.1. $E[G(Z; \theta) | X]$ has up to J -th total derivatives with respect to X , and the J -th derivative is uniformly bounded over $X \in X$ and $\theta \in A_n$;
- A.1.2. there exists a measurable function $C_1(Z) : \mathcal{Z} \rightarrow [0, 1)$ such that $E[C_1(Z)^p] < 1$ for some $p > 2$, and $|G(Z; \theta)| \leq C_1(Z)$ for all $\theta \in A_n$;
- A.1.3. there exists a measurable function $C_{2n}(Z) : \mathcal{Z} \rightarrow [0, 1)$ such that $E[C_{2n}(Z) | X] = c_{2n} < 1$ and $|G(Z; \theta_1) - G(Z; \theta_2)| \leq C_{2n}(Z) \|\theta_1 - \theta_2\|_2$ for any $\theta_1, \theta_2 \in A_n$;
- A.1.4. $\{\epsilon_n\}$ is a sequence of small positive values satisfying

$$\frac{na_n^{s+2}}{\ln[N(A_n; \frac{\epsilon_n}{c_{2n}}) \leq a_n^{s(s+1)} \epsilon_n^s] \leq \max\{1, \frac{1}{\epsilon_n^{2-p}} a_n^{2(s+1)-p}\} \epsilon_n} \rightarrow 1 \text{ as } n \rightarrow \infty$$

Then uniformly over $(X_i; \otimes) \in X_n \in A_n$:

$$g(X_i; \otimes) = g(X_i; \otimes) + o_p(\pm_n) + O_p(a_n^J) :$$

Proof. Let $E[g(X_i; \otimes) | X_i]$ denote the conditional expectation given X_i : Write

$$g(X_i; \otimes) - E[g(X_i; \otimes) | X_i] = E[g(X_i; \otimes) | X_i] - E[g(X_i; \otimes) | X_i] + E[g(X_i; \otimes) | X_i] - g(X_i; \otimes) :$$

Under Assumptions 1, 5(ii), using the variable transformation $t = \frac{x_i - X_i}{a_n}$, we have:

$$\begin{aligned} E[g(X_i; \otimes) | X_i] &= \int_{\mathcal{Z}} \frac{1}{a^s} E[G(Z; \otimes(t)) | X = x] K\left(\frac{X_i - x}{a_n}\right) f_X(x) dx \\ &= E[G(Z; \otimes(t)) | X = X_i + a_n t] K(t) f_X(X_i + a_n t) dt : \end{aligned}$$

Under Assumptions 4, 5 and Condition A.1.1, Taylor expansion of $E[G(Z; \otimes(t)) | X = X_i + a_n t] f_X(X_i + a_n t)$ around X_i to the J -th order yields: $E[g(X_i; \otimes) | X_i] - g(X_i; \otimes) = O_p(a_n^J)$ uniformly over $(X_i; \otimes) \in X_n \in A_n$.

Next, consider the residual $g(X_i; \otimes) - E[g(X_i; \otimes) | X_i]$. Let c denote a generic constant which may have different values in different expressions. Let $f_{M_n g}$ be a sequence of positive values that increase with the sample size at certain rates to be specified later. Define $d_{j_n} = 1 + c_1(Z_j) / M_n g$. Then we can write

$$g(X_i; \otimes) - E[g(X_i; \otimes) | X_i] = (g_1(X_i; \otimes) - E[g_1(X_i; \otimes) | X_i]) + (g_2(X_i; \otimes) - E[g_2(X_i; \otimes) | X_i])$$

where

$$\begin{aligned} g_1(X_i; \otimes) &= \frac{1}{(n-1)a^s} \sum_{j \neq i, j=1}^n d_{j_n} G(Z_j; \otimes) K\left(\frac{X_i - X_j}{a_n}\right) ; \\ g_2(X_i; \otimes) &= \frac{1}{(n-1)a^s} \sum_{j \neq i, j=1}^n [1 - d_{j_n}] G(Z_j; \otimes) K\left(\frac{X_i - X_j}{a_n}\right) : \end{aligned}$$

For any small value $\epsilon > 0$; and a positive sequence $f_{\pm n g}$ satisfying Condition A.1.4, we have:

$$\begin{aligned} &P \sup_{(X_i; \otimes) \in X_n \in A_n} |g(X_i; \otimes) - E[g(X_i; \otimes) | X_i]| > \epsilon_{\pm n} \\ &P \sup_{(X_i; \otimes) \in X_n \in A_n} |g_1(X_i; \otimes) - E[g_1(X_i; \otimes) | X_i]| > \epsilon_{\pm n} \\ &+ P \sup_{(X_i; \otimes) \in X_n \in A_n} |g_2(X_i; \otimes) - E[g_2(X_i; \otimes) | X_i]| > \epsilon_{\pm n} \\ &= P_1 + P_2 : \end{aligned}$$

Applying the Chebyshev inequality and under Assumption 5(i) and Condition A.1.2, we obtain

$$P_2 \leq \frac{2E \sup_{(X_i; \otimes) \in A_n} |g_2(X_i; \otimes) - E[g_2(X_i; \otimes)]|^2}{4E [1 - d_{jn}] \sup_{(X_i; \otimes) \in A_n} |G(Z_j; \otimes) - E[G(Z_j; \otimes)]|^2} \leq \frac{cE [1 - d_{jn}] C_1(Z_j)}{E [1 - d_{jn}] E [C_1(Z_j)^2]} = \frac{c}{M_n^{p-2} a_n^s} = O(a_n);$$

where the last inequality follows from Condition A.1.2 and

$$E[1 - d_{jn}] = P(C_1(Z_j) > M_n) \leq \frac{E[C_1(Z_j)^p]}{M_n^p};$$

and the last equality is obtained by letting $M_n = (\pm_n a_n^{s+1})^{2=p}$. Then $P_2 \rightarrow 0$ as $n \rightarrow \infty$.

Let $W_n = X_n \in A_n$. For any pair $(X_i^1; \otimes^1) \in W_n$ and $(X_i^2; \otimes^2) \in W_n$, Assumption 5(i), Conditions A.1.2 and A.1.3 imply:

$$\begin{aligned} & \frac{1}{a_n^s} d_{jn} G(Z_j; \otimes^1) K \frac{\|X_i^1 - X_j\|}{a_n} + \frac{1}{a_n^s} d_{jn} G(Z_j; \otimes^2) K \frac{\|X_i^2 - X_j\|}{a_n} \\ & \leq \frac{1}{a_n^s} d_{jn} \left[|G(Z_j; \otimes^1) - G(Z_j; \otimes^2)| K \frac{\|X_i^1 - X_j\|}{a_n} \right. \\ & \quad \left. + |G(Z_j; \otimes^2)| K \frac{\|X_i^2 - X_j\|}{a_n} \right] \\ & \leq \frac{c}{a_n^s} C_{2n}(Z_j) K \frac{\|X_i^1 - X_j\|}{a_n} + \frac{c}{a_n^{s+1}} C_1(Z_j) \|X_i^1 - X_j^2\|; \end{aligned}$$

where $\|X_i^1 - X_j^2\|$ denotes the Euclidean norm of $X_i^1 - X_j^2$. Since $E[C_1(Z)]$ is finite, by the Law of Large Number for i.i.d. processes, we have:

$$\frac{1}{(n-1) \sum_{j \neq i} 1} C_1(Z_j) = O_p(1);$$

Since $E[C_{2n}(Z)|X] = c_{2n} < 1$, we have uniformly over $X_i^1 \in X$:

$$\frac{1}{c_{2n}} \frac{1}{(n-1)a_n^s} \sum_{j \neq i, j=1}^n C_2(Z_j) \in jK \frac{\prod_{i=1}^n X_i^1}{a_n} \quad j = O_p(1):$$

Thus, for some arbitrarily small value $\epsilon > 0$; there exists a sufficiently large n such that, with $\prod_{i=1}^n X_i^1 \in a_n^{s+1} \pm_n = (12c)$ and $k^{\otimes 2} \in \otimes^1 k_2 \pm_n = (12c_{2n})$, we have with probability at least $1 - \epsilon$:

$$\prod_{i=1}^n g(X_i^2; \otimes^2) \in \prod_{i=1}^n g(X_i^1; \otimes^1) \pm_n = \epsilon:$$

Partition W_n into b_n mutually exclusive subsets $W_{nm}; m = 1; 2; \dots; b_n$; where $(X_i^1; \otimes^1) \in W_{nm}$ and $(X_i^2; \otimes^2) \in W_{nm}$ imply $\prod_{i=1}^n X_i^1 \in a_n^{s+1} \pm_n = (12c)$ and $k^{\otimes 2} \in \otimes^1 k \pm_n = (12c_{2n})$: Since X_n is a compact subset in R^s , we have

$$b_n = O \left((a_n^{s+1} \pm_n)^s N(A_n; \frac{\pm_n}{c_{2n}}) \right):$$

Let $(X_i^m; \otimes^m)$ denote an arbitrarily fixed point in W_{nm} . We obtain:

$$\begin{aligned} P_1 &< \epsilon + \sum_{m=1}^{b_n} P \left(\prod_{i=1}^n g_1(X_i^m; \otimes^m) \in E[g_1(X_i^m; \otimes^m)|X_i^m] \right) > \pm_n = \epsilon \\ &= \epsilon + \sum_{m=1}^{b_n} P \left(\prod_{i=1}^n g_1^a(X_i^m; \otimes^m) \in E[g_1^a(X_i^m; \otimes^m)|X_i^m] \right) > a_n^s \pm_n = \epsilon \end{aligned}$$

where $g_1^a(X_i^m; \otimes^m) = a_n^s g_1(X_i^m; \otimes^m) = \frac{1}{(n-1)} \prod_{j \neq i, j=1}^n d_{jn} G(Z_j; \otimes^m) K \frac{\prod_{i=1}^n X_i^m}{a_n}$; which is bounded by $O(M_n)$. Applying the Bernstein inequality for i.i.d. processes, see e.g. Ichimura and Lee (1991, p.26), we obtain:

$$\begin{aligned} P \left(\prod_{i=1}^n g_1^a(X_i^m; \otimes^m) \in E[g_1^a(X_i^m; \otimes^m)|X_i^m] \right) > a_n^s \pm_n = \epsilon \\ 2 \exp \left\{ - \frac{a_n^{2s} \pm_n^2}{c_{3m}^2 + 8M_n a_n^s \pm_n} \right\} \end{aligned}$$

where c_{3m}^2 is the conditional variance of $d_{jn} G(Z_j; \otimes^m) K \frac{\prod_{i=1}^n X_i^m}{a_n}$ given X_i^m . Under Assumptions 1, 4, 5 and Condition A.1.2, it is easy to see that c_{3m}^2 is $O(a_n^s \int K^2(t) f_X(X_i^m + a_n t) dt) = O(a_n^s)$. Hence,

$$P_1 < \epsilon + 2b_n \exp \left\{ - \frac{a_n^{2s} \pm_n^2}{c + 8M_n \pm_n} \right\}$$

which is arbitrarily small if $\frac{na_n^{s+2}}{\max\{1; M_n\}ng} \leq \ln(b_n) + 1$: By the construction of M_n and b_n , we have $b_n \leq \ln(b_n) + 1$ and

$$= O\left(\frac{na_n^{s+2}}{\ln(b_n) \max\{1; M_n\}ng}\right) = O\left(\frac{na_n^{s+2}}{\ln\left[N\left(\mathbf{A}_n; \frac{\pm_n}{c_{2n}}\right) \leq a_n^{i^{s(s+1)}} \pm_n^i\right] \max\{1; \pm_n^{1-2p} a_n^{2(s+1)-p}\}g}\right)$$

where the limit to infinity is by Condition A.1.4. This completes the proof. ■

The uniform convergence rate of the kernel estimator $\hat{g}(X_i; \mathbb{R})$ has two components. The first term, $O(a_n^J)$, is the convergence rate of the bias term. This rate depends on the order of the kernel and on the smoothness condition of $g(X; \mathbb{R})$. Clearly, higher the order of the kernel function and the degree of smoothness, faster the bias term shrinks to zero. This component is obtained under Assumptions 1, 4, 5 and Condition A.1.1, while Conditions A.1.2, A.1.3 and A.1.4 are not needed. The second component, $O_p(\pm_n)$, is the uniform convergence rate of the residual term. This rate is limited by: (i) p , the number of the moments of $G(Z; \mathbb{R})$, (ii) s , the dimension of X , and (iii) $N(\mathbf{A}_n; \cdot)$; the size of the sieve space. Higher the number of moments, lower the dimension of X ; and the smaller the size of the sieve space, faster the residual term converges to zero. This component is obtained under Assumptions 1, 4, 5 and Conditions A.1.2, A.1.3 and A.1.4 without using Condition A.1.1. It must be noted that Lemma A.1 holds only on the sieve parameter space \mathbf{A}_n , not on the original parameter space \mathbf{A} :

Lemma A.1 still holds when we weaken Conditions A.1.1-A.1.2. For example, if Condition A.1.1 is weakened to assume that the J_j th total derivative of $g(X; \mathbb{R})$ with respect to X is bounded by B_n over the sieve space where B_n could grow slowly with the sample size, then Lemma A.1 holds with the convergence rate of the bias term replaced by $O(B_n a_n^J)$. Likewise, we can relax Condition A.1.2 by replacing the envelope function $C_1(Z)$ by $C_{1n}(Z)$ where $c_{1n}^p = E[C_{1n}(Z)^p]$ may increase with the sample size. Then Lemma A.1 still holds with the convergence rate of the residual term $O_p(\pm_n)$ if the term

$$\ln\left[N\left(\mathbf{A}_n; \frac{\pm_n}{c_{2n}}\right) \leq a_n^{i^{s(s+1)}} \pm_n^i\right] \max\{1; \pm_n^{1-2p} a_n^{2(s+1)-p}\}g$$

in Condition A.1.4 is replaced by

$$\ln\left[N\left(\mathbf{A}_n; \frac{\pm_n}{c_{2n}}\right) \leq \frac{C_{1n}^s}{a_n^{s(s+1)} \pm_n^s}\right] \max\{1; c_{1n}^{2-p} \pm_n^{1-2p} a_n^{2(s+1)-p}\}g$$

Corollary A.1: Under Assumptions 1, 4, 5 and 8, we have: $f_{X_i}^{\pm} = f_{X_i} + o_p(n^{i-4})$ uniformly over $X_i \in X_n$:

Proof. We apply Lemma A.1 to $G(Z; \mathbb{R}) = 1$. Conditions A.1.1 - A.1.3 are trivially satisfied with $p = +1$. Because $G(Z; \mathbb{R})$ does not depend on the parameters \mathbb{R} , we have $N(\mathbf{A}_n; \pm_n) = 1$. Because $p = +1$, we have $\max\{1; \pm_n^{1-2p} a_n^{2(s+1)-p}\} = \max\{1; \pm_n\} = 1$. Hence, with $\pm_n = n^{i-4}$, Condition A.1.4 is satisfied if

$$\frac{p_n a_n^s}{\ln[a_n^{s(s+1)} n^{s-4}]} \rightarrow +1 \text{ as } n \rightarrow +1$$

which is implied by Assumption 8. Assumption 8 also implies $n^{i-4} a_n^j \rightarrow 0$. ■

Corollary A.2: Under Assumptions 1, 4, 5, 6, 7 and 8, we obtain uniformly over $(X; \mathbb{R}) \in X_n \in \mathbf{A}_n$: $m(X; \mathbb{R}) = m(X; \mathbb{R}) + o_p(n^{i-4})$:

Proof. We apply Lemma A.1 to $G(Z; \mathbb{R}) = \frac{1}{2}(Z; \mathbb{R})$. Condition A.1.1 is implied by Assumption 6, Conditions A.1.2 and A.1.3 are implied by Assumption 7. Assumption 8 implies Condition A.1.4 with $\pm_n = n^{i-4}$ and that $n^{i-4} a_n^j \rightarrow 0$. The result now follows from Lemma A.1 and Corollary A.1. ■

Define

$$m_{\mathbb{R}}[X_i; \mathbb{C}^{\mathbb{R}}] = \frac{[(n_i - 1) a_n^s]^{i-1} \prod_{j \in i; j=1}^n \frac{1}{2} [Z_j; \mathbb{C}^{\mathbb{R}}] K \frac{X_{ii} X_j}{a_n}}{f_{X_i}^{\pm}}$$

Corollary A.3: Under Assumptions 1, 4, 5, 8 and 11, we obtain uniformly over $X \in X_n$ and $\mathbb{A} \in U_n$: $m_{\mathbb{R}_o}[X; \mathbb{A}] = m_{\mathbb{R}_o}[X; \mathbb{A}] + o_p(n^{i-4})$:

Proof. We apply Lemma A.1 to $G(Z; \mathbb{A}) = \frac{1}{2}_{\mathbb{R}_o}[Z; \mathbb{A}]$ and $\pm_n = n^{i-4}$. Assumption 11 implies Conditions A.1.1 - A.1.3 hold in U_n . Restricting the parameter space to U_n and Noting that $N(U_n; \pm_n) < N(\mathbf{A}_n; \pm_n)$, Condition A.1.4 is implied by Assumption 8. The result follows from Lemma A.1 and Corollary A.1. ■

For some sequence of positive values $\pm_n \rightarrow 0$ as $n \rightarrow 1$, Corollary A.3 implies that uniformly over $X \in X_n$ and $k_{\mathbb{R}_o}^{\pm} \in \mathbb{R}_o^{\pm_n}$,

$$m_{\mathbb{R}_o}[X; \mathbb{R}_o^{\pm} | \mathbb{R}_o] = m_{\mathbb{R}_o}[X; \mathbb{R}_o^{\pm} | \mathbb{R}_o] + o_p(n^{i-4} \pm_n).$$

Recall that $r_{\frac{1}{2}}[Z; \mathbb{R}_o^{\pm} | \mathbb{R}_o]$ is the remainder term in the linear approximation of $\frac{1}{2}(Z; \mathbb{R}_o)$. Define

$$r_m[X; \mathbb{R}_o^{\pm} | \mathbb{R}_o] = E[r_{\frac{1}{2}}[Z; \mathbb{R}_o^{\pm} | \mathbb{R}_o] | X]$$

and the kernel estimator as

$$\hat{b}_m[X_i; \mathbb{R}_i, \mathbb{R}_o] = \frac{[(n_i - 1)a_n^s]^{-1} \sum_{j \in \mathbb{R}_i} r_{\frac{1}{2}}[Z_j; \mathbb{R}_i, \mathbb{R}_o] K\left(\frac{X_i - X_j}{a_n}\right)}{f_{X_i}}$$

Corollary A.4: Under Assumptions 1, 4, 5, 6, 8, 11(i) and 12, we obtain: (i) uniformly over $X \in \mathcal{X}_n$ and $\mathbb{R} \in \mathcal{R}_o$:

$$\frac{\hat{b}_m[X; \mathbb{R}_i, \mathbb{R}_o]}{k_{\mathbb{R}_i, \mathbb{R}_o} k_2} = \frac{r_m[X; \mathbb{R}_i, \mathbb{R}_o]}{k_{\mathbb{R}_i, \mathbb{R}_o} k_2} + o_p(n_i^{-1/4})$$

(ii) uniformly over $X \in \mathcal{X}_n$, $\mathbb{R} \in \mathcal{R}_n$ with $k_{\mathbb{R}_i, \mathbb{R}_o} k_2 \asymp \pm_n$: $\hat{b}_m[X; \mathbb{R}_i, \mathbb{R}_o] = O_p(\pm_n^2) + o_p(n_i^{-1/4} \pm_n)$:

Proof. (i) We apply Lemma A.1 with $G(Z_j; \mathbb{R}) = \frac{r_{\frac{1}{2}}[Z_j; \mathbb{R}_i, \mathbb{R}_o]}{k_{\mathbb{R}_i, \mathbb{R}_o} k_2}$ and $\pm_n = n_i^{-1/4}$, Condition A.1.1 is implied by definition of $r_{\frac{1}{2}}[Z; \mathbb{R}_i, \mathbb{R}_o]$ and Assumptions 6 and 11(i); Conditions A.1.2 and A.1.3 are implied by Assumption 12, Condition A.1.4 is implied by Assumption 8. The result follows from Lemma A.1 and Corollary A.1.

(ii) Immediately follows from Result (i) and Assumption 12(i). ■

Corollary A.4(i) implies that uniformly over $X \in \mathcal{X}_n$ and $k_{\mathbb{R}_i, \mathbb{R}_o} k_2 \asymp \pm_n$,

$$\hat{b}_m[X; \mathbb{R}_i, \mathbb{R}_o] = r_m[X; \mathbb{R}_i, \mathbb{R}_o] + o_p(n_i^{-1/4} \pm_n)$$

Appendix B: Convergence Rates for Sieve Estimators

We now prove the consistency and compute the convergence rates of the proposed estimator. Define

$$\hat{L}_n(\mathbb{R}) = \frac{1}{n} \sum_{i=1}^n 1_{\{X_i \in \mathcal{X}_n\}} g_{\mathbb{R}}(X_i; \mathbb{R})^0 [S(X_i)]^{-1} m(X_i; \mathbb{R})^2$$

$$L_n(\mathbb{R}) = \frac{1}{n} \sum_{i=1}^n 1_{\{X_i \in \mathcal{X}_n\}} g_m(X_i; \mathbb{R})^0 [S(X_i)]^{-1} m(X_i; \mathbb{R})^2$$

$$\hat{L}_{n\mathbb{R}_o}[\mathbb{R}_i, \mathbb{R}_o] = \frac{1}{n} \sum_{i=1}^n 1_{\{X_i \in \mathcal{X}_n\}} g_{\mathbb{R}_o}[X_i; \mathbb{R}_i, \mathbb{R}_o]^0 [S(X_i)]^{-1} m(X_i; \mathbb{R}_o);$$

$$L_{n\mathbb{R}_o}[\mathbb{R}_i, \mathbb{R}_o] = \frac{1}{n} \sum_{i=1}^n 1_{\{X_i \in \mathcal{X}_n\}} g_{\mathbb{R}_o}[X_i; \mathbb{R}_i, \mathbb{R}_o]^0 [S(X_i)]^{-1} m(X_i; \mathbb{R}_o)$$

$\hat{L}_{n\mathbb{R}_o}[\mathbb{R}_i, \mathbb{R}_o] \rightarrow 0$ since $m(X_i; \mathbb{R}_o) \rightarrow 0$ by definition.

Also define:

$$R[\theta_i, \theta_o] = L_n(\theta_i) - L_n(\theta_o) - L_{n\theta_o}[\theta_i, \theta_o] = L_n(\theta);$$

$$\dot{R}[\theta_i, \theta_o] = \dot{L}_n(\theta_i) - \dot{L}_n(\theta_o) - \dot{L}_{n\theta_o}[\theta_i, \theta_o];$$

By definition of $r_m[X; \theta_i, \theta_o]$, and $m(X; \theta_o) = 0$; we have:

$$m(X; \theta) = m_{\theta_o}[X; \theta_i, \theta_o] + r_m[X; \theta_i, \theta_o];$$

Write

$$\begin{aligned} & \sum_{i=1}^n 2R[\theta_i, \theta_o] \\ = & \frac{1}{n} \sum_{i=1}^n 1fX_i - 2 X_n g m_{\theta_o}[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 m_{\theta_o}[X_i; \theta_i, \theta_o] + \\ & \frac{2}{n} \sum_{i=1}^n 1fX_i - 2 X_n g m_{\theta_o}[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 r_m[X_i; \theta_i, \theta_o] + \\ & \frac{1}{n} \sum_{i=1}^n 1fX_i - 2 X_n g r_m[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 r_m[X_i; \theta_i, \theta_o] \\ \sim & R1 + 2 \times R3 + R4; \end{aligned}$$

$$\begin{aligned} & \sum_{i=1}^n 2\dot{R}[\theta_i, \theta_o] \\ = & \frac{1}{n} \sum_{i=1}^n 1fX_i - 2 X_n g \dot{m}_{\theta_o}[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 \dot{m}_{\theta_o}[X_i; \theta_i, \theta_o] + \\ & \frac{2}{n} \sum_{i=1}^n 1fX_i - 2 X_n g \dot{m}(X_i; \theta_o) [S(X_i)]^i - 1 \dot{m}_m[X_i; \theta_i, \theta_o] + \\ & \frac{2}{n} \sum_{i=1}^n 1fX_i - 2 X_n g \dot{m}_{\theta_o}[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 \dot{m}_m[X_i; \theta_i, \theta_o] + \\ & \frac{1}{n} \sum_{i=1}^n 1fX_i - 2 X_n g \dot{m}_m[X_i; \theta_i, \theta_o] [S(X_i)]^i - 1 \dot{m}_m[X_i; \theta_i, \theta_o] \\ \sim & \dot{R}1 + 2 \times \dot{R}2 + 2 \times \dot{R}3 + \dot{R}4; \end{aligned}$$

The following corollaries are needed for deriving the asymptotic results of the proposed estimator.

Corollary A.5: (i) Under Assumptions 1, 3-8, we have uniformly over $\mathbb{R}^2 \mathbf{A}_n$:

$$\mathbf{L}_n(\mathbb{R}) - \mathbf{L}_n(\mathbb{R}) = o_p(n^{1-4}).$$

(ii) Under Assumptions 1-8, 11, we have uniformly over $\mathbb{R}^2 \mathbf{A}_n, k^{\mathbb{R}} - \mathbb{R}_0 k_2 \pm n$:

$$\mathbf{L}_{n^{\mathbb{R}_0}}[\mathbb{R} - \mathbb{R}_0] = o_p(n^{1-4} \pm n):$$

Proof. (i) A direct consequence of Assumptions 1, 3, Corollary A.2 and the following decomposition:

$$\begin{aligned} & \mathbf{m}(X_i; \mathbb{R})^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}) - \mathbf{m}(X_i; \mathbb{R})^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}) \\ = & \mathbf{f} \mathbf{m}(X_i; \mathbb{R}) - \mathbf{m}(X_i; \mathbb{R}) \mathbf{g}^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}_0(X_i; \mathbb{R}) + \\ & + \mathbf{m}(X_i; \mathbb{R})^0 [\mathbf{S}(X_i)]^{i-1} - [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}) + \\ & + \mathbf{m}(X_i; \mathbb{R})^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{f} \mathbf{m}(X_i; \mathbb{R}) - \mathbf{m}(X_i; \mathbb{R}) \mathbf{g} \end{aligned}$$

(ii) By definition, we have:

$$\begin{aligned} & \mathbf{L}_{n^{\mathbb{R}_0}}[\mathbb{R} - \mathbb{R}_0] \\ = & \sum_{i=1}^n \frac{1}{n} \mathbf{X}^{\mathbb{R}_0} \mathbf{1}(X_i \in \mathbf{X}_n) \mathbf{m}_{\mathbb{R}_0}[X_i; \frac{\mathbb{R} - \mathbb{R}_0}{k^{\mathbb{R}} - \mathbb{R}_0 k_2}]^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0): \end{aligned}$$

Now by the decomposition for any $\lambda \in \mathbf{U}_n$:

$$\begin{aligned} & \mathbf{m}_{\mathbb{R}_0}[X_i; \lambda]^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0) \\ = & \mathbf{m}_{\mathbb{R}_0}[X_i; \lambda]^0 [\mathbf{S}(X_i)]^{i-1} - [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0) \\ & + (\mathbf{m}_{\mathbb{R}_0}[X_i; \lambda] - \mathbf{m}_{\mathbb{R}_0}[X_i; \lambda])^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0) \\ & + \mathbf{m}_{\mathbb{R}_0}[X_i; \lambda]^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0) \end{aligned}$$

and Assumptions 1, 2, 3, Corollaries A.2 and A.3, we have uniformly over $\lambda \in \mathbf{U}_n$:

$$\frac{1}{n} \sum_{i=1}^n \mathbf{X}^{\mathbb{R}_0} \mathbf{1}(X_i \in \mathbf{X}_n) \mathbf{m}_{\mathbb{R}_0}[X_i; \lambda]^0 [\mathbf{S}(X_i)]^{i-1} \mathbf{m}(X_i; \mathbb{R}_0) = o_p(n^{1-4}):$$

Hence uniformly over $\mathbb{R}^2 \mathbf{A}_n; k^{\mathbb{R}} - \mathbb{R}_0 k_2 \pm n; \mathbf{L}_{n^{\mathbb{R}_0}}[\mathbb{R} - \mathbb{R}_0] = \pm n \in o_p(n^{1-4}): \blacksquare$

Corollary A.6: Under Assumptions 1-8, 11 and 12, we have uniformly over $k^{\otimes i} \otimes_0 k_2 \pm_n$ for any $\pm_n = o(n^{1-8})$:

- (i) $\hat{R}[\otimes_i \otimes_0] = R[\otimes_i \otimes_0] + o_p(n^{1-2})$;
- (ii) $\hat{L}_n(\otimes_i) - \hat{L}_n(\otimes_0) - fL_n(\otimes_i) - L_n(\otimes_0)g = o_p(n^{1-4} \pm_n)$;

Proof. (i) Recall $\hat{R}[\otimes_i \otimes_0] = \hat{R}1 + 2\hat{R}2 + 2\hat{R}3 + \hat{R}4$, and $R[\otimes_i \otimes_0] = R1 + 2R3 + R4$. Write

$$\begin{aligned} & \hat{R}1 - R1 \\ &= \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) (\mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0] - \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0])^0 [S(X_i)]^{i-1} \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0] \\ & \quad + \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0]^3 [S(X_i)]^{i-1} - [S(X_i)]^{i-1} \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0] \\ & \quad + \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0]^0 [S(X_i)]^{i-1} (\mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0] - \mathbb{m}_{\otimes_0}[X_i; \otimes_i \otimes_0]) : \end{aligned}$$

Corollary A.3 and Assumption 3 imply $\hat{R}1 - R1 = o_p(\pm_n^2 n^{1-4}) = o_p(n^{1-2})$ uniformly over $k^{\otimes i} \otimes_0 k_2 \pm_n$ for any $\pm_n = o(n^{1-8})$.

Similarly it is easy to check that Corollaries A.2 and A.4 imply $\hat{R}2 = o_p(n^{1-2})$ uniformly over $k^{\otimes i} \otimes_0 k_2 \pm_n$; Corollaries A.3 and A.4(ii) imply $\hat{R}3 - R3 = o_p(n^{1-2})$ uniformly over $k^{\otimes i} \otimes_0 k_2 \pm_n$; Corollary A.4(ii) implies $\hat{R}4 - R4 = o_p(n^{1-2})$ uniformly over $k^{\otimes i} \otimes_0 k_2 \pm_n$.

(ii) A direct consequence of Corollaries A.5(ii) and A.6(i). ■

We now are ready to prove Theorem 3.1.

Proof. (Theorem 3.1): Let $\eta = 1-4$; $\hat{\tau}_{0n} = o(n^\eta)$ and $\pm_{0n} = 2\frac{\hat{\tau}_{0n}}{0n} = o(n^{\eta-2})$: For any fixed $k > 1$,

$$\begin{aligned} & P_{\otimes} (k^{\otimes i} \otimes_0 k_2 \pm_{0n}) & 1 \\ & P_{\otimes} \sup_{f \in k^{\otimes i} \otimes_0 k_2 \pm_{0n}; \otimes_0 2A_n g} \hat{L}_n(\otimes_i) - \hat{L}_n(\otimes_0) & \hat{A} \\ & P \sup_{f \in \otimes_0 2A_n g} \hat{L}_n(\otimes_i) - L_n(\otimes_i) > \hat{\tau}_{0n} & ! \end{aligned}$$

$$\begin{aligned}
& P_1 + P_2: \\
& P_1 = \sup_{f \in \mathbb{R}^{2A_n g}} \left| \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right] \right| \\
& P_2 = \sup_{f \in \mathbb{R}^{2A_n g}} \left| \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right] \right|
\end{aligned}$$

Corollary A.5(i) implies $P_1 \rightarrow 0$ as $n \rightarrow \infty$. To show $P_2 \rightarrow 0$, we notice that all conditions A.1-A.4 in Chen and Shen (1998) are trivially satisfied given our Assumptions 1, 2, 3(ii), 7(ii), 9, 10 and our definition of k_{k_2} and τ_n . Hence $P_2 \rightarrow 0$ by the theorem 1 in Chen and Shen (1998). This proves: $k_{k_2} \in \mathbb{R}^{2A_n g} = O_p(\tau_n) = o_p(n^{1/2})$:

Next, we refine the convergence rate by exploiting the local curvature of $\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right]$ around μ_0 . Let $\tau_n = n^{1/4} \tau_{0n} = o(n^{1/2})$ and $\tau_{1n} = 2^{1/k_2} \tau_n = o(n^{1/4})$. For any fixed $k > 1$, we have:

$$\begin{aligned}
& P_3 = \sup_{f \in \mathbb{R}^{2A_n g}} \left| \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right] \right| \\
& P_4 = \sup_{f \in \mathbb{R}^{2A_n g}} \left| \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right] \right| \\
& P_3 + P_4:
\end{aligned}$$

Corollary A.6(ii) implies $\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(\frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(\frac{y - \mu}{\sigma} \right)^2 \right) dy \right) \right] = o_p(\tau_n)$, hence $P_3 \rightarrow 0$ as $n \rightarrow \infty$. Now under Assumptions 1, 2, 3(ii), 7(ii), 9

and 10, the Theorem 1 of Chen and Shen again implies $P_4 \rightarrow 0$ as $n \rightarrow \infty$. This proves $k_{n, \alpha}^2 = O_p(\pm_{1n}) = o_p(n^{-(1+1=2+1=4)})$.

Repeat the above proof j times, we obtain $\hat{\tau}_{jn} = o(n^{-(1+1=2+\dots+1=2^j)})$ and $k_{n, \alpha}^2 = O_p(\pm_{jn})$ with $\pm_{jn} = o(n^{-(1=2+1=2^2+\dots+1=2^{j+1})})$. Repeat the above proof infinite many times, we obtain $k_{n, \alpha}^2 = O_p(\pm_{1n})$ with $\pm_{1n} = o(n^{-(1=2+1=2^2+\dots)}) = o(n^{-1=4})$. This completes the proof. ■

Appendix C: Asymptotic Normality

In this section, we derive the asymptotic distribution of the proposed estimator. We begin by first introducing notation, then proving several Lemmas, and finally proving Theorem 4.1.

Define a pseudo-norm on A_n as:

$$\begin{aligned} k_{n, \alpha}^2 &= E \int_0^1 (X_i - X_n) m_{\alpha}^0 [X_i; \alpha] [S(X_i)]^{i-1} m_{\alpha}^0 [X_i; \alpha] \\ &= (\mu_i - \mu_0) E \int_0^1 (X_i - X_n) D_w(X) [S(X)]^{i-1} D_w(X) (\mu_i - \mu_0) \end{aligned}$$

Note that the difference between $k_{n, \alpha}^2$ and the Fisher-like norm $k_{n, \alpha}^2$ defined in section 4 is the trimming function. It is easy to show $k_{n, \alpha}^2 \rightarrow 0$ implies $k_{n, \alpha}^2 \rightarrow 0$. Let $\langle \cdot, \cdot \rangle_n$ denote the inner product induced by $k_{n, \alpha}^2$. For the bounded linear functional $f(\alpha) = \int_0^1 (\mu_i - \mu_0)$, there is again a Riesz representer under the inner product $\langle \cdot, \cdot \rangle_n$. Let $w_n^{\alpha}(\cdot) = (w_n^{\alpha 1}(\cdot); \dots; w_n^{\alpha b}(\cdot))$ be the solution to

$$\inf_{w(\cdot) \in \mathcal{W}; i=1, \dots, b} E [1(X \in X_n) D_w(X) [S(X)]^{i-1} D_w(X)]:$$

Since $1(X \in X_n) \rightarrow 1$, it is easy to show that $D_{w_n^{\alpha}}(X) \rightarrow D_{w^{\alpha}}(X)$ for each $X \in X_n$. Let $v_n^{\alpha} = (v_{n\mu}^{\alpha}; v_{nh}^{\alpha})$, where

$$\begin{aligned} v_{n\mu}^{\alpha} &= E \int_0^1 D_{w_n^{\alpha}}(X) [S(X)]^{i-1} D_{w_n^{\alpha}}(X) \mathbf{i}^{\top} \mathbf{i}^{-1} \mathbf{s}; \\ v_{nh}^{\alpha} &= \int_0^1 w_n^{\alpha}(\cdot) E \int_0^1 D_{w_n^{\alpha}}(X) [S(X)]^{i-1} D_{w_n^{\alpha}}(X) \mathbf{i}^{\top} \mathbf{i}^{-1} \mathbf{s}; \end{aligned}$$

Then, by the Riesz representation theorem, we have $f(\alpha) - f(\alpha_0) = \mathbf{h}_{n, \alpha}^{\top} v_n^{\alpha}$. In the following we shall link $\mathbf{h}_{n, \alpha}^{\top} v_n^{\alpha}$ to the derivative of the sample criterion function, and then replace v_n^{α} by v^{α} and removes the trimming function.

Remark A.1: Under Assumptions 1, 11(ii) and 14(i), we have: (i) uniformly over $X_i \in X_n$: $m_{\alpha}^0 [X_i; v_n^{\alpha}] = m_{\alpha}^0 [X_i; v^{\alpha}] + o_p(1)$; (ii) $k_{n, \alpha}^2 \rightarrow k^2$.

Corollary A.7: Under Assumptions 1-8, 11 and 14(i), we obtain:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} m(X_i; \theta_0) \\ = \frac{1}{n} \sum_{i=1}^n m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{1}{f_{X_i}} + o_p\left(\frac{1}{n}\right); \end{aligned}$$

Proof. Write

$$\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} m(X_i; \theta_0) = A + B + C$$

where

$$A = \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{1}{f_{X_i}} m(X_i; \theta_0);$$

$$B = \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) (m_{\theta_0}[X_i; v_n^a] - m_{\theta_0}[X_i; v_n^a]) [S(X_i)]^{i-1} m(X_i; \theta_0);$$

$$C = \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} m(X_i; \theta_0);$$

Applying Corollaries A.2 and A.3, and Assumptions 1 and 3, and Remark A.1, we have: $A = o_p(n^{-1/2})$ and $B = o_p(n^{-1/2})$: Denote $g(x; \theta_0) = m(x; \theta_0) f_{X_i}$, we have $g(x; \theta_0) = m(x; \theta_0) f_{X_i} = 0$ by Assumption 2. Then

$$\begin{aligned} C &= \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{g(X_i; \theta_0)}{f_{X_i}} \\ &= \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{1}{f_{X_i}} \frac{1}{f_{X_i}} g(X_i; \theta_0) \\ &+ \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{g(X_i; \theta_0)}{f_{X_i}}; \end{aligned}$$

Corollaries A.1, A.3, and Assumptions 1, 2 and 4 imply that

$$g(X_i; \theta_0) = o_p(n^{-1/4}), \quad \frac{1}{f_{X_i}} \frac{1}{f_{X_i}} = o_p(n^{-1/4}) \quad \text{uniformly over } X_i \in X_n;$$

Hence

$$C = \frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) m_{\theta_0}[X_i; v_n^a] [S(X_i)]^{i-1} \frac{g(X_i; \theta_0)}{f_{X_i}} + o_p\left(\frac{1}{n}\right);$$

Note that

$$g(X_i; \otimes_0) = [(n_i - 1)a_n^s]^{i-1} \sum_{j \in i; j=1}^n \frac{1}{2}(Z_j; \otimes_0(\mathfrak{t})) K \frac{\mu_{X_i, i, X_j}^{\mathfrak{t}}}{a_n}.$$

By exchanging summation (since $K(\mathfrak{t})$ is symmetric), we obtain

$$C = \frac{1}{n} \sum_{j=1}^n \frac{1}{(n_i - 1)a_n^s} \sum_{i \in j; i=1}^n \frac{1(X_i \in X_n) m_{\otimes_0}[X_i; v_n^s][S(X_i)]^{i-1}}{f_{X_i}} K \frac{\mu_{X_i, i, X_j}^{\mathfrak{t}}}{a_n} + o_p(n^{i-2})$$

Notice that the term in the \mathfrak{t} g bracket,

$$\frac{1}{(n_i - 1)a_n^s} \sum_{i \in j; i=1}^n \frac{1(X_i \in X_n) m_{\otimes_0}[X_i; v_n^s][S(X_i)]^{i-1}}{f_{X_i}} K \frac{\mu_{X_i, i, X_j}^{\mathfrak{t}}}{a_n};$$

is a kernel estimator of $1(X_j \in X_n) m_{\otimes_0}[X_j; v_n^s][S(X_j)]^{j-1}$: This kernel estimator is consistent (by Assumptions 1, 3, 4, 5, 11 and Remark A.1). Because $E[\frac{1}{2}(Z_j; \otimes_0(\mathfrak{t})) | X_1, X_2, \dots, X_n] = 0$, it follows from applying the Chebyshev inequality that

$$C = \frac{1}{n} \sum_{j=1}^n 1(X_j \in X_n) m_{\otimes_0}[X_j; v_n^s][S(X_j)]^{j-1} \frac{1}{2}(Z_j; \otimes_0) + o_p\left(\frac{1}{n}\right).$$

Also, by Remark A.1, $1(X_j \in X_n) \rightarrow 1$ and $E[\frac{1}{2}(Z_j; \otimes_0(\mathfrak{t})) | X_j] = 0$ and applying the Chebyshev inequality, we obtain:

$$C = \frac{1}{n} \sum_{j=1}^n m_{\otimes_0}[X_j; v_n^s][S(X_j)]^{j-1} \frac{1}{2}(Z_j; \otimes_0) + o_p\left(\frac{1}{n}\right).$$

This completes the proof. ■

In the rest of the proofs, we denote $\rho_n = o(n^{i-2})$; $\rho_n(\rho_n) = f^{\otimes 2} A_n : k_{\otimes_0}^{\otimes 2} \rho_n$, and $u^s = Sv_n^s$. Define $\rho_n = (1 - \rho_n)^{\otimes 2} + \rho_n(u^s + \otimes_0)$:

Remark A.2: (i) For any $\rho \in \rho_n(\rho_n)$, we have: $k_{\otimes_0}^{\otimes 2} \rho = O(\rho_n)$; $k_{\otimes_0}^{\otimes 2} \rho = O(\rho_n)$; $k_{\otimes_0}^{\otimes 2} \rho = O(\rho_n)$; $k_{\otimes_0}^{\otimes 2} \rho = O(\rho_n)$.
(ii) Under Assumption 9, for any $\rho \in \rho_n(\rho_n)$ we have: $k_{\otimes_0}^{\otimes 2} \rho = \rho_n O(\rho_n)$

Lemma A.2: Under Assumptions 1 - 9, 11 - 14(i), we obtain uniformly over $\mathbb{P}_{n, \sigma}^2$,

$$\mathbb{R}_n[\hat{\beta}_n] - \mathbb{R}_n[\beta_n] = \mathbb{R}_n[\hat{\beta}_n] - \mathbb{R}_n[\beta_n] + O_p(n^{-2});$$

Proof. By definition, we can write:

$$\begin{aligned} & \mathbb{R}_n[\hat{\beta}_n] - \mathbb{R}_n[\beta_n] = (\mathbb{R}_n[\hat{\beta}_n] - \mathbb{R}_n[\beta_n]) \\ & = \frac{1}{2n} \sum_{i=1}^n \mathbb{1}_{\{X_i \in X_n\}} (A_i + B_i + C_i + D_i); \end{aligned}$$

where $A_i \leq A_{1i} + A_{2i}$, with

$$\begin{aligned} A_{1i} &= m_{\sigma_0}[X_i; \beta_n]^{-1} \mathbb{1}_{\{X_i \in X_n\}} [S(X_i)]^{-1} m_{\sigma_0}[X_i; \beta_n] \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} \mathbb{1}_{\{X_i \in X_n\}} [S(X_i)]^{-1} m_{\sigma_0}[X_i; \beta_n]; \\ A_{2i} &= m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} m_{\sigma_0}[X_i; \beta_n] \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} m_{\sigma_0}[X_i; \beta_n]; \end{aligned}$$

$$B_i = 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} (b_m[X_i; \beta_n] - b_m[X_i; \beta_n]);$$

$C_i \leq C_{1i} + C_{2i}$, with

$$\begin{aligned} C_{1i} &= 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} (b_m[X_i; \beta_n] - b_m[X_i; \beta_n]) \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} b_m[X_i; \beta_n] \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} (b_m[X_i; \beta_n] - b_m[X_i; \beta_n]); \\ C_{2i} &= 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} (r_m[X_i; \beta_n] - r_m[X_i; \beta_n]) \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} r_m[X_i; \beta_n] \\ &\quad + 2m_{\sigma_0}[X_i; \beta_n]^{-1} [S(X_i)]^{-1} (r_m[X_i; \beta_n] - r_m[X_i; \beta_n]); \end{aligned}$$

and $D_i \leq D_{1i} + D_{2i}$, with

$$\begin{aligned} D_{1i} &= (b_m[X_i; \beta_n] - b_m[X_i; \beta_n])^{-1} [S(X_i)]^{-1} (b_m[X_i; \beta_n] - b_m[X_i; \beta_n]) \\ &\quad + 2(b_m[X_i; \beta_n] - b_m[X_i; \beta_n])^{-1} [S(X_i)]^{-1} b_m[X_i; \beta_n]; \\ D_{2i} &= (r_m[X_i; \beta_n] - r_m[X_i; \beta_n])^{-1} [S(X_i)]^{-1} (r_m[X_i; \beta_n] - r_m[X_i; \beta_n]) \\ &\quad + 2(r_m[X_i; \beta_n] - r_m[X_i; \beta_n])^{-1} [S(X_i)]^{-1} r_m[X_i; \beta_n]; \end{aligned}$$

Assumptions 1, 3(ii) and 11(ii) and Remark A.2 imply:

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] = O_p(n^{-2});$$

Corollary A.3 and Assumptions 1, 3(ii) and 11(ii) and Remark A.2 imply:

$$\frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] = O_p(n^{-2}).$$

Notice that

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] \\ & \quad - \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] \\ & = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g (A_{1i} + A_{2i} + A_{3i}); \end{aligned}$$

where

$$A_{1i} = (m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] - m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0]) [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0];$$

$$A_{2i} = m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} (m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] - m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0]);$$

$$A_{3i} = m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] [S(X_i)]^{i-1} [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0];$$

Applying Corollary A.3 and Assumptions 1, 3(ii) and 11(ii) and Remark A.2, we obtain by Holder inequality,

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g A_{1i} \\ & \leq \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g (m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] - m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0]) [S(X_i)]^{i-1} m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] \\ & \leq \frac{1}{n} \sum_{i=1}^n \frac{1}{2} \frac{1}{X_i} \frac{1}{2} X_n g m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] m_{\mathbb{R}_0} [X_i; \mathbb{R}_i | \mathbb{R}_0] \\ & = k_1 n^{-2} k_2 \leq O_p(n^{-2}) \leq k_1 n^{-2} k_2 \leq O_p(1) = O_p(n^{-2}). \end{aligned}$$

Using the same argument, we obtain $\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g A_{2i} = O_p(n^{-2})$ and $\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g A_{3i} = O_p(n^{-2})$. Hence, $\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g A_i = O_p(n^{-2})$:

Assumptions 1, 2, 3, 13 and Corollaries A.2, A.4 and Remark A.2 imply:

$$\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g B_i = O_p(n^{-2});$$

Assumptions 1, 3, 11, 12, 13 and Corollaries A.3, A.4, and Remark A.2 imply:

$$\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g C_i = O_p(n^{-2}).$$

Assumptions 1, 3, 12 - 14 and Corollary A.4 and Remark A.2 imply

$$\frac{1}{n} \sum_{i=1}^n 1(X_i \in X_n) g D_i = O_p(n^{-2} P_n^{-1}).$$

The lemma now follows. ■

For any function $g(Z)$, define $\frac{1}{n} \sum_{i=1}^n g(Z_i) = E[g(Z_i)]$:

Corollary A.8: Under Assumptions 1, 3, 8, 11 and 14(i), we have: uniformly over $X_i \in X_n$ and $\Delta \in U_n$,

$$\frac{1}{n} (1(X_i \in X_n) m_{\otimes_0}[X_i; \Delta]^0 [S(X_i)]^i m_{\otimes_0}[X_i; u^a]) = o_p(n^{-1/4}):$$

Proof. We can still apply Lemma A.1 with $G(Z; \Delta) = m_{\otimes_0}[X_i; \Delta]^0 [S(X_i)]^i m_{\otimes_0}[X_i; u^a]$ but without the bias term. Notice that Conditions A.1.2-A.1.4 are trivially satisfied given Assumptions 1, 3, 8, 11 and 14(i), and with $\pm_n = n^{-1/4}$, hence the result. ■

Lemma A.3: Under Assumptions 1 - 9, 11 - 14(i), we have uniformly over $\otimes \in \otimes_0(P_n^{-1})$,

$$\begin{aligned} & R[\otimes; \otimes_0] - R[\otimes; \otimes_0] = E \{ R[\otimes; \otimes_0] - R[\otimes; \otimes_0] \} \\ & L_n(\otimes) - L_n(\otimes) = E \{ L_n(\otimes) - L_n(\otimes) \} = O_p(n^{-2}); \end{aligned}$$

Proof. Using the same notation as that in the proof of Lemma A.2, we have:

$$R[\otimes; \otimes_0] - R[\otimes; \otimes_0] = \frac{1}{2n} \sum_{i=1}^n 1(X_i \in X_n) (A_{2i} + C_{2i} + D_{2i})$$

From the results (for A2_i; C2_i; D2_i) in the proof of Lemma A.2, we have:

$$\begin{aligned}
& i \left(R[\otimes_i \otimes_o] \right) R[\otimes_i \otimes_o] E f R[\otimes_i \otimes_o] R[\otimes_i \otimes_o] g \\
&= \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] + O_p(n^{-2}) \\
& i \in \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] \\
&= \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] + O_p(n^{-2}) \\
& i \in \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] \\
&= \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] + O_p(n^{-2}) \\
&= \frac{1}{n} \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] + O_p(n^{-2})
\end{aligned}$$

where the last three equalities follow from Remark A.2. The result now follows by applying Corollary A.8. ■

Lemma A.4: Under Assumptions 1, 2, 3, 9, 11(ii), 12, 13 and 14, we have uniformly over $\otimes \geq \otimes_o(P_{\otimes_n})$,

$$E[L_n(\otimes) - L_n(\otimes_o)] = \frac{1}{2} k_{\otimes}^2 - k_{\otimes_o}^2 + O(n^{-2});$$

Proof. Since $m(X; \otimes_o) = 0$; we can write

$$\begin{aligned}
i E[L_n(\otimes)] &= \frac{1}{2} E \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} m_{\otimes_o}[X_i; \otimes_i \otimes_o] \\
&+ E \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} r_m[X_i; \otimes_i \otimes_o] \\
&+ \frac{1}{2} E \sum_{i=1}^n 1(X_i \leq X_n) r_m[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} r_m[X_i; \otimes_i \otimes_o] :
\end{aligned}$$

Then

$$E[L_n(\otimes) - L_n(\otimes_o)] = \frac{1}{2} k_{\otimes}^2 - k_{\otimes_o}^2 = A_n + \frac{1}{2} B_n$$

where after some manipulation, we have

$$\begin{aligned}
& A_n \\
&= E \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} f r_m[X_i; \otimes_i \otimes_o] r_m[X_i; \otimes_i \otimes_o] g \\
&+ E \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} f r_m[X_i; \otimes_i \otimes_o] r_m[X_i; \otimes_i \otimes_o] g \\
&+ E \sum_{i=1}^n 1(X_i \leq X_n) m_{\otimes_o}[X_i; \otimes_i \otimes_o] [S(X_i)]^{i-1} r_m[X_i; \otimes_i \otimes_o]
\end{aligned}$$

and

$$B_n = E \left(\frac{1(X_i \neq X_n) \text{fr}_m[X_i; \lfloor n^{\otimes \alpha} \rfloor; \otimes_0] \text{r}_m[X_i; \otimes; \otimes_0] g^0 \varepsilon}{E[S(X_i)]^{i-1} \text{fr}_m[X_i; \lfloor n^{\otimes \alpha} \rfloor; \otimes_0] \text{r}_m[X_i; \otimes; \otimes_0] g} \right) \\ + 2E \frac{1(X_i \neq X_n) \text{r}_m[X_i; \otimes; \otimes_0] [S(X_i)]^{i-1} \text{fr}_m[X_i; \lfloor n^{\otimes \alpha} \rfloor; \otimes_0] \text{r}_m[X_i; \otimes; \otimes_0] g^0}{\dots}$$

Now Assumptions 11(ii), 12 - 14(i) together with Remark A.2 imply:

$$A_n = O_p(n^{-2}); \quad B_n = O_p(n^{-\frac{p}{n}}):$$

The result now follows. ■

Lemma A.5: Under Assumptions 1-9, 11, we have, uniformly over $\otimes \in \otimes_0(n^{-\frac{p}{n}})$

$$\hat{L}_{n^{\otimes_0}}[\otimes; \lfloor n^{\otimes \alpha} \rfloor] = O_p(n^{-2}) \text{ and } \hat{L}_{n^{\otimes_0}}[\otimes; \otimes_0] = O_p(n^{-1}):$$

Proof.

$$\hat{L}_{n^{\otimes_0}}[\otimes; \lfloor n^{\otimes \alpha} \rfloor] \\ = \frac{1}{n} \sum_{i=1}^n \mathfrak{m}_{\otimes_0}[X_i; \otimes; \lfloor n^{\otimes \alpha} \rfloor] [S(X_i)]^{i-1} \mathfrak{m}(X_i; \otimes_0) 1(X_i \neq X_n) \\ = k^{\otimes \alpha} \lfloor n^{\otimes \alpha} \rfloor k_2^{-\alpha} \frac{1}{n} \sum_{i=1}^n 1(X_i \neq X_n) \mathfrak{m}_{\otimes_0}[X_i; \frac{\otimes; \lfloor n^{\otimes \alpha} \rfloor}{k^{\otimes \alpha} \lfloor n^{\otimes \alpha} \rfloor k_2^{-\alpha}}] [S(X_i)]^{i-1} \mathfrak{m}(X_i; \otimes_0) \\ = O_p(n^{-\frac{p}{n}}) \in O_p(n^{-\frac{p}{n}}) = O_p(n^{-2})$$

where the last equality follows from applying Remark A.2 and using the same decomposition as that in the proof of Corollary A.5(ii). Also Corollary A.5(ii) directly implies $\hat{L}_{n^{\otimes_0}}[\otimes; \otimes_0] = O_p(n^{-1})$ uniformly over $\otimes \in \otimes_0(n^{-\frac{p}{n}})$. ■

Proof. (Theorem 4.1) Define $\mathfrak{G}^{\otimes} = (1; \lfloor n^{\otimes \alpha} \rfloor) \mathfrak{G} + \lfloor n^{\otimes \alpha} \rfloor (u^{\otimes} + \otimes_0)$: Simple calculations give:

$$\hat{L}_n(\mathfrak{G}) = \hat{L}_n(\lfloor n^{\otimes \alpha} \rfloor \mathfrak{G}^{\otimes}) + \hat{L}_{n^{\otimes_0}}[\mathfrak{G}; \lfloor n^{\otimes \alpha} \rfloor] + \mathfrak{R}[\mathfrak{G}; \otimes_0] \text{r}_m[\lfloor n^{\otimes \alpha} \rfloor; \otimes_0]:$$

By Lemma A.2, we have:

$$\hat{L}_n(\mathfrak{G}) = \hat{L}_n(\lfloor n^{\otimes \alpha} \rfloor \mathfrak{G}^{\otimes}) + \hat{L}_{n^{\otimes_0}}[\mathfrak{G}; \lfloor n^{\otimes \alpha} \rfloor] + \mathfrak{R}[\mathfrak{G}; \otimes_0] \text{r}_m[\lfloor n^{\otimes \alpha} \rfloor; \otimes_0] + o_p\left(\frac{1}{n}\right) \\ = \hat{L}_n(\lfloor n^{\otimes \alpha} \rfloor \mathfrak{G}^{\otimes}) + \hat{L}_{n^{\otimes_0}}[\mathfrak{G}; \lfloor n^{\otimes \alpha} \rfloor] + E[L_n(\mathfrak{G}) \text{r}_m[\lfloor n^{\otimes \alpha} \rfloor; \otimes_0]] + \\ + \mathfrak{R}[\mathfrak{G}; \otimes_0] \text{r}_m[\lfloor n^{\otimes \alpha} \rfloor; \otimes_0] + E \text{fr}_m[\mathfrak{G}; \otimes_0] \text{r}_m[\lfloor n^{\otimes \alpha} \rfloor; \otimes_0] + o_p\left(\frac{1}{n}\right)$$

where the last equality follows from rearranging terms. Lemmas A.3 and A.4 imply

$$\begin{aligned} \mathbb{E}_n(\mathbb{E}) &= \mathbb{E}_n(\mathbb{I}_n \mathbb{E}^\pi) + \frac{1}{2} k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 + k_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{I}_n \mathbb{E}^\pi}^2 \\ &\quad + \mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] + o_p\left(\frac{1}{n}\right); \end{aligned}$$

Write

$$\begin{aligned} k_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{I}_n \mathbb{E}^\pi}^2 &= k_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{I}_n \mathbb{E}^\pi}^2 + (1 - \pi)^2 k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 + \pi^2 k_{\mathbb{U}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{U}}^2 \\ &\quad + 2h_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi; (1 - \pi)(\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi) \mathbb{I}_n \\ &\quad + 2h_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi; \pi \mathbb{U} \mathbb{I}_n \mathbb{E}^\pi + 2h(1 - \pi)(\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi); \pi \mathbb{U} \mathbb{I}_n \mathbb{E}^\pi; \end{aligned}$$

By Theorem 3.1 and that $k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2$ is dominated by $k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2$, we have

$$jh_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi; \mathbb{E} \mathbb{I}_n \mathbb{E}^\pi; \mathbb{E} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 \in k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 = O_p(\pi^2)$$

Hence

$$k_{\mathbb{I}_n \mathbb{E}^\pi} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{I}_n \mathbb{E}^\pi}^2 = O_p(\pi^2) + (1 - \pi)^2 k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 + 2\pi(1 - \pi)h(\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi); \mathbb{U} \mathbb{I}_n \mathbb{E}^\pi$$

and

$$\begin{aligned} &\mathbb{E}_n(\mathbb{E}) + \mathbb{E}_n(\mathbb{I}_n \mathbb{E}^\pi) + \mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] \\ &= \mathbb{I}_n \frac{1}{2} \mathbb{I}_n (1 - \pi)^2 k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 + \pi(1 - \pi)h(\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi); \mathbb{U} \mathbb{I}_n \mathbb{E}^\pi + O_p(\pi^2); \end{aligned}$$

By definition of \mathbb{E} , $\mathbb{E}_n(\mathbb{E}) + \mathbb{E}_n(\mathbb{I}_n \mathbb{E}^\pi) \rightarrow 0$. And by Theorem 3.1, $k_{\mathbb{E}} \mathbb{I}_n \mathbb{E}^\pi k_{\mathbb{E}}^2 = O_p(\pi^2)$: It follows that

$$0 = \mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] + \pi(1 - \pi)h(\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi); \mathbb{U} \mathbb{I}_n \mathbb{E}^\pi + O_p(\pi^2);$$

Note that $\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi = \mathbb{I}_n (\mathbb{U}^\pi + \mathbb{E} \mathbb{I}_n \mathbb{E}^\pi) + \mathbb{E}^\pi \mathbb{I}_n \mathbb{E}^\pi$. Write

$$\mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] = \mathbb{I}_n \mathbb{E}_{n^{\otimes_0}}[\mathbb{U}^\pi] + \pi \mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] + \mathbb{E}_{n^{\otimes_0}}[\mathbb{E}^\pi \mathbb{I}_n \mathbb{E}^\pi];$$

Applying Lemma A.5, we obtain

$$\mathbb{E}_{n^{\otimes_0}}[\mathbb{E} \mathbb{I}_n \mathbb{E}^\pi] = \mathbb{I}_n \mathbb{E}_{n^{\otimes_0}}[\mathbb{U}^\pi] + o_p\left(\frac{1}{n}\right)$$

and

$$0 \leq \frac{1}{n} \mathbb{E}_{n^{\otimes o}}[u^{\otimes}] + \frac{1}{n} (1 - \frac{1}{n}) h_{\otimes}^{\otimes} i^{\otimes o}; u^{\otimes} i_n + O_p(\frac{1}{n^2}):$$

Therefore,

$$0 \leq \mathbb{E}_{n^{\otimes o}}[u^{\otimes}] i (1 - \frac{1}{n}) h_{\otimes}^{\otimes} i^{\otimes o}; u^{\otimes} i_n + O_p(\frac{1}{n}):$$

Since this holds for $u^{\otimes} = \mathbb{S}v_n^{\otimes}$, we obtain:

$$\mathbb{E}_{n^{\otimes o}}[v_n^{\otimes}] i (1 - \frac{1}{n}) h_{\otimes}^{\otimes} i^{\otimes o}; v_n^{\otimes} i_n = O_p(\frac{1}{n}):$$

This proves that

$$\mathbb{E}_{n^{\otimes o}}[v_n^{\otimes}] = h_{\otimes}^{\otimes} i^{\otimes o}; v_n^{\otimes} i_n + o_p(\frac{1}{n}):$$

Applying Corollary A.7, we have:

$$\mathbb{E}_{n^{\otimes o}}[v_n^{\otimes}] = \frac{1}{n} \sum_{i=1}^{\infty} m_{\otimes o}[X_i; v^{\otimes}]^{\otimes} [S(X_i)]^i \frac{1}{2}(Z_i; \otimes_o) + o_p(\frac{1}{n})$$

Hence for any fixed non-zero $\mu \in \mathbb{R}^b$,

$$\mathbb{E}_{\otimes}(\beta_n i \mu_o) = \frac{1}{n} \sum_{i=1}^{\infty} m_{\otimes o}[X_i; v^{\otimes}]^{\otimes} S(X_i)^i \frac{1}{2}(Z_i; \otimes_o) + o_p(\frac{1}{n}):$$

Substituting for v^{\otimes} , we obtain

$$\begin{aligned} & \mathbb{E}_{\otimes}(\beta_n i \mu_o) \\ &= \frac{1}{n} \mathbb{E} \left[D_{w^{\otimes}}(X)^{\otimes} (S(X))^i \frac{1}{2} D_{w^{\otimes}}(X)^{\otimes} i \frac{1}{n} \sum_{i=1}^{\infty} D_{w^{\otimes}}(X_i)^{\otimes} (S(X_i))^i \frac{1}{2}(Z_i; \otimes_o) + o_p(1) \right] \end{aligned}$$

The theorem now follows from applying a standard CLT for i.i.d. data. ■

Appendix D: Asymptotic Variance and Efficiency

In this appendix, we first prove the consistency of $\hat{\psi}$; then show the efficiency of the optimally weighted minimum distance estimator, and finally establish the asymptotic properties of the three-step estimator.

Proof. (Theorem 5.1): To prove the consistency of the estimated covariance matrix, it suffices to show that $\hat{S}_o(X_i)$ and $\hat{D}_{w^k}(X_i)$ are consistent estimators of $S_o(X_i)$ and $D_{w^k}(X_i)$ uniformly over $X_i \in X_n$. Define

$$\hat{S}_o(X_i) = \frac{[(n-1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{1}{2}(Z_j; \otimes) \frac{1}{2}(Z_j; \otimes)^0 K \frac{X_i X_j}{a_n}}{f_{X_i}}$$

We apply Lemma A.1 with $G(Z_j; \otimes) = \frac{1}{2}(Z_j; \otimes) \frac{1}{2}(Z_j; \otimes)^0$ and $\pm_n = 1$, Conditions A.1.1-A.1.4 are satisfied by Assumptions 6, 8 and 16. Hence by Lemma A.1 and Corollary A.1 we have uniformly over $X_i \in X_n$ and $\otimes \in A_n$:

$$\hat{S}_o(X_i) = S_o(X_i) + o_p(1);$$

This together with Assumptions 1, 5, 7, Corollary A.1 and $k_{\otimes_n i} \otimes_o k_2 = o_p(n^{1-4})$ implies that uniformly over $X_i \in X_n$:

$$\hat{S}_o(X_i) = S_o(X_i) + o_p(1);$$

To prove the consistency of $\hat{D}_{w^k}(X_i)$, note that Assumption 15(i) implies uniformly over $X_i \in X_n$:

$$\begin{aligned} & \frac{[(n-1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{\tilde{A} \otimes \frac{1}{2}(Z_j; \otimes_n)}{\otimes_{\mu_k}} \frac{\otimes \frac{1}{2}(Z_j; \otimes_o)}{\otimes_{\mu_k}} K \frac{\mu X_i X_j}{a_n}}{k_{\otimes_n i} \otimes_o k_2} \approx \frac{[(n-1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{\tilde{A}_8(Z_j)}{\otimes_{\mu_k}} K \frac{\mu X_i X_j}{a_n}}{k_{\otimes_n i} \otimes_o k_2} = o_p(n^{1-4}) \end{aligned}$$

where the last equality follows from $k_{\otimes_n i} \otimes_o k_2 = o_p(n^{1-4})$: Applying Lemma A.1 with $G(Z_j; \otimes) = \frac{\otimes \frac{1}{2}(Z_j; \otimes)}{\otimes_{\mu_k}}$ and $\pm_n = 1$, we have uniformly over $X_i \in X_n$:

$$[(n-1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{\otimes \frac{1}{2}(Z_j; \otimes_o)}{\otimes_{\mu_k}} K \frac{\mu X_i X_j}{a_n} = E\left[\frac{\otimes \frac{1}{2}(Z; \otimes_o)}{\otimes_{\mu_k}} \middle| X = X_i\right] + o_p(1);$$

Hence

$$[(n-1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{\otimes \frac{1}{2}(Z_j; \otimes_n)}{\otimes_{\mu_k}} K \frac{\mu X_i X_j}{a_n} = E\left[\frac{\otimes \frac{1}{2}(Z; \otimes_o)}{\otimes_{\mu_k}} \middle| X = X_i\right] + o_p(1);$$

Since the objective function

$$\sum_{i=1}^n 1(X_i \in X_n) \hat{D}_{w^k}(X_i) [\hat{S}(X_i)]^{-1} \hat{D}_{w^k}(X_i)$$

is globally convex in $w^k(\cdot)$, the solution $\hat{w}^k(\cdot)$ is bounded by $\hat{w}^k(\cdot) \leq C$ for some constant C . Thus, we only need to concern with the subset $\{A \in \overline{W}_n : \|A\|_2 \leq C\}$. Assumption 15(ii) implies uniformly over $X_i \in X_n$, $A \in \overline{W}_n$ and $\|A\|_2 \leq C$:

$$[(n_i - 1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{1}{h_n} [Z_j; A(\cdot)] \frac{1}{h_0} [Z_j; A(\cdot)] \leq \frac{\sum_{i=1}^n X_i X_j}{a_n} = o_p(n^{i-4}):$$

Moreover, Corollary A.3 implies uniformly over $X_i \in X_n$, $A \in \overline{W}_n$ and $\|A\|_2 \leq C$:

$$[(n_i - 1)a_n^s]^{-1} \sum_{j \in i; j=1}^n \frac{1}{h_0} [Z_j; A(\cdot)] \leq \frac{\sum_{i=1}^n X_i X_j}{a_n} = E[\frac{1}{h_0} [Z; A(\cdot)] | X = X_i] + o_p(1):$$

these results imply $\mathcal{D}_{w^k}(X_i) = D_{w^k}(X_i) + o_p(1)$ uniformly over $X_i \in X_n$. Employing the same proof as of Theorem 3.1, we show that $\hat{w}^k(\cdot) - w^k(\cdot) = o_p(1)$. Hence, $\mathcal{D}_{\hat{w}^k}(X_i) = D_{w^k}(X_i) + o_p(1)$ uniformly over $X_i \in X_n$. The theorem now follows immediately. ■

Proof. (Theorem 6.1): Since the semiparametric efficiency bound is the lower bound of the variances of all \sqrt{n} regular (consistent) semiparametric estimators and V_0^{-1} is the asymptotic variance of a semiparametric estimator, it follows that V_0^{-1} must be no smaller than the efficiency bound in positive semi-definite matrix sense. To show that it also is the semiparametric efficiency bound, we only need to show that there is a sequence of parametric submodels whose asymptotic covariance matrices converge to V_0^{-1} : Recall that $w^\pi(\cdot) \in \overline{W}$ and $\int w^\pi(\cdot)(\mu - \mu_0) \in H \perp h_0$. Because \overline{W} is a completion of W under the Hilbert norm, there is a sequence $w_s(\cdot) \in W$ for $s = 1, 2, \dots$ such that $w_s(\cdot) \rightarrow w^\pi(\cdot)$ as $s \rightarrow \infty$. Correspondingly, we have $\int w_s(\cdot)(\mu - \mu_0) \in H \perp h_0$. Hence, $g_s(\cdot; \mu) = h_0 + \int w_s(\cdot)(\mu - \mu_0) \in H$ is a parametric submodel passing through h_0 . For this particular submodel, Chamberlain (1987) showed that the semiparametric efficiency bound is the inverse of $V_s = E[D_{w_s}(X)^T [S_0(X)]^{-1} D_{w_s}(X)]$, where $D_w(X)$ is given in (4.1). Letting $s \rightarrow \infty$ and noting that $w_s \rightarrow w^\pi(\cdot)$ and that $D_w(X)$ is a continuous functional of $w(\cdot)$, it follows that $V_s \rightarrow V_0$. This proves that V_0^{-1} is the semiparametric efficiency bound for model (1.3). ■

Proof. (Theorem 6.2): (i) We apply Lemma A.1 with $G(Z_j; \theta) = \frac{1}{2}(Z_j; \theta) \frac{1}{2}(Z_j; \theta)^T$ and $\pm_n = n^{i-4}$, Conditions A.1.1-A.1.4 are satisfied by Assumptions 8, 16 and 17. Hence by Lemma A.1 and Corollary A.1 we have uniformly over $X_i \in X_n$ and

② A_n :

$$\hat{S}_n(X_i) = S_n(X_i) + o_p(n^{-1/4}):$$

Notice that the estimator $\hat{\theta}_n$ in Step 1 has the convergence rate $\|\hat{\theta}_n - \theta_0\|_2 = o_p(n^{-1/4})$ under Assumptions 1, 2, 4-12. These together with Assumptions 1, 14(ii) and 16(ii) imply that uniformly over $X_i \in \mathcal{X}_n$:

$$\hat{S}_n(X_i) = S_n(X_i) + o_p(n^{-1/4}):$$

(ii) A direct consequence of result (i) and Theorem 4.1. ■