

A Simple Ordered Data Estimator for Inverse Density Weighted Expectations

Arthur Lewbel* and Susanne M. Schennach
Boston College and University of Chicago

Abstract

We consider estimation of means of functions that are scaled by an unknown density, or equivalently, integrals of conditional expectations. The “ordered data” estimator we provide is root n consistent, asymptotically normal, and is numerically extremely simple, involving little more than ordering the data and summing the results. No sample size dependent smoothing is required. A similarly simple estimator is provided for the limiting variance. The proofs include new limiting distribution results for functions of nearest neighbor spacings. Potential applications include endogenous binary choice, willingness to pay, selection, and treatment models.

JEL codes: C14, C21, C25, J24

Keywords: Semiparametric, Conditional Expectation, Density Estimation, Binary Choice, Binomial Response.

*Corresponding Author: Arthur Lewbel, Department of Economics, Boston College, 140 Commonwealth Ave., Chestnut Hill, MA, 02467, USA. phone: 617-552-3678, fax: 617-552-2308, lewbel@bc.edu, <http://www2.bc.edu/~lewbel/>

1 Introduction

Given an iid dataset $(x_1, y_1), \dots, (x_n, y_n)$, on a bivariate random variable (x, y) , we propose to estimate an expectation of the inverse density weighted form

$$\theta = E \left[\frac{y}{f(x)} \right] \quad (1)$$

where $f(x)$ is the unknown density of the continuously distributed x . The “ordered data” estimator we provide possesses the rather surprising property of achieving root n consistency and asymptotic normality without requiring sample-size-dependent smoothing. It also offers the advantage of being numerically extremely simple, requiring little more than ordering the data and summing the results, thereby avoiding issues regarding the selection of smoothers such as bandwidths, kernels, etc. A similarly simple estimator is provided for the limiting variance.

Inverse density weighted estimation applies generically to the estimation of definite integrals of conditional expectations. Suppose

$$\theta = \int_{x \in \mathcal{X}} E(w|x) dx \quad (2)$$

for some random variable w , where $\mathcal{X} \subset \text{supp}(x)$. Then θ can be rewritten as $\theta = E[y/f(x)]$ with $y = wI(x \in \mathcal{X})$, and where I is the indicator function that equals one if its argument is true and zero otherwise.

A number of existing semiparametric estimators make use of inverse density weighted expectations either directly or indirectly, via their relationship with integrated conditional expectations. Examples include density weighted least squares (Newey and Ruud, 1984), average derivative estimation (Härdle and Stoker, 1993), estimators of willingness-to-pay models and general estimators of moments from binomial data (Lewbel, 1997, McFadden, 1999, Lewbel, Linton and McFadden, 2002), semiparametric estimators of consumer surplus (Hausman and Newey, 1995, Newey, 1997), some discrete choice, sample selection, and other latent variable model estimators (Lewbel, 1998, 2000, 2002), entropy measures of dependence (Hong and White, 2000) and semiparametric functional tests (Hall and Yatchew, 2005).

Let $(y_{[i]}, x_{[i]})$ denote the i -th observation when the data are sorted in increasing order of x , so $x_{[i]}$ is the i -th order statistic and $y_{[i]}$ is the concomitant statistic to $x_{[i]}$. Let (y_i, x_i) denote the i -th observation when the data are left unsorted. Let $F(x)$ denote the unknown distribution function of x . We show that the

numerically trivial “ordered data” estimator

$$\hat{\theta} = \sum_{i=1}^{n-1} (y_{[i+1]} + y_{[i]}) (x_{[i+1]} - x_{[i]}) / 2 \quad (3)$$

is root n consistent and asymptotically normal. Specifically, $n^{1/2}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 3\sigma^2/2)$ where

$$\sigma^2 = E \left[\frac{\text{Var}(y|x)}{f^2(x)} \right] \quad (4)$$

and a consistent estimator of σ^2 is the simple expression

$$\hat{\sigma}^2 = \frac{n}{4} \sum_{i=1}^{n-1} (y_{[i+1]} - y_{[i]})^2 (x_{[i+1]} - x_{[i]})^2. \quad (5)$$

For some intuition for this estimator, let x be a point that lies between $x_{[i]}$ and $x_{[i+1]}$ for some i , and let y be a corresponding point that lies between $y_{[i]}$ and $y_{[i+1]}$. Then $y \approx (y_{[i+1]} + y_{[i]}) / 2$ and

$$\frac{1}{f(x)} = \left(\frac{dF(x)}{dx} \right)^{-1} \approx \frac{x_{[i+1]} - x_{[i]}}{F(x_{[i+1]}) - F(x_{[i]})} \approx \frac{x_{[i+1]} - x_{[i]}}{1/n} \quad (6)$$

where the last step results from replacing F with the corresponding empirical distribution function. The estimator $\hat{\theta}$ is then just an average of $y/f(x)$, using these approximations for y and $f(x)$. This differencing of the empirical distribution function for x does not yield a consistent estimator for $1/f(x)$, but averaging over x speeds the rate of convergence to yield root n consistency. This result is rather exceptional among semiparametric estimators of nonlinear functionals of the data generating process in that it attains root n consistency without the use of explicit smoothing techniques.

The ability to bypass kernel and bandwidth selection procedures constitutes a substantial benefit of the proposed estimator. The well-known bandwidth selection rules used in nonparametric estimation, such as cross-validation, are not generally applicable to semiparametric settings because the optimal bandwidth in nonparametric settings typically fails to undersmooth at the rate needed to achieve a $o(n^{-1/2})$ bias in the semiparametric functional. In practice, a bandwidth “slightly” smaller than the one given by cross-validation is often used instead, in an attempt to undersmooth. As the required amount of undersmoothing is a sample-size-dependent quantity, this informal method does not provide much guidance and may lead to a nonnegligible bias. The derivation and application of formal data-driven bandwidth selection rules in semiparametric settings is rarely done, as it involves technical higher-order asymptotic analyses (see,

for instance, Härdle *et al.*, 1992, Hall and Horowitz, 1990) that must be rederived for each semiparametric estimator of interest, and may depend strongly on the unknown precise degree of smoothness of the estimand.

The method for approximating the reciprocal of a density using longer spacings to achieve consistency is known (see, e.g., Bloch and Gastwirth, 1968). Our contribution consists mainly of providing limiting distribution theory for an average of such estimators, for fixed or increasing spacings. The complication that arises in doing so is accounting for the fact that the spacings $x_{[i+1]} - x_{[i]}$ are not independent. In fact, the statistical dependence among them is of a form that is not covered by standard central limit theorems for dependent processes. One spacing depends equally strongly on arbitrarily distant spacings and not only on its neighbors. To handle this difficulty, we substantially extend a technique of proof proposed by Weiss (1958) in the case of homogenous functions of spacings of a sample drawn from a uniform distribution to cover the more general functional (1).

While the use of nearest neighbors in differencing schemes in semiparametric¹ settings is not new (e.g., Yatchew, 1997), our use of a statistic based on nearest neighbor *spacings* is innovative. The asymptotic analysis of differencing techniques only relies on the fact that spacings converge to zero as sample size increases, while our results require an analysis of their asymptotic distribution.

In addition to providing the limiting distribution of an estimator of (1), we provide an extension of the limiting distribution theory to the case where x itself is estimated and consider a generalizations to multivariate y . We also show how the ordered data estimator could be used in some semiparametric models, we provide a small Monte Carlo analysis of the estimator, and we apply the estimator in a small empirical study. Proofs are in the Appendix.

2 Asymptotics

2.1 Main results

The derivation of the asymptotic properties of our “ordered data” estimator relies on a few standard assumptions.

Condition 1 (y_i, x_i) is a sequence of *i.i.d.* random variables.

¹The use of differencing using k -nearest neighbors, where $k \rightarrow \infty$ as $n \rightarrow \infty$ is also well-established (see, e.g. Robinson, 1987) but our asymptotic analysis under fixed k is fundamentally different.

Condition 2 The support of $f(x)$, denoted \mathcal{F} , is a finite interval and $\inf_{x \in \mathcal{F}} f(x) > 0$.

Condition 3 $f(x)$ is uniformly Hölder continuous of exponent $h_f > 1/2$ over the interior of its support (i.e. $|f(x) - f(\xi)| \leq H_f |x - \xi|^{h_f}$ for $x, \xi \in \mathcal{F}$).

Condition 4 $g(x) = E[y|x]$ is uniformly Hölder continuous of exponent $h_g > 1/2$ over the interior of the support of $f(x)$ (i.e. $|g(x) - g(\xi)| \leq H_g |x - \xi|^{h_g}$ for $x, \xi \in \mathcal{F}$).

Condition 5 $E \left[\frac{\text{Var}[y|x]}{f^2(x)} \right] < \infty$.

Assumption 1 is very common in cross-sectional analysis. Assumptions 3 and 4 impose smoothness constraints on $f(x)$ and $E[y|x]$ that are slightly stronger than continuity and weaker than assuming that $f(x)$ and $E[y|x]$ are Lipschitz (which would correspond to $h_f = h_g = 1$). Assumption 5 is crucial in order to obtain root n consistency but could probably be relaxed using sample size-dependent trimming if only consistency is desired. Assumption 2 is frequently made in the literature focusing on inverse density weighted estimators. It can probably be relaxed at the expense of substantial complications in the proofs, however, we do not pursue such extensions in the present work, because Assumption 5 is rarely satisfied when $f(x)$ is not bounded away from zero.²

Theorem 1 Under Assumptions 1 through 5, $n^{1/2}(\tilde{\theta} - \theta) \xrightarrow{d} N(0, 2\sigma^2)$ where

$$\tilde{\theta} = \sum_{i=1}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i]}) \tag{7}$$

$$\theta = E \left[\frac{y}{f(x)} \right] = \int E[y|x] dx \tag{8}$$

$$\sigma^2 = E \left[\frac{\text{Var}[y|x]}{f^2(x)} \right]. \tag{9}$$

²Another paper that employs this support restriction, for similar reasons, is Abadie and Imbens (2002). This support restriction can likely be relaxed using trimming or boundary methods similar to those used for kernel estimators, though doing so may require the introduction of data dependent parameters (such as asymptotic trimming terms). One of the main advantages of the estimator is that it does not require selection of data dependent parameters like bandwidths. Alternatively, it is always possible to map a variable with unbounded support onto one with a bounded support by using a fixed nonlinear mapping. Our estimator will be root n consistent for this transformed variable, provided that the Jacobian of this known transformation is included in our inverse-density weighted expectation. This would just amount to redefining the y variable and verifying that the resulting $E[y|x]$ is sufficiently smooth and bounded.

Corollary 2 Under Assumptions 1 through 5, $n^{1/2} (\hat{\theta} - \theta) \xrightarrow{d} N(0, (3/2) \sigma^2)$ where

$$\hat{\theta} = \sum_{i=1}^{n-1} (y_{[i]} + y_{[i+1]}) (x_{[i+1]} - x_{[i]}) / 2. \quad (10)$$

and where θ and σ^2 are defined in Theorem 1.

The fundamental reason why root n consistency is possible without smoothing is because spacings directly estimate the reciprocal of the density and θ is a linear functional of the reciprocal of the density, thus allowing substantial undersmoothing without introducing nonnegligible nonlinear remainder terms. Kernel-based estimation would first proceed by estimating the density and then take its reciprocal, a nonlinear operation that precludes the use of such substantial undersmoothing. The need for smoothing in nonlinear functional estimation, but not in linear functional estimation, parallels Newey's linearization condition (Newey, 1994, Assumption 5.1), which places different constraints on the rate of convergence of plugged-in nonparametric estimates, depending on whether the functional of interest is linear or not.

The factor 2 in the variance of $\tilde{\theta}$ was also noted by Mack and Müller (1988) in the related problem of estimating a nonparametric conditional expectation using the Priestley and Chao (1972) estimator. See also Wand and Jones (1995, p. 131). Corollary 2 reduces this factor to 3/2. This factor can be further reduced to arbitrarily close to one by considering longer spacings. For instance, it can be shown along the same lines as Theorem 1 that

$$\tilde{\theta}_k = \sum_{i=1}^{n-k} y_{[i]} (x_{[i+k]} - x_{[i]}) / k. \quad (11)$$

has an asymptotic variance equal to $\sigma^2 (1 + \frac{1}{k})$. Note that for $k = 2$, Equation (11) is asymptotically equivalent to Equation (10).³ Centered differences can also be used, e.g.

$$\hat{\theta}_c = \sum_{i=2}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i-1]}) / 2. \quad (12)$$

is asymptotically equivalent to both $\hat{\theta}$ and $\tilde{\theta}_2$.

The following theorem shows that when the length of the spacings grows with sample size, the asymptotic variance reaches the value σ^2 , which can be shown to be semiparametrically efficient using the variance expressions in Newey (1994).⁴

³This can be shown along the following lines $\sum_{i=1}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i]}) / 2 + \sum_{i=1}^{n-1} y_{[i+1]} (x_{[i+1]} - x_{[i]}) / 2 = \sum_{i=1}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i]}) / 2 + \sum_{i=2}^n y_{[i]} (x_{[i]} - x_{[i-1]}) / 2 \approx \sum_{i=2}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i-1]}) \approx \sum_{i=1}^{n-2} y_{[i]} (x_{[i+2]} - x_{[i]})$.

⁴Let the quantities defined in Newey (1994) be denoted by a N subscript. Our estimator can be written in terms of a nonparametric estimate $h_N(x)$ of the conditional expectation $E[y|x]$ as $\hat{\theta} = \int h_N(x) dx = \int \frac{h_N(x)}{f(x)} f(x) dx = E \left[\frac{h_N(x)}{f(x)} \right] =$

Theorem 3 Let k_n be a deterministic sequence of integers such that $k_n = o(\ln n)$ and $k_n \rightarrow \infty$. Under Assumptions 1 through 5, $n^{1/2} (\tilde{\theta}_{k_n} - \theta) \xrightarrow{d} N(0, \sigma^2)$ where $\tilde{\theta}_k$ is given in Equation (11) and where θ and σ^2 are defined in Theorem 1.

While Theorem 3 requires an upper bound on the rate at which k_n can grow, it imposes *no lower bound* (i.e. k_n can diverge arbitrarily slowly). This is unusual among semiparametric estimators, which typically require both an upper bound and a lower bound, in order to control both the bias and the variance. Here, the variance is already finite even when no smoothing is performed and hence, the mild requirement that $k_n \rightarrow \infty$ is sufficient to reduce the asymptotic variance term $2\sigma^2$ to its efficient value of σ^2 .

The variance term σ^2 can be consistently estimated by⁵

$$\hat{\sigma}^2 = \frac{n}{4} \sum_{i=1}^{n-1} (y_{[i+1]} - y_{[i]})^2 (x_{[i+1]} - x_{[i]})^2, \quad (13)$$

as established in the following Theorem.

Theorem 4 Under Assumptions 1 through 5, $\hat{\sigma}^2 \xrightarrow{p} \sigma^2$.

One important point needs to be emphasized regarding the asymptotic variance σ^2 . If f were known, the most obvious estimator of θ would be $\hat{\theta}_0 \equiv n^{-1} \sum_{i=1}^n y_i / f(x_i)$. An interesting feature of this model is that $\hat{\theta}_0$ is inefficient, that is, it is more efficient to plug an appropriate estimator for $f(x_i)$ into this sample average than to use the true f . This can be immediately verified by

$$\begin{aligned} \text{Var}[y/f(x)] &= E[\text{Var}[y|x]/f^2(x)] + \text{Var}[E(y|x)/f(x)] \\ &= \sigma^2 + \text{Var}[E(y|x)/f(x)] \geq \sigma^2. \end{aligned} \quad (14)$$

Similar efficiency gains from using estimated rather than true functions has been observed in a model that is scaled by a conditional density function (see Magnac and Maurin, 2003) and in a model scaled by a propensity score function (see Hirano, Imbens and Ridder, 2003).

$E[m_N(z_N, h_N)]$ where $z_N = (x, y)$ and $m_N(z_N, h_N) = \frac{h_N(x)}{f(x)}$. Since $m_N(z_N, h_N)$ is linear in $h_N(x)$, Equation 4.4 in Newey (1994) holds for $\delta_N(x) = (f(x))^{-1}$ and Proposition 4 yields the influence function $\alpha_N(z) = (f(x))^{-1}(y - E[y|x])$, the variance of which is $E[\text{Var}[y|x]/(f(x))^2] \equiv \sigma^2$.

⁵The idea of employing $(y_{[i+1]} - y_{[i]})$ to “differentiate out” the contribution of $E[y|x]$ has been used for instance by Yatchew (1997) to estimate the partially linear model and by Abadie and Imbens (2002) to estimate the variance of a matching estimator of treatment effects models. However, our variance estimator combines differencing and spacings, making it impossible to use these existing results.

2.2 Extensions

2.2.1 Multivariate case

The result presented so far can be extended in various directions. First, the scalar y can be replaced with a vector, as follows.

Corollary 5 *Let v_i for $i = 1, \dots, n$ denote observations from random vector v taking value in \mathbb{R}^k and let $v_{[i]}$ denote the concomitant statistic to $x_{[i]}$. If Assumptions 1 through 5 hold for $y_i = v_{i1}, \dots, v_{ik}$, then $n^{1/2}(\hat{\psi} - \psi) \xrightarrow{d} N(0, (3/2)V)$ where $\hat{\psi} = \sum_{i=1}^{n-1} (v_{[i]} + v_{[i+1]})(x_{[i+1]} - x_{[i]})/2$, $\psi = E[v/f(x)] = \int E[v|x] dx$ and $V = E[(E[vv'|x] - E[v|x]E[v'|x])f^{-2}(x)]$. Moreover, $\hat{V} \xrightarrow{p} V$, where*

$$\hat{V} = \frac{n}{4} \sum_{i=1}^{n-1} (v_{[i+1]} - v_{[i]})(v_{[i+1]} - v_{[i]})'(x_{[i+1]} - x_{[i]})^2. \quad (15)$$

2.2.2 Estimated x

In some applications, it is useful to be able to use an estimated value of x . While the ordered data estimator requires no modification *per se* to handle this generated regressor, its asymptotic variance needs to account for the error in the estimation of x in the first step. For this purpose, let x be a random scalar and w be a random vector taking a value in \mathbb{R}^{N_w} that are related through $x = X(w, \gamma)$ for some function $X : \mathbb{R}^{N_w} \times \mathbb{R}^{N_\gamma} \mapsto \mathbb{R}$ and for some parameter vector γ whose true value is γ^* . Letting $\theta(\gamma) = E[y/f_{x|\gamma}(x|\gamma)]$, we can then establish the following result.

Condition 6 *The support of $f_{x|\gamma}(x|\gamma)$, denoted \mathcal{F}_γ , is a finite interval and $\inf_{\gamma \in \Gamma} \inf_{x \in \mathcal{F}_\gamma} f_{x|\gamma}(x|\gamma) > 0$, where Γ is some neighborhood of γ^* .*

Condition 7 $|f_{x|\gamma}(x|\gamma) - f_{x|\gamma}(\xi|\gamma)| \leq H_f |x - \xi|^{h_f}$ for $x, \xi \in \mathcal{F}_\gamma$ and any $\gamma \in \Gamma$.

Condition 8 $|g(x) - g(\xi)| \leq H_g |x - \xi|^{h_g}$ for $x, \xi \in \cup_{\gamma \in \Gamma} \mathcal{F}_\gamma$.

Condition 9 $E[\sup_{\gamma \in \Gamma} \text{Var}[y|x]/f_{x|\gamma}^2(x|\gamma)] < \infty$.

Condition 10 $E[\sup_{\gamma \in \Gamma} |\partial X(w, \gamma)/\partial \gamma|] < \infty$.

Condition 11 $\partial E[y/f_{x|\gamma}(X(w, \gamma)|\gamma)]/\partial \gamma = E[\partial(y/f_{x|\gamma}(X(w, \gamma)|\gamma))/\partial \gamma]$ for all $\gamma \in \Gamma$ and $\partial E[y/f_{x|\gamma}(X(w, \gamma)|\gamma)]/\partial \gamma$ is continuous in γ at γ^* .

Theorem 6 Let $\tilde{\theta}_{k_n}(\gamma)$ and k_n respectively denote the estimator and the sequence described in Theorem 3.

Let $\hat{\gamma}$ be a consistent estimate of γ^* with influence function⁶ Ψ_i . Under Assumptions 1 and 6 through 11,

$n^{1/2} \left(\tilde{\theta}_{k_n}(\hat{\gamma}) - \theta(\gamma^*) \right) \rightarrow N(0, V)$ where

$$V = E \left[\frac{(y - E[y|w, \gamma^*])^2}{f_{x|\gamma}^2(x|\gamma^*)} \right] + 2B'E \left[\frac{(y - E[y|w, \gamma^*])}{f_{x|\gamma}(x|\gamma^*)} \Psi_i \right] + B'E [\Psi_i \Psi_i'] B \quad (16)$$

$$B = E \left[\frac{y}{f_{x|\gamma}(x|\gamma^*)} \frac{\partial}{\partial x} E \left[\frac{\partial X(w, \gamma^*)}{\partial \gamma} | x \right] \right] + \\ - E \left[\frac{y}{(f_{x|\gamma}(x|\gamma^*))^2} \frac{\partial f_{x|\gamma}(x|\gamma^*)}{\partial x} \left(\frac{\partial X(w, \gamma^*)}{\partial \gamma} - E \left[\frac{\partial X(w, \gamma^*)}{\partial \gamma} | x \right] \right) \right]. \quad (17)$$

Moreover, if $E[\partial X(w, \gamma^*)/\partial \gamma | x, y] = E[\partial X(w, \gamma^*)/\partial \gamma]$, then $B = 0$.

Assumptions 6 through 9 replace Assumptions 2 through 5 to ensure that the regularity conditions hold uniformly for $\gamma \in \Gamma$, since $\tilde{\theta}_{k_n}(\gamma)$ is a random function instead of a random scalar.

Theorem 6 assumes $k_n \rightarrow \infty$, omitting the case where the spacing width k is fixed. With fixed spacings the function $\tilde{\theta}_k(\gamma)$ exhibits a peculiar mode of convergence to $\theta(\gamma)$, in that the derivative $d\tilde{\theta}_k(\gamma)/d\gamma$ diverges as $n \rightarrow \infty$ almost everywhere and changes discontinuously every time the order of the $x_{[i]}$ changes. As a result, we do not know if $n^{1/2} \left(\tilde{\theta}_k(\hat{\gamma}) - \theta(\hat{\gamma}) \right) - n^{1/2} \left(\tilde{\theta}_k(\gamma^*) - \theta(\gamma^*) \right) \xrightarrow{p} 0$ with fixed k . The function $\tilde{\theta}_{k_n}(\gamma)$ with $k_n \rightarrow \infty$ does not exhibit these problems.

3 Monte Carlo

Here we provide a Monte Carlo analysis to assess the small sample behavior of the estimator. We draw x_i and e_i as independent standard normals and let $y_i = 2x_i(1 + e_i)I(0 < x_i < 1)$. Table 1 reports results from estimating $\theta = E[y/f(x)] = 1$ using this simulated data. The sample size is $n = 100$, the number of replications is 10,000, and the reported summary statistics are, respectively, the mean, standard deviation, quartiles (lower, median, upper), root mean squared error, mean absolute error, and median absolute error.

We report results for seven estimators. The first is $\hat{\theta}_0 = n^{-1} \sum_{i=1}^n y_i/f(x_i)$, an estimator that uses the true, normal, density function $f(x)$, which in a typical application would be unknown. The next is $\hat{\theta}_1 = \sum_{i=1}^{n-1} (y_{[i+1]} + y_{[i]}) (x_{[i+1]} - x_{[i]}) / 2$, the ordered data estimator of Corollary 2. We also report $\hat{\theta}_k = \sum_{i=1}^{n-k} (y_{[i+k]} + y_{[i]}) (x_{[i+k]} - x_{[i]}) / (2k)$ for $k = 2$ and $k = 3$, which are ordered data estimators with wider

⁶That is, $n^{1/2}(\hat{\gamma} - \gamma^*) = n^{1/2} \sum_{i=1}^n \Psi_i + o_p(1)$ where Ψ_i is i.i.d with mean zero finite variance.

spacings. The next estimator is $\hat{\theta}_4 = n^{-1} \sum_{i=1}^n y_i / \hat{f}(x_i, b)$ where $\hat{f}(x, b)$ is a kernel density estimator, using a quartic kernel and bandwidth b given by Silverman's rule of thumb. Finally, $\hat{\theta}_5$ and $\hat{\theta}_6$ are the same kernel density based estimators, except using bandwidths $b/2$ and $2b$, respectively.

In terms of mean squared, mean absolute, or median absolute error, the estimator having the best fit is $\hat{\theta}_5$, though this kernel estimator also has the most mean and median bias. This estimator is undersmoothed relative to (Silverman's approximation of) pointwise optimality. Efficiency and root n convergence of plug-in kernel estimators require undersmoothing relative to pointwise optimality of the nonparametric density estimator (see, e.g., Newey, 1994).

The ordered data estimator $\hat{\theta}_1$ had the smallest mean and median bias of all the estimators, but somewhat larger mean squared errors. These errors decrease as expected as the spacings increase. The estimator using the true density, $\hat{\theta}_0$, has largest errors, at least in part reflecting the inefficiency of that estimator as discussed earlier.

4 Examples

4.1 Latent Moments From Binomial Data

Consider a model of the form $d_i = I(w_i > x_i)$, where d is an observed dummy variable, x is an observed continuously distributed random variable, and w is an unobserved latent random variable that is drawn from a distribution that is independent of x . The goal is estimation of moments of w .

Problems like this arise in survey research, where w is an attribute of an individual such as wealth or willingness to pay for a public good. The individual is asked if w exceeds some randomly chosen value x , and d is the response. This form of survey design is used, because it is likely to produce less biased responses than directly asking for w (see, e.g., McFadden, 1999 and reference therein).

This model may also be applied in destructive testing, e.g., w could be the speed at which a car safety device fails, x would be the speed at which the car was tested, and d an indicator of outcome failure, such as whether a test dummy was injured. Similarly, in bioassay w might be the time required for an animal to suffer an abnormality, x is the time at which the animal is sacrificed to test for the abnormality, d is indicator of the test result for abnormality at time x .

Let $\theta = E(w^\lambda) - c^\lambda$ where λ is a moment chosen for estimation and c is any chosen element of the support

of x (e.g., c could be the median of x). A special case of results in Lewbel, Linton, and McFadden (2002) is

$$\theta = E\left(\frac{\lambda x^{\lambda-1}[d - I(x < c)]}{f(x)}\right), \quad (18)$$

assuming that $\text{supp}(w) \subset \text{supp}(x)$.⁷ Therefore, letting $y_i = \lambda x_i^{\lambda-1}[d_i - I(x_i < c)]$, a simple estimator of $E(w^\lambda)$ is $\hat{\theta} + c^\lambda$ where $\hat{\theta} = \sum_{i=1}^{n-1} (y_{[i+1]} + y_{[i]}) (x_{[i+1]} - x_{[i]}) / 2$, $n^{1/2}(\hat{\theta} - \theta) \xrightarrow{d} N(0, 3\sigma^2/2)$, and $\hat{\sigma}^2 = \sum_{i=1}^{n-1} (y_{[i+1]} - y_{[i]})^2 (x_{[i+1]} - x_{[i]})^2 n/4$.

4.2 Selection and Treatment Effects

Consider the model $y_i = y_i^* d_i$, $d_i = I(0 \leq w_i + x_i \leq a)$, where a is a constant (which could equal infinity), y_i is an individual's observed outcome, d is an observed dummy variable that indicates if the individual is selected or treated, x is an observed continuously distributed random variable, and y^* and w are unobserved latent random variables that are drawn from a distribution that is independent of x . The goal is estimation of moments of the potential outcome y^* .

An example is a wage model, where y^* is an individual's wage if employed, y is the individual's observed wage, d is the indicator of whether an individual is employed, a is infinite, $-x$ is some form of nonwage income such as a government defined benefit, and w is a latent variable such that the individual chooses to work if $w + x$ is sufficiently large. More generally, in a treatment context y^* is an individual's outcome if treated (the potential outcome), d is a treatment indicator, and x is a variable that only affects the decision to treat but not the outcome if treated. An example with a finite a (two sided censoring) would be ordered treatment, where the latent $w_i + x_i$ determines the treatment, which if negative would indicate a lesser (or no) treatment and if greater than a would indicate a stronger treatment, where the possible treatments are, e.g., dosages of a drug or years of schooling.

Define ω by

$$\omega = \frac{E[y/f(x)]}{E[d/f(x)]} \quad (19)$$

Theorem 1 and Corollary 1 in Lewbel (2002) show that, if a is finite and x has a sufficiently large support, then $\omega = E(y^*)$, and that even without these assumptions, $\omega \approx E(y^*)$.

Both the numerator and denominator of ω are inverse density weighted means, so we propose the simple

⁷See also Lewbel (1997) and McFadden (1999).

estimator

$$\hat{\omega} = \frac{\sum_{i=1}^{n-1} (y_{[i+1]} + y_{[i]}) (x_{[i+1]} - x_{[i]})}{\sum_{i=1}^{n-1} (d_{[i+1]} + d_{[i]}) (x_{[i+1]} - x_{[i]})} \quad (20)$$

The limiting distribution for $\hat{\omega}$ is obtained by applying Corollary 5 with $v = (y, d)$ to obtain the joint distribution of the numerator and denominator of $\hat{\omega}$, then applying the delta method. The result is $n^{1/2} (\hat{\omega} - \omega) \xrightarrow{d} N(0, (3/2)s^2)$, where

$$s^2 = \left[E \left(\frac{d}{f(x)} \right) \right]^{-2} E \left[\frac{\text{Var}[(y - \omega d)|x]}{f^2(x)} \right] \quad (21)$$

which can be consistently estimated by

$$\hat{s}^2 = \frac{n \sum_{i=1}^{n-1} [(y_{[i+1]} - \hat{\omega} d_{[i+1]}) - (y_{[i]} - \hat{\omega} d_{[i]})]^2 (x_{[i+1]} - x_{[i]})^2}{\left[\sum_{i=1}^{n-1} (d_{[i]} + d_{[i+1]}) (x_{[i+1]} - x_{[i]}) \right]^2} \quad (22)$$

4.3 Endogeneous Binary Choice Models

Consider the binary choice or binomial response model

$$d_i = I(x_i + z_i' \beta + e_i \geq 0) \quad (23)$$

Where for each individual i , d_i is an observed zero or one outcome, z_i is a vector of possibly endogeneous observed regressors, x_i is an observed scalar regressor with coefficient normalized to equal one, β is a vector of coefficients to be estimated, and e_i is an unobserved error term. If e is independent of x, z then β can be estimated either parametrically by maximum likelihood if the distribution of e is known or semiparametrically using, e.g., Klein and Spady (1993).

Suppose that the joint distribution of e, z is unknown, but we observe a vector of instrumental variables r_i such that $E(re) = 0$. If x is independent of z, r, e then, given some regularity and support assumptions⁸, Lewbel (2000) shows that

$$E(rz')\beta = E \left(r \frac{d - I(x > 0)}{f(x)} \right) \quad (24)$$

so β can be consistently estimated by a linear two stage least squares regression of $[d - I(x > 0)]/f(x)$ on z using instruments r . Note that r and z may contain a constant term, so the regression includes estimation of location in $z'\beta$ assuming e has unconditional mean zero.

⁸In particular, x should have a continuous distribution, be demeaned or otherwise located to contain zero in its support, and contain the support of $-(z'\beta + e)$ in its support. Magnac and Maurin (2003) discuss alternative restrictions.

Corollary 7 below provides a simple method for implementing this result, by applying linear two stage least squares and using Corollary 5 to deal with the density in the right side of equation (24).

Make the following definitions. Let $v_i = r_i[d_i - I(x_i > 0)]$, $\psi = E[v/f(x)]$, $\Sigma_{zr} = E(zr')$, $\Sigma_{rr} = E(rr')$, and $\Delta = (\Sigma_{zr}\Sigma_{rr}^{-1}\Sigma'_{zr})^{-1}\Sigma_{zr}\Sigma_{rr}^{-1}$. Assuming the inverses in the definition of Δ exists, it follows from equation (24) that $\beta = \Delta\psi$. The corresponding finite sample expressions are $\widehat{\Sigma}_{zr} = \sum_{i=1}^{n-1} z_{[i]}r'_{[i]}/n$, $\widehat{\Sigma}_{rr} = \sum_{i=1}^{n-1} r_{[i]}r'_{[i]}/n$, $\widehat{\psi} = \sum_{i=1}^{n-1} (v_{[i]} + v_{[i+1]}) (x_{[i+1]} - x_{[i]}) / 2$, $\widehat{\beta} = \widehat{\Delta}\widehat{\psi}$.and

$$\widehat{\Delta} = (\widehat{\Sigma}_{zr}\widehat{\Sigma}_{rr}^{-1}\widehat{\Sigma}'_{zr})^{-1}\widehat{\Sigma}_{zr}\widehat{\Sigma}_{rr}^{-1}. \quad (25)$$

Corollary 7 *Given the above definitions, if Assumptions 1 through 5 hold for $y_i = v_{i1}, \dots, v_{ik}$, then $n^{1/2}(\widehat{\beta} - \beta) \xrightarrow{d} N(0, \Omega)$ and $\widehat{\Omega} \xrightarrow{p} \Omega$ where $\Omega = \Delta E[(W + U)(W + U)'] \Delta'$, $W = (rz' - \Sigma'_{zr})\beta$, $U = y - E[y|x]$, $\widehat{\Omega} = n^{-1}\widehat{\Delta} \left(\sum_{i=1}^{n-1} \frac{3}{8}\widehat{U}_{[i]}\widehat{U}'_{[i]} + 2\widehat{U}_{[i]}\widehat{W}'_{[i]} + \widehat{W}_{[i]}\widehat{W}'_{[i]} \right) \widehat{\Delta}'$, $\widehat{W}_{[i]} = (r_{[i]}z'_{[i]} - \widehat{\Sigma}'_{zr})\widehat{\beta}$ and $\widehat{U}_{[i]} = (y_{[i]} - y_{[i+1]}) (x_{[i+1]} - x_{[i]}) n$.*

As noted by Lewbel (2000), the binary choice or binomial response model just described is identified under more general settings. In particular, the binary outcome can be generated from

$$d_i = I(\tilde{x}_i + z'_i\beta + e_i \geq 0) \quad (26)$$

where \tilde{x} does not need to be independent from z and r . In this case, a consistent estimation method consists in constructing a new variable x from the residuals of the least-squares projection of \tilde{x} on z and r . Provided that the residual x is independent from z and r , the normal equation 24 still holds and the methodology outlined in the beginning of this section applies. However, the asymptotic variance and this estimator is affected by the fact that the variable x is estimated in a first step. Theorem 6 can be used to handle this situation, by making the following identifications between the quantities used in the statement of the Theorem (on the left-hand side) and the quantities introduced in the present section (on the right-hand side):

$$w_i = (z'_i, r'_i) \quad (27)$$

$$\hat{\gamma} = \left(\sum_{i=1}^n w_i w'_i \right)^{-1} \sum_{i=1}^n w_i \tilde{x}_i \quad (28)$$

$$X(w, \gamma) = w'_i \gamma. \quad (29)$$

5 Empirical Application

We now apply the results of the previous section to estimation of a model by Cogneau and Maurin (2002) on the effects of parental income on school attendance in Madagascar. The data set consists of a representative sample of 1401 children aged six to eight, from a World Bank survey conducted in 1993-1994. The model is equation (23) where d_i equals one if child i is enrolled in school on time (that is, by age 6), and zero otherwise, x_i is the date of birth of the child in the relevant year, normalized to vary from $-1/2$ to $1/2$, and the other regressors z_i are a constant, the child's sex, parents' income, and the mother's education level (Cogneau and Maurin report virtually no difference between using mother's or father's education level).

The latent error e_i may be correlated with parent's income and education level, because of family fixed effects (common unobserved determinants of parent's resources and decisions) and measurement errors, in that parent's observed income and education level are rough proxies for permanent income and other measures of total household resources. To control for this endogeneity, instruments r_i are defined as the child's sex, differences between parents and grandparents education levels, and the difference between fathers and grandfathers sector of labor activity (agriculture vs nonagriculture). For more details about the data, the model, and alternative estimators, see Cogneau and Maurin (2002).

Table 2 below reports results using two different estimators, both based on equation (24) and the implication that $\beta = \Delta\psi$. The first estimator is from Lewbel (2000), which uses the same estimator of Δ as equation (25), except that observation $[n]$ is not omitted, and estimates ψ as $n^{-1} \sum_{i=1}^n v_i / \hat{f}(v_i)$, where $\hat{f}(v_i)$ is a quartic kernel estimator. Bandwidth selection and kernel based standard error estimates for this estimator are constructed using the methods described in Lewbel (2000). The second estimator in Table 2 is the ordered data estimator in Corollary 7.

Both the kernel and ordered data estimators are equivalent to a linear two stage least squares regression of an estimate of $[d - I(x > 0)]/f(x)$ on regressors z using instruments. Each estimator is applied twice in Table 2. The first application of each uses z as instruments as well as regressors, and so is equivalent to an ordinary least squares instead of a two stage least squares, which fails to control for endogeneity. The second application of each estimator uses r as described above as instruments.⁹

⁹The first-step F statistics of the regression of each element of x on z are relatively large (greater than 50), indicating that our results do not suffer from a weak instrument problem (see Equations (3.5) and (3.7) in Staiger and Stock, 1997).

The kernel and ordered data based estimates are generally quite similar. The kernel based two stage least squares estimates are all within two standard errors of the ordered data based two stage least squares estimates, and vice versa. The main empirical finding is that controlling for endogeneity more than quadruples the estimated effect of parental income on the decision to start children’s schooling on time. After controlling for income and endogeneity, the effects of mother’s education level and the sex of the child are small and statistically insignificant

6 Conclusion

We provide the limiting root n distribution for a simple “ordered data” estimator of means of functions that are scaled by an unknown density, or equivalently, integrals of conditional expectations. We show that the ordered data estimator is a viable alternative to more complicated estimators that require smoothing parameters such as kernels and bandwidths.

Our asymptotic distribution theory is complicated by the fact that the dependence among sorted data spacings $x_{[i+1]} - x_{[i]}$ is of a form that is not covered by standard central limit theorems for dependent processes. Each spacing depends equally strongly on arbitrarily distant spacings and not only on its neighbors. We substantially extend Weiss (1958) to derive asymptotic distribution theory for these spacings.

Although our “ordered data” approach only covers the case where x is scalar, it should be noted that the majority of empirical applications involving nonparametric density estimation are likely to be univariate, e.g., the popular STATA econometrics package has built in commands for univariate, but not multivariate, nonparametric kernel density and kernel regression estimation.

It is possible to extend the approach presented here to obtain consistent estimators for K -dimensional x by replacing spacings by suitable functions of first-nearest neighbors distances (details are available from the authors upon request). However, achieving root n consistency involves handling some of the technical difficulties that also plague semiparametric multivariate kernel estimation. In conventional kernel estimators, boundary effects introduce a $O(h)$ bias, where h is the bandwidth (Cheng *et al.*, 1997), regardless of the order of the kernel. By analogy, since the effective bandwidth in a nearest-neighbor estimator is $O(n^{-1/K})$, the bias is $O(n^{-1/K})$, which is sufficiently large to prevent root n consistent estimation with no bias in the

asymptotic distribution¹⁰ for $K > 1$. In the context of kernel estimation, this bias is dealt with using some form of asymptotic trimming or kernel refinement, such as a local polynomial kernel smoother (Stone, 1977 and Cleveland, 1979). An analogous local nearest-neighbor polynomial estimator would consist of regressing y_j on a polynomial in x_j for all x_j in some $O(n^{-1/K})$ neighborhood of x_i and using the resulting predicted value \hat{y}_i in a sample average weighted by nearest-neighbor distances to the power K . It would be interesting to see if such extensions can attain root n consistency with a fixed number of neighbors in a multivariate setting.

Appendix A

Lemma 8 *Let (i) g_i be an iid sequence drawn from a $e^{-g_i}1(g_i \geq 0)$ density, (ii) u_i be an iid sequence drawn from a uniform density on $[0, 1]$, and (iii) z_i be a sequence such that the following quantities are defined: $\mu = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i]$, $\rho^2 = \lim_{n \rightarrow \infty} \text{Var} [n^{-1/2} \sum_{i=1}^n z_i]$, and $\tau^2 = \lim_{n \rightarrow \infty} n^{-1} (\sum_{i=1}^n E[z_i^2])$. If (i) z_i is independent from g_i and $(u_{[i+1]} - u_{[i]})$ and (ii) if $n^{-1/2} (\sum_{i=1}^n z_i g_i - n\mu)$ is asymptotically normal, then*

$$n^{-1/2} \left(\sum_{i=1}^{n-1} z_{[i]} n (u_{[i+1]} - u_{[i]}) - n\mu \right) \xrightarrow{d} N(0, \rho^2 + \tau^2 - \mu^2). \quad (30)$$

Proof. This proof is analogous to Weiss' (1958) derivation of the asymptotic distribution of homogenous functions of spacings. Let $G_{[i]} = \sum_{j=0}^{i+1} g_{[j]}$ and $s_{[i]} = g_{[i]}/G_n$ for $i = 1, \dots, n-1$. It can be shown (see Weiss, 1958) that the joint distribution of the $s_{[i]}$ is identical to the one of the $u_{[i+1]} - u_{[i]}$. The desired result can thus be established by relating the distribution of $n^{-1} \sum_{i=1}^{n-1} z_{[i]} g_{[i]}$ to the one of $n^{-1} \sum_{i=1}^{n-1} z_{[i]} s_{[i]}$.

Let us first calculate the mean and the variance of $n^{-1} \sum_{i=1}^{n-1} z_{[i]} g_{[i]}$. We have $\lim_{n \rightarrow \infty} E \left[n^{-1} \sum_{i=1}^{n-1} z_{[i]} g_{[i]} \right]$
 $= \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^{n-1} E[z_{[i]} g_{[i]}] = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i g_i] = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i] E[g_i] = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i] \cdot 1 = \mu$. Also, $\lim_{n \rightarrow \infty} \text{Var} \left[n^{-1/2} \sum_{i=1}^{n-1} z_{[i]} g_{[i]} \right] = \lim_{n \rightarrow \infty} n^{-1} E \left[\left(\sum_{i=1}^n z_i g_i \right)^2 \right] - \mu^2 = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i^2] E[g_i^2] + \sum_{i=1}^n \sum_{j \neq i} E[z_i z_j] E[g_i] E[g_j] - \mu^2 = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[z_i^2] \cdot 2 + \sum_{i=1}^n \sum_{j \neq i} E[z_i z_j] \cdot 1 \cdot 1 - \mu^2 = \lim_{n \rightarrow \infty} n^{-1} (\sum_{i=1}^n E[z_i^2]) + n^{-1} \left(\sum_{i=1}^n E[z_i^2] + \sum_{i=1}^n \sum_{j \neq i} E[z_i z_j] \right) - \mu^2 = \tau^2 + \lim_{n \rightarrow \infty}$

¹⁰The same issue of $O(n^{-1/K})$ bias in an otherwise root n consistent estimator was noted in a nearest neighbor matching estimator investigated by Abadie and Imbens (2002), although their estimator does not make use of nearest neighbor distances to evaluate an inverse density.

$E \left[\left(n^{-1/2} \sum_{i=1}^n z_i \right)^2 \right] - \mu^2 = \tau^2 + \rho^2$. Thus,

$$X_0 = \frac{\sum_{i=1}^{n-1} z_{[i]} g_{[i]} - n\mu}{n^{1/2}\gamma} \quad (31)$$

where $\gamma = \sqrt{\tau^2 + \rho^2}$, has a $N(0, 1)$ asymptotic distribution. However, X_0 can also be written (see Weiss, 1958) as

$$\begin{aligned} X_0 &= \frac{\sum_{i=1}^{n-1} z_{[i]} g_{[i]} - G_n \mu + G_n \mu - n\mu}{n^{1/2}\gamma} \\ &= \frac{\sum_{i=1}^{n-1} z_{[i]} g_{[i]} - G_n \mu}{\left(\frac{G_n}{n}\right) n^{1/2}\gamma} + \frac{G_n \mu - n\mu}{\left(\frac{G_n}{n}\right) n^{1/2}\gamma} + O_p(n^{-1}) \\ &= \frac{\sum_{i=1}^{n-1} z_{[i]} n (g_{[i]}/G_n) - n\mu}{n^{1/2}\gamma} + \frac{G_n \mu - n\mu}{\left(\frac{G_n}{n}\right) n^{1/2}\gamma} + O_p(n^{-1}) \\ &= \frac{\sum_{i=1}^{n-1} z_i n s_{[i]} - n\mu}{n^{1/2}\gamma} + \frac{G_n \mu - n\mu}{\left(\frac{G_n}{n}\right) n^{1/2}\gamma} + O_p(n^{-1}) \\ &\equiv X_1 + X_2 + O_p(n^{-1}) \end{aligned} \quad (32)$$

where X_1 and X_2 are independent (because it can be shown (see Weiss, 1958) that d_i and G_n are independent). Moreover,

$$X_2 = \frac{G_n \mu - n\mu}{\left(\frac{G_n}{n}\right) n^{1/2}\gamma} = \left(\frac{G_n - n}{n^{1/2}}\right) \left(\frac{\mu}{\left(\frac{G_n}{n}\right) \gamma}\right) \xrightarrow{p} \left(\frac{G_n - n}{n^{1/2}}\right) \left(\frac{\mu}{\gamma}\right). \quad (33)$$

It follows that $X_2 \xrightarrow{d} N\left(0, (\mu/\gamma)^2\right)$. Since (i) $X_0 = X_1 + X_2$, (ii) $X_0 \xrightarrow{d} N(0, 1)$, (iii) $X_2 \xrightarrow{d} N\left(0, (\mu/\gamma)^2\right)$ and (iv) X_1 and X_2 are independent, it follows that the asymptotic distribution of X_1 is the one of X_0 , “deconvoluted” by the one of X_2 . For independent normals, the deconvolution operation simply amounts to subtracting the variances. Thus, $X_1 \xrightarrow{d} N\left(0, 1 - (\mu/\gamma)^2\right)$. Since $X_1 = \left(\sum_{i=1}^{n-1} z_{[i]} n s_{[i]} - n\mu\right) / (n^{1/2}\gamma)$, then $n^{-1/2} \left(\sum_{i=1}^{n-1} z_{[i]} n s_{[i]} - n\mu\right) \xrightarrow{d} N\left(0, \gamma^2 - \mu^2\right)$, where $\gamma^2 - \mu^2 = \rho^2 + \tau^2 - \mu^2$. ■

Lemma 9 *Let g_i be an iid sequence drawn from a $e^{-g_i} 1(g_i \geq 0)$ density and let u_i be an iid sequence drawn from a uniform density on $[0, 1]$. If (i) z_i is independent from g_i , and $(u_{[i+1]} - u_{[i]})$ and (ii) if $n^{-1} \sum_{i=1}^n z_i g_i^\alpha \xrightarrow{p} \mu$, then $n^{-1} \sum_{i=1}^n z_{[i]} n^\alpha (u_{[i+1]} - u_{[i]})^\alpha \xrightarrow{p} \mu$.*

Proof. Let $G_{[i]} = \sum_{j=0}^{i+1} g_{[j]}$. As noted by Weiss (1958), $n^{-1} \sum_{i=1}^{n-1} z_{[i]} n^\alpha (u_{[i+1]} - u_{[i]})^\alpha$ has the same distribution as $n^{-1} \sum_{i=1}^n z_{[i]} n^\alpha (g_{[i]}/G_{[n]})^\alpha = n^{-1} \sum_{i=1}^n z_i n^\alpha (g_i/G_{[n]})^\alpha = (n/G_{[n]})^\alpha n^{-1} \sum_{i=1}^{n-1} z_i g_i^\alpha = (1 + o_p(1)) n^{-1} \sum_{i=1}^n z_i g_i^\alpha$, where $n^{-1} \sum_{i=1}^n z_i g_i^\alpha \xrightarrow{p} \mu$, by assumption. ■

Lemma 10 *If x_i is an iid sequence drawn from a continuous density $f(x)$ satisfying assumption 2, then*

$$\sum_{i=1}^{n-1} (x_{[i+1]} - x_{[i]})^\alpha = O_p(n^{1-\alpha}) \text{ for } \alpha > 0.$$

Proof. Let $\underline{f} = \inf_x f(x)$ and $u_i = F(x_i)$. By Assumption 2, $\underline{f} > 0$ and the inverse function $F^{-1}(\cdot)$ is uniquely defined. We then have, by the mean value theorem and the continuity of $f(x)$,

$$\begin{aligned} & \sum_{i=1}^{n-1} (x_{[i+1]} - x_{[i]})^\alpha = \sum_{i=1}^{n-1} (F^{-1}(u_{[i+1]}) - F^{-1}(u_{[i]}))^\alpha \\ &= \sum_{i=1}^{n-1} \left(f(\xi_{[i]}) \right)^{-\alpha} (u_{[i+1]} - u_{[i]})^\alpha \text{ for some } \xi_{[i]} \in [u_{[i]}, u_{[i+1]}] \\ &\leq \underline{f}^{-\alpha} \sum_{i=1}^{n-1} (u_{[i+1]} - u_{[i]})^\alpha = \underline{f}^{-\alpha} n^{1-\alpha} \left(n^{-1} \sum_{i=1}^{n-1} n^\alpha (u_{[i+1]} - u_{[i]})^\alpha \right) \\ &= \underline{f}^{-\alpha} n^{1-\alpha} O_p(1) = O_p(n^{1-\alpha}) \end{aligned} \tag{34}$$

where the second to last equality follows from Lemma 9. ■

Lemma 11 *If $\text{Var}[a_i] < \infty$, then $\sup_{i \in \{1, \dots, n\}} a_i = O_p(n^{1/2})$.*

Proof. Combining $P[A \cup B] \leq P[A] + P[B]$ with Tschebychev's inequality, we have that $P[a_i \geq Cn^{1/2} \text{ for some } i \leq n] \leq \sum_{i=1}^n P[a_i \geq Cn^{1/2}] \leq \sum_{i=1}^n \text{Var}[a_i] C^{-2} n^{-1} = \text{Var}[a_i] C^{-2}$, which can be made arbitrarily small for all n by choosing a C sufficiently large. ■

Proof. (of Theorem 1) Let $g(x) = E[y|x]$, $\Delta y_{[i]} = y_{[i]} - g(x_{[i]})$, $d_{[i]} = F(x_{[i+1]}) - F(x_{[i]})$ and write $\tilde{\theta} - \theta = N_1 + R_1 + R_2$ where

$$N_1 = n^{-1} \sum_{i=1}^{n-1} \frac{\Delta y_{[i]}}{f(x_{[i]})} n d_{[i]} \tag{35}$$

$$R_1 = \sum_{i=1}^{n-1} g(x_{[i]}) (x_{[i+1]} - x_{[i]}) - \int g(x) dx \tag{36}$$

$$R_2 = \sum_{i=1}^{n-1} \Delta y_{[i]} (x_{[i+1]} - x_{[i]}) - \sum_{i=1}^{n-1} \frac{\Delta y_{[i]}}{f(x_{[i]})} d_{[i]}. \tag{37}$$

R_1 can be bounded in probability using the mean value theorem, the Hölder properties of $g(x)$, and Lemma

$$\begin{aligned} 10: |R_1| &= \left| \sum_{i=1}^{n-1} g(x_{[i]}) (x_{[i+1]} - x_{[i]}) - \int g(x) dx \right| = \left| \sum_{i=1}^{n-1} \left(g(x_{[i]}) (x_{[i+1]} - x_{[i]}) - \int_{x_{[i]}}^{x_{[i+1]}} g(x) dx \right) \right| \\ &= \left| \sum_{i=1}^{n-1} \left(g(x_{[i]}) (x_{[i+1]} - x_{[i]}) - g(\xi_{[i]}) (x_{[i+1]} - x_{[i]}) \right) \right| \text{ for some } \xi_{[i]} \in [x_{[i]}, x_{[i+1]}]. \\ &\text{We thus have } |R_1| \leq \sum_{i=1}^{n-1} \left| g(x_{[i]}) - g(\xi_{[i]}) \right| (x_{[i+1]} - x_{[i]}) \leq \sum_{i=1}^{n-1} H_g \left| \xi_{[i]} - x_{[i]} \right|^{h_g} (x_{[i+1]} - x_{[i]}) \leq H_g \sum_{i=1}^{n-1} (x_{[i+1]} - x_{[i]})^{1+h_g} \\ &= O_p(n^{-h_g}) = o_p(n^{-1/2}). \end{aligned}$$

The second remainder term R_2 can be similarly bounded with the help of the Cauchy-Schwartz inequality:

$$\begin{aligned}
|R_2| &= \left| \sum_{i=1}^{n-1} \Delta y_{[i]} (x_{[i+1]} - x_{[i]}) - \sum_{i=1}^{n-1} \frac{\Delta y_{[i]}}{f(x_{[i]})} d_{[i]} \right| \\
&= \left| \sum_{i=1}^{n-1} \frac{\Delta y_{[i]}}{f(x_{[i]})} ((x_{[i+1]} - x_{[i]}) f(x_{[i]}) - (F(x_{[i+1]}) - F(x_{[i]}))) \right| \\
&= \left| \sum_{i=1}^{n-1} \frac{\Delta y_{[i]}}{f(x_{[i]})} (x_{[i+1]} - x_{[i]}) (f(x_{[i]}) - f(\xi_{[i]})) \right| \text{ for some } \xi_{[i]} \in [x_{[i]}, x_{[i+1]}] \\
&\leq \sum_{i=1}^{n-1} \frac{|\Delta y_{[i]}|}{f(x_{[i]})} (x_{[i+1]} - x_{[i]}) |f(x_{[i]}) - f(\xi_{[i]})| \leq H_f n^{-1} \sum_{i=1}^{n-1} \frac{|\Delta y_{[i]}|}{f(x_{[i]})} n (x_{[i+1]} - x_{[i]})^{1+h_f} \\
&\leq H_f \left(n^{-1} \sum_{i=1}^{n-1} \frac{(\Delta y_{[i]})^2}{f^2(x_{[i]})} \right)^{1/2} \left(n^{-1} \sum_{i=1}^{n-1} n^2 (x_{[i+1]} - x_{[i]})^{2+2h_f} \right)^{1/2} \\
&= H_f \left(n^{-1} \sum_{i=1}^{n-1} \frac{(\Delta y_{[i]})^2}{f^2(x_{[i]})} \right)^{1/2} \left(n \sum_{i=1}^{n-1} (x_{[i+1]} - x_{[i]})^{2+2h_f} \right)^{1/2} \\
&= O_p(1) O_p(n^{-h_f}) = o_p(n^{-1/2}). \tag{38}
\end{aligned}$$

We now show that the N_1 term is asymptotically normal and root n consistent. One cannot simply use the Lindeberg-Levy CLT to determine the asymptotics of this sum, because the $d_{[i]}$ are dependent. However, by Lemma 8, one can still achieve asymptotic normality and root n consistency, if two requirements are met: (i) $\Delta y_{[i]}/f(x_{[i]})$ is independent from $d_{[i]}$ and (ii) $n^{-1} \sum_{i=1}^n \Delta y_i g_i / f(x_i)$ is asymptotically normal and root n consistent, where g_i is an iid sequence independent from $\Delta y_i / f(x_i)$ and drawn from a $e^{-g_i} 1 (g_i \geq 0)$ density.

To show that $(\Delta y_{[i]}, x_{[i]})$ is asymptotically independent (a.i.) from $d_{[i]}$, we use the fact, shown by Barbe (1994), that $u_{[i]}$ and $d_{[i]}$ are a.i. First note that $x_{[i]}$ and $d_{[i]}$ are also a.i. since $x_{[i]} = F^{-1}(u_{[i]})$. Then, to see that $\Delta y_{[i]}$ and $d_{[i]}$ are a.i., observe that, asymptotically, $P[\Delta y_{[i]} | x_{[i]}, d_{[i]}] = P[\Delta y_{[i]} | x_{[i]}, x_{[i+1]}] = P[\Delta y_{[i]} | x_{[i]}]$.

Now, $n^{-1} \sum_{i=1}^{n-1} \Delta y_{[i]} g_{[i]} / f(x_{[i]}) \xrightarrow{p} n^{-1} \sum_{i=1}^n \Delta y_i g_i / f(x_i) = n^{-1} \sum_{i=1}^n \Delta y_i g_i / f(x_i)$, which can be shown to be asymptotically normal by the Lindeberg-Levi CLT, since all the variables are iid and since $E[(\Delta y g / f(x))^2] = E[(\Delta y / f(x))^2] E[g^2] = E[\text{Var}(y|x) / f^2(x)] \cdot 2 < \infty$ by assumption.

Since both conditions of Lemma 8 are met, we can conclude that the sum $n^{-1} \sum_{i=1}^{n-1} (\Delta y_{[i]} / f(x_{[i]})) n d_{[i]}$ is asymptotically normal with mean $\mu = \lim_{n \rightarrow \infty} n^{-1} \sum_{i=1}^n E[\Delta y_i / f(x_i)] = 0$ and variance $\rho^2 + \tau^2 - \mu^2$ where $\rho^2 = \lim_{n \rightarrow \infty} \text{Var}[n^{-1/2} \sum_{i=1}^n \Delta y_i / f(x_i)] = E[(\Delta y_i / f(x_i))^2] = E[\text{Var}[y_i | x_i] / f^2(x_i)] = \sigma^2$ and $\tau^2 = \lim_{n \rightarrow \infty} n^{-1} \left(\sum_{i=1}^n E[(\Delta y_i / f(x_i))^2] \right) = \sigma^2$. ■

Proof. (of Corollary 2) First observe that

$$\hat{\theta} = \sum_{i=1}^{n-1} y_{[i]} (x_{[i+1]} - x_{[i]}) / 2 + \sum_{i=1}^{n-1} y_{[i+1]} (x_{[i+1]} - x_{[i]}) / 2 \equiv T_1 + T_2. \quad (39)$$

Theorem 1 directly implies that the T_1 term is asymptotically normal. After reversing the order of the data, Theorem 1 also implies that the T_2 term is asymptotically normal. We then need to compute the asymptotic variance, which can be done by applying Lemma 8 with $z_{[i]} = \Delta y_{[i]} / (2f(x_{[i]})) + \Delta y_{[i+1]} / (2f(x_{[i+1]}))$, since it allows for dependent sequences such as z_i . We simply have to derive the new values of ρ^2 and τ^2 : $\rho^2 = \lim_{n \rightarrow \infty} \text{Var} \left[n^{-1/2} \sum_{i=1}^{n-1} (\Delta y_{[i]} / (2f(x_{[i]})) + \Delta y_{[i+1]} / (2f(x_{[i+1]}))) \right] = \lim_{n \rightarrow \infty} n^{-1} \text{Var} [\sum_{i=1}^n \Delta y_i / f(x_i)] = E [\text{Var} [y|x] / f^2(x)] = \sigma^2$ and $\tau^2 = \lim_{n \rightarrow \infty} n^{-1} \left(\sum_{i=1}^{n-1} \text{Var} (y_{[i]} / (2f(x_{[i]})) + y_{[i+1]} / (2f(x_{[i+1]}))) \right) = \lim_{n \rightarrow \infty} (2n)^{-1} \sum_{i=1}^n \text{Var} (y_{[i]} / f(x_{[i]})) = \text{Var} (\Delta y_i / f(x_i)) / 2 = E [\text{Var} [y|x] / f^2(x)] / 2 = \sigma^2 / 2$, implying that $n^{1/2} (\hat{\theta} - \theta) \xrightarrow{d} N(0, \rho^2 + \tau^2)$, where $\rho^2 + \tau^2 = 3\sigma^2 / 2$. ■

Proof. (of Theorem 3) Let $g(x) = E[y|x]$, $\Delta y_{[i]} = y_{[i]} - g(x_{[i]})$, $u_{[i]} = F(x_{[i]})$, and $d_{k_n, [i]} = u_{[i+k_n]} - u_{[i]}$. Then, $\tilde{\theta}_{k_n} - \theta = \sum_{i=1}^{n-k_n} y_{[i]} (x_{[i+k_n]} - x_{[i]}) / k_n - \int g(x) dx = N_1 - R_0 + R_1 + R_2 + R_3$, where

$$N_1 = n^{-1} \sum_{i=1}^n \frac{\Delta y_{[i]}}{f(x_{[i]})} \quad (40)$$

$$R_0 = n^{-1} \sum_{i=n-k_n+1}^n \frac{\Delta y_{[i]}}{f(x_{[i]})} \quad (41)$$

$$R_1 = \sum_{i=1}^{n-k_n} g(x_{[i]}) (x_{[i+k_n]} - x_{[i]}) / k_n - \int g(x) dx \quad (42)$$

$$R_2 = \sum_{i=1}^{n-k_n} \Delta y_{[i]} (x_{[i+k_n]} - x_{[i]}) / k_n - \sum_{i=1}^{n-k_n} \frac{\Delta y_{[i]}}{f(x_{[i]})} d_{k_n, [i]} / k_n \quad (43)$$

$$R_3 = n^{-1} \sum_{i=1}^{n-k_n} \frac{\Delta y_{[i]}}{f(x_{[i]})} (n d_{k_n, [i]} / k_n - 1). \quad (44)$$

First, we have $N_1 = n^{-1} \sum_{i=1}^n \Delta y_{[i]} / f(x_{[i]}) = n^{-1} \sum_{i=1}^n \Delta y_i / f(x_i)$, which is asymptotically normal by the i.i.d. assumption and the assumption that $E [\text{Var} [\Delta y_i | x_i] / f^2(x_i)] < \infty$. Next, we bound the remainder terms. $R_0 \equiv r_0(x_{[n-k_n+1]})$ where

$$r_0(\xi) = n^{-1} \sum_{i=1}^n \frac{\Delta y_{[i]}}{f(x_{[i]})} 1(x_{[i]} \geq \xi) = n^{-1} \sum_{i=1}^n \frac{\Delta y_i}{f(x_i)} 1(x_i \geq \xi) \quad (45)$$

By Tschebychev's inequality, for any $\varepsilon > 0$, $P[|n^{1/2} r_0(\xi)| \geq \varepsilon |\xi|] \leq \varepsilon^{-2} E[1(x_i \geq \xi) (\Delta y_i / f(x_i))^2]$. As $\text{Var}(y_i | x_i) / f^2(x_i)$ is positive and $E[\text{Var}(y_i | x_i) / f^2(x_i)]$ is finite, $E[1(x_i \geq \xi) (\Delta y_i / f(x_i))^2] \rightarrow 0$ as $\xi \rightarrow$

$\bar{x} \equiv \sup_{x \in \mathcal{F}} x$. Let ξ_n be a deterministic sequence such that $\xi_n \rightarrow \bar{x}$ and such that $x_{[n-k_n+1]} \geq \xi_n$ with probability approaching one, which is possible since $\text{plim } x_{[n-k_n+1]} = \text{plim } x_{[n]} = \bar{x}$. Then, w.p.a. 1, $P \left[|n^{1/2} r_0(x_{[n-k_n+1]})| \geq \varepsilon \right] \leq P \left[|n^{1/2} r_0(\xi_n)| \geq \varepsilon \right] \rightarrow 0$, implying that $P \left[|n^{1/2} r_0(x_{[n-k_n+1]})| \geq \varepsilon \right] \rightarrow 0$ and that $r_0(x_{[n-k_n+1]}) = o_p(n^{-1/2})$.

Some of the changes of variables in the summations below introduce ‘‘boundary’’ terms, denoted by B , which can be shown to be $o_p(n^{-1/2})$ (The proof is available upon request). Next, we consider, $|R_1| = \left| \sum_{i=1}^{n-k_n} g(x_{[i]}) (x_{[i+k_n]} - x_{[i]}) / k_n - \sum_{i=1}^{n-1} \int_{x_{[i]}}^{x_{[i+1]}} g(x) dx \right| = \left| \sum_{i=1}^{n-k_n} g(x_{[i]}) (x_{[i+k_n]} - x_{[i]}) / k_n - k_n^{-1} \sum_{i=1}^{n-k_n} \int_{x_{[i]}}^{x_{[i+k_n]}} g(x) dx + B \right| = \left| \sum_{i=1}^{n-k_n} g(x_{[i]}) (x_{[i+k_n]} - x_{[i]}) / k_n - k_n^{-1} \sum_{i=1}^{n-k_n} (x_{[i+k_n]} - x_{[i]}) g(\xi_{[i]}) + B \right|$ for $\xi_{[i]} \in [x_{[i]}, x_{[i+k_n]}]$.

Next, $|R_1| \leq H_g |k_n^{-1} \sum_{i=1}^{n-k_n} (x_{[i+k_n]} - x_{[i]})|^{1+h_g} + |B| \leq H_g \left| \sum_{i=1}^{n-k_n} (k_n^{-1} \sum_{j=1}^{k_n} (x_{[i+j]} - x_{[i+j-1]})) \right|^{1+h_g} + |B| \leq H_g \left| \sum_{i=1}^{n-k_n} k_n^{-1} \sum_{j=1}^{k_n} (x_{[i+j]} - x_{[i+j-1]}) \right|^{1+h_g} + |B| \leq H_g \left| \sum_{i=1}^{n-k_n} (x_{[i+1]} - x_{[i]}) \right|^{1+h_g} + |B| = O_p(n^{-h_g}) = o_p(n^{-1/2})$

The R_2 remainder term can be similarly shown to be $o_p(n^{-1/2})$. Finally, $R_3 = n^{-1} \sum_{i=1}^{n-k_n} f^{-1}(x_{[i]}) \Delta y_{[i]} (nd_{k, [i]} / k_n - 1) = n^{-1} \sum_{i=1}^{n-k_n} f^{-1}(x_{[i]}) \Delta y_{[i]} \left(\left(k_n^{-1} \sum_{j=1}^{k_n} nd_{1, [i+j-1]} \right) - 1 \right)$. Proceeding as in Lemma 9 and introducing $g_{[i]}$ and $G_{[n]}$ defined therein, R_3 has the same distribution as

$$\begin{aligned} & n^{-1} \sum_{i=1}^{n-k_n} \frac{\Delta y_{[i]}}{f(x_{[i]})} \left(\left(\frac{n}{G_{[n]}} \right) \left(k_n^{-1} \sum_{j=1}^{k_n} g_{[i+j-1]} \right) - 1 \right) \\ &= \left(\frac{n}{G_{[n]}} \right) n^{-1} \sum_{i=1}^{n-k_n} \frac{\Delta y_{[i]}}{f(x_{[i]})} \left(\left(k_n^{-1} \sum_{j=1}^{k_n} g_{[i+j-1]} \right) - 1 \right) + \left(1 - \left(\frac{n}{G_{[n]}} \right) \right) n^{-1} \sum_{i=1}^{n-k_n} \frac{\Delta y_{[i]}}{f(x_{[i]})} \\ &= (1 + o_p(1)) R_{31} + o_p(1) O_p(n^{-1/2}) \end{aligned} \quad (46)$$

where $R_{31} = n^{-1} \sum_{i=1}^{n-k_n} f^{-1}(x_{[i]}) \Delta y_{[i]} \left(k_n^{-1} \sum_{j=1}^{k_n} (g_{[i+j-1]} - 1) \right)$. Employing well-known techniques used for the study of U -statistics, it can be shown that $E[R_{31}^2] = O(k_n^{-1} n^{-1})$, implying that $R_3 = O_p(k_n^{-1/2} n^{-1/2}) = o_p(n^{-1/2})$. ■

Proof. (of Theorem 4) We first observe that $\hat{\sigma}^2 = \frac{n}{4} \sum_{i=1}^{n-1} (\Delta y_{[i+1]} - \Delta y_{[i]} + r_{[i]})^2 (x_{[i+1]} - x_{[i]})^2$ where $|r_{[i]}| = |(y_{[i+1]} - y_{[i]}) - (\Delta y_{[i+1]} - \Delta y_{[i]})| = |(g(x_{[i+1]}) - g(x_{[i]}))| \leq H_g |x_{[i+1]} - x_{[i]}|^{h_g}$, by the Hölder property of $g(x)$. The remainder of the proof (available upon request) is a straightforward but tedious extension of the techniques used in Theorem 1. ■

Proof. (of Corollary 5) This result can be shown along the same lines as Theorems 1, 4 and Corollary 2, with the Cramer-Wold device, letting $y_i = \sum_{j=1}^k \alpha_j v_{ij}$, where $(\alpha_1, \dots, \alpha_k)$ is a vector of arbitrary constants, and noting that $E[\text{Covar}[v_{ij}, v_{i'j'}|x_i]/f^2(x_i)] \leq (E[\text{Var}[v_{ij}|x_i]/f(x_i)] E[\text{Var}[v_{i'j'}|x_i]/f(x_i)])^{1/2}$ by Cauchy-Schwartz, thus implying that all elements of the matrix V are bounded. ■

Lemma 12 *Let x be a random scalar and w be a random vector taking value in \mathbb{R}^{N_w} related through $x = X(w, \gamma)$ for some function $X : \mathbb{R}^{N_w} \times \mathbb{R}^{N_\gamma} \mapsto \mathbb{R}$ and for some parameter vector γ . If $\partial X(w, \gamma)/\partial\gamma$ exists and is such that $E[|\partial X(w, \gamma)/\partial\gamma|] < \infty$ then*

$$\frac{\partial f_{x|\gamma}(x|\gamma)}{\partial\gamma} = -\frac{\partial}{\partial x} \left(f_{x|\gamma}(x|\gamma) E \left[\frac{\partial X(w, \gamma)}{\partial\gamma} | x \right] \right). \quad (47)$$

Proof. Let us introduce the sequence

$$S_m(\xi) = \begin{cases} 1 & \text{if } \xi < -m^{-1} \\ (1 - m\xi)/2 & \text{if } \xi \in [-m^{-1}, m^{-1}] \\ 0 & \text{if } \xi > m^{-1} \end{cases}, \quad (48)$$

which converges pointwise to the indicator function $1(\xi \leq 0)$, except at $\xi = 0$. By the Dominated Convergence Theorem, we have $F_{x|\gamma}(x|\gamma) = \int 1[X(w, \gamma) \leq x] dF_w(w) = \int \lim_{m \rightarrow \infty} S_m(X(w, \gamma) - x) dF_w(w) = \lim_{m \rightarrow \infty} \int S_m(X(w, \gamma) - x) dF_w(w)$ since the integrand is dominated, for all m , by the absolutely integrable measure $dF_w(w)$. Differentiating yields $\partial F_{x|\gamma}(x|\gamma)/\partial\gamma = \lim_{m \rightarrow \infty} \int S'_m(X(w, \gamma) - x) \partial X(w, \gamma)/\partial\gamma dF_w(w)$, where $S'_m(\xi) = -2m1(|x| \leq m^{-1})$ and where the operator $\partial/\partial\gamma$ commutes with the limit and integral because $S'_m(\xi)$ is absolutely integrable by construction and so is $(\partial X(w, \gamma)/\partial\gamma) dF_w(w)$, by assumption. Noting that $S'_m(\xi)$ forms a sequence of functions converging to minus the Dirac delta distribution $-\delta(\xi)$, we have

$$\begin{aligned} \frac{\partial}{\partial\gamma} F_{x|\gamma}(x|\gamma) &= - \int \delta(X(w, \gamma) - x) \frac{\partial X(w, \gamma)}{\partial\gamma} dF_w(w) \\ &= -f_{x|\gamma}(x|\gamma) \frac{\int \delta(X(w, \gamma) - x) \frac{\partial X(w, \gamma)}{\partial\gamma} dF_w(w)}{f_{x|\gamma}(x|\gamma)} \\ &= -f_{x|\gamma}(x|\gamma) E \left[\frac{\partial X(w, \gamma)}{\partial\gamma} | x \right], \end{aligned} \quad (49)$$

where the last equality uses the definition of $E[\partial X(w, \gamma)/\partial\gamma|x]$. Finally, differentiating Equation (49) with respect to x yields Equation (47). ■

Proof. (of Theorem 6) We can decompose the estimation error as $n^{1/2} \left(\tilde{\theta}_{k_n}(\hat{\gamma}) - \theta(\gamma^*) \right) = N_1 + R_1 + N_2 + R_2$ where

$$N_1 = n^{1/2} \left(\tilde{\theta}_{k_n}(\gamma^*) - \theta(\gamma^*) \right) \quad (50)$$

$$N_2 = n^{1/2} \frac{\partial \theta(\gamma^*)}{\partial \gamma'} (\hat{\gamma} - \gamma^*) \quad (51)$$

$$R_1 = n^{1/2} \left(\tilde{\theta}_{k_n}(\hat{\gamma}) - \theta(\hat{\gamma}) \right) - n^{1/2} \left(\tilde{\theta}_{k_n}(\gamma^*) - \theta(\gamma^*) \right) \quad (52)$$

$$R_2 = n^{1/2} \left(\frac{\partial \theta(\hat{\gamma})}{\partial \gamma'} - \frac{\partial \theta(\gamma^*)}{\partial \gamma'} \right) (\hat{\gamma} - \gamma^*) \quad (53)$$

where $\hat{\gamma}$ is a mean value located along the segment joining γ^* and $\hat{\gamma}$. We will evaluate, in turn, (i) R_1 , (ii) R_2 , (iii) N_1 (iv) N_2 and (v) the asymptotic variance of $N_1 + N_2$.

(i) Theorem 3 shows that $n^{1/2} \left(\tilde{\theta}_{k_n}(\gamma) - \theta(\gamma) \right)$ for a given γ as asymptotically equivalent to $n^{-1/2} \sum_{i=1}^n (y_i - E[y_i|x_i]) / f_{x|\gamma}(x_i|\gamma)$. In order to handle the fact that all the quantities are function of γ , Theorem 3 needs to be adapted for the order of the remainder terms to hold jointly for γ^* and $\hat{\gamma}$. This is achieved by replacing Assumptions 2 through 5 by corresponding Assumptions 6 through 9 that hold uniformly for $\gamma \in \Gamma$. It follows that R_1 is such that

$$\begin{aligned} R_1 &= n^{-1/2} \sum_{i=1}^n \frac{y_i - E[y_i|x_i = X(w_i, \hat{\gamma})]}{f_{x|\gamma}(x_i|\hat{\gamma})} - n^{-1/2} \sum_{i=1}^n \frac{y_i - E[y_i|x_i = X(w_i, \gamma^*)]}{f_{x|\gamma}(x_i|\gamma^*)} + o_p(1) \\ &= R_{11} + R_{12} + o_p(1) \end{aligned} \quad (54)$$

where

$$R_{11} = n^{-1/2} \sum_{i=1}^n \frac{f_{x|\gamma}(x_i|\gamma^*) - f_{x|\gamma}(x_i|\hat{\gamma})}{f_{x|\gamma}(x_i|\hat{\gamma}) f_{x|\gamma}(x_i|\gamma^*)} \Delta y_i \quad (55)$$

$$R_{12} = n^{-1/2} \sum_{i=1}^n \left(\frac{E[y_i|x_i = X(w_i, \hat{\gamma})] - E[y_i|x_i = X(w_i, \gamma^*)]}{f_{x|\gamma}(x_i|\hat{\gamma})} \right) \quad (56)$$

and where $\Delta y_i = y_i - E[y_i|x_i = X(w_i, \gamma^*)]$. Tschebychev's inequality lets us write, for any $\varepsilon > 0$,

$$P[|R_{11}| \geq \varepsilon | \hat{\gamma}] \leq \varepsilon^{-2} E \left[\left(n^{-1/2} \sum_{i=1}^n \frac{(f_{x|\gamma}(x_i|\gamma^*) - f_{x|\gamma}(x_i|\hat{\gamma}))}{f_{x|\gamma}(x_i|\hat{\gamma}) f_{x|\gamma}(x_i|\gamma^*)} \Delta y_i \right)^2 | \hat{\gamma} \right] \quad (57)$$

$$= \varepsilon^{-2} (n^{-1}n) E \left[\frac{(f_{x|\gamma}(x_i|\gamma^*) - f_{x|\gamma}(x_i|\hat{\gamma}))^2}{f_{x|\gamma}^2(x_i|\hat{\gamma}) f_{x|\gamma}^2(x_i|\gamma^*)} (\Delta y_i)^2 | \hat{\gamma} \right] \quad (58)$$

where $E \left[f_{x|\gamma}^{-2}(x_i|\hat{\gamma}) f_{x|\gamma}^{-2}(x_i|\gamma^*) (f_{x|\gamma}(x_i|\gamma^*) - f_{x|\gamma}(x_i|\hat{\gamma}))^2 y_i^2 | \hat{\gamma} \right] \xrightarrow{p} 0$ as¹¹ $\hat{\gamma} \xrightarrow{p} \gamma^*$, by the continuity of

¹¹This expectation converges *in probability* because it is a function of $\hat{\gamma}$, which is random.

$f_{x|\gamma}(x_i|\gamma)$ in γ (from Lemma 12), the fact that $f_{x|\gamma}(x_i|\gamma)$ is bounded away from zero for $\gamma \in \Gamma$ and that

$$E \left[\frac{(f_{x|\gamma}(x_i|\gamma^*) - f_{x|\gamma}(x_i|\hat{\gamma}))^2}{f_{x|\gamma}^2(x_i|\hat{\gamma}) f_{x|\gamma}^2(x_i|\gamma^*)} (\Delta y_i)^2 | \hat{\gamma} \right] \leq \frac{4 \sup_{\gamma \in \Gamma} \sup_{x \in \mathcal{F}_\gamma} f_{x|\gamma}^2(x|\gamma)}{\inf_{\gamma \in \Gamma} \inf_{x \in \mathcal{F}_\gamma} f_{x|\gamma}^2(x|\gamma)} E \left[\frac{(\Delta y_i)^2}{f_{x|\gamma}^2(x|\gamma^*)} \right] < \infty$$

(since $\sup_{\gamma \in \Gamma} \sup_{x \in \mathcal{F}_\gamma} f_{x|\gamma}^2(x|\gamma)$ must be finite by the Hölder continuity of $f_{x|\gamma}(x|\gamma)$). Equation (58) then implies that $R_{11} \xrightarrow{p} 0$.

Similarly, $P[|R_{12}| \geq \varepsilon | \hat{\gamma}] \leq \varepsilon^{-2} E[(E[y_i|x_i = X(w, \hat{\gamma})] - E[y_i|x_i = X(w, \gamma^*)])^2 / f_{x|\gamma}^2(x_i|\hat{\gamma}) | \hat{\gamma}]$. Since $E[y_i|x_i = X(w, \gamma)] = \int f_{x|\gamma}^{-1}(X(w, \gamma) | \gamma) E[y|w] \delta(x_i - X(w, \gamma)) dw$, continuity of $f_{x|\gamma}(x|\gamma)$ and $X(w, \gamma)$ in γ and $\inf_{\gamma \in \Gamma} \inf_x f_{x|\gamma}(x|\gamma) > 0$ implies that $E[y_i|x_i = X(w, \gamma)]$ is continuous in γ . Then, as for R_{11} , we have $R_{12} \xrightarrow{p} 0$.

(ii) The remainder R_2 can then be bounded using the continuity of $\partial\theta(\gamma)/\partial\gamma'$ (from Assumption 11), the assumed root n consistency of $\hat{\gamma}$ and the fact that $\hat{\gamma} \xrightarrow{p} \gamma^*$: $R_2 = (\partial\theta(\gamma^*)/\partial\gamma' - \partial\theta(\hat{\gamma})/\partial\gamma') n^{1/2}(\hat{\gamma} - \gamma) = o_p(1) n^{1/2} O_p(n^{-1/2}) = o_p(1)$.

(iii) From the proofs of Theorem 3, the N_1 term can be written as $N_1 = n^{-1/2} \sum_{i=1}^n f_{x|\gamma}^{-1}(x|\gamma^*) \Delta y_i + o_p(1)$, where the sum is asymptotically normal.

(iv) By the assumption of the existence of an influence representation for $\hat{\gamma}$, the N_2 term can be written as $N_2 = B' n^{-1/2} \sum_{i=1}^n \Psi_i + o_p(1)$, where $B = \partial\theta(\gamma^*)/\partial\gamma$ and where the sum is asymptotically normal.

(v) The estimation error can then be written as $n^{1/2} (\tilde{\theta}_{k_n}(\hat{\gamma}) - \theta(\gamma^*)) = n^{-1/2} \sum_{i=1}^n (f_{x|\gamma}^{-1}(x|\gamma^*) \Delta y_i + B' \Psi_i) + o_p(1)$ and straightforward calculations provide the expression of the asymptotic variance V of this sum. An explicit expression for B can be given. First note that

$$\theta(\gamma) = E \left[\frac{y}{f_{x|\gamma}(X(w, \gamma) | \gamma)} \right] = \int \frac{E[y|w]}{f_{x|\gamma}(X(w, \gamma) | \gamma)} f_w(w) dw. \quad (59)$$

This shows that the only dependence of this expression on γ comes from the denominator $f_{x|\gamma}(X(w, \gamma) | \gamma)$.

We can then write (the required expectations and derivatives commute by Assumption 11):

$$\frac{\partial\theta(\gamma)}{\partial\gamma} = -E \left[\frac{y}{(f_{x|\gamma}(X(w, \gamma) | \gamma))^2} \frac{\partial f_{x|\gamma}(X(w, \gamma) | \gamma)}{\partial\gamma} \right]. \quad (60)$$

The quantity $f_{x|\gamma}(X(w, \gamma) | \gamma)$ depends on γ through two different paths: (i) the function $f_{x|\gamma}(\cdot | \gamma)$ depends on γ and (ii) the point of evaluation $X(w, \gamma)$ depends on γ . We then have

$$\frac{\partial\theta(\gamma)}{\partial\gamma} = -E \left[\frac{y}{(f_{x|\gamma}(X(w, \gamma) | \gamma))^2} \left(\left[\frac{\partial f_{x|\gamma}(x|\gamma)}{\partial\gamma} \right]_{x=X(w, \gamma)} + \left[\frac{\partial f_{x|\gamma}(x|\gamma)}{\partial x} \right]_{x=X(w, \gamma)} \frac{\partial X(w, \gamma)}{\partial\gamma} \right) \right]. \quad (61)$$

Using Lemma 12 to evaluate $\partial f_{x|\gamma}(x|\gamma)/\partial\gamma$, we obtain

$$\begin{aligned} \frac{\partial\theta(\gamma)}{\partial\gamma} &= E \left[\frac{y}{f_{x|\gamma}(x|\gamma)} \frac{\partial}{\partial x} E \left[\frac{\partial X(w, \gamma)}{\partial\gamma} | x \right] \right] + \\ &\quad - E \left[\frac{y}{f_{x|\gamma}(x|\gamma)} \frac{\partial (\ln(f_{x|\gamma}(x|\gamma)))}{\partial x} \left(\frac{\partial X(w, \gamma)}{\partial\gamma} - E \left[\frac{\partial X(w, \gamma)}{\partial\gamma} | x \right] \right) \right] \end{aligned} \quad (62)$$

In the special case mentioned in the second part of the theorem, it can be verified that $E[\partial X(w, \gamma^*)/\partial\gamma|x, y] = E[\partial X(w, \gamma^*)/\partial\gamma]$ implies that $\partial\theta(\gamma)/\partial\gamma|_{\gamma=\gamma^*} = 0$. ■

Proof. (of Corollary 7) Since the estimator $\tilde{\theta} = \sum_{i=1}^{n-1} v_{[i]}(x_{[i+1]} - x_{[i]})$, cannot be written as a differentiable functional of the joint cdf of y and x , an “influence function” cannot be defined. However, this estimator can be written as a differentiable functional of the joint density of y , x and an auxiliary random variable g_i , which is an iid sequence independent from y_i and x_i and drawn from a $e^{-g_i}1(g_i \geq 0)$ distribution. We can then consider a “pseudo influence function” $\tilde{\Psi}(y_i, x_i) = (y_i - E[y_i|x_i])g_i/f(x_i)$. The “pseudo influence function” for the more efficient estimator $\hat{\theta} = \sum_{i=1}^{n-1} (v_{[i]} + v_{[i+1]})(x_{[i+1]} - x_{[i]})/2$ is $\hat{\Psi}(y_i, x_i) = (y_i - E[y_i|x_i])(g_i + h_i)/(2f(x_i))$ where h_i is independent from g_i and has the same distribution.

For the binary choice estimator $\hat{\beta}$, the appropriate pseudo influence function is

$$n^{1/2}(\hat{\beta} - \beta) = n^{-1/2} \Delta \sum_{i=1}^n \left(\frac{(y_i - E[y_i|x_i])(g_i + h_i)}{f(x_i)} - (r_i z'_i - \overline{r z'}) \beta \right) + o_p(1) \quad (63)$$

(with g_i and h_i defined above) and so the asymptotic variance is $\Delta E[(U_i - W_i)(U_i - W_i)'] \Delta'$ where

$$U_i = \frac{(y_i - E[y_i|x_i])(g_i + h_i)}{f(x_i)} \quad (64)$$

$$W_i = (r_i z'_i - \overline{r z'}) \beta. \quad (65)$$

Note that $E[W_i W_i']$ is just a standard variance expression. However,

$$\begin{aligned} E[U_i U_i'] &= E \left[((y_i - E[y_i|x_i])(y_i - E[y_i|x_i])' / f^2(x_i)) (g_i + h_i)^2 / 4 \right] \\ &= E \left[((y_i - E[y_i|x_i])(y_i - E[y_i|x_i])' / f^2(x_i)) \right] E \left[(g_i + h_i)^2 / 4 \right] \\ &= E \left[((y_i - E[y_i|x_i])(y_i - E[y_i|x_i])' / f^2(x_i)) \right] (E[g_i^2] + 2E[g_i]E[h_i] + [h_i^2]) / 4 \\ &= E \left[((y_i - E[y_i|x_i])(y_i - E[y_i|x_i])' / f^2(x_i)) \right] \frac{6}{4} \\ &= (3/2) E[\text{Var}(y_i|x_i) / f^2(x_i)] \end{aligned} \quad (66)$$

where $E [\text{Var} (y_i|x_i) / f^2 (x_i)]$ can be estimated by $\frac{n}{4} \sum_{i=1}^{n-1} (y_{[i+1]} - y_{[i]}) (y_{[i+1]} - y_{[i]})' (x_{[i+1]} - x_{[i]})^2$. Finally,

$$E [U_i W_i'] = E [(y_i - E [y_i|x_i]) W_i' / f (x_i)] E [(g_i + h_i) / 2] = E [(y_i - E [y_i|x_i]) W_i' / f (x_i)] \quad (67)$$

and $E [U_i W_i']$ can be estimated by $\sum_{i=1}^n (y_{[i]} - y_{[i+1]}) (x_{[i+1]} - x_{[i]}) n W_{[i]}$ since,

$$\begin{aligned} & E [g_{[i]} W_{[i]} (y_{[i]} - y_{[i+1]}) / f (x_{[i]})] \\ &= E [g_{[i]} W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}] - (y_{[i+1]} - E [y_{[i]}|x_{[i]}])) / f (x_{[i]})] \\ &= E [g_{[i]} W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}] - (y_{[i+1]} - E [y_{[i+1]}|x_{[i+1]}])) / f (x_{[i]})] + o_p (n^{-1/2}) \\ &= E [g_{[i]} W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}]) / f (x_{[i]})] + \\ &\quad - E [E [(y_{[i+1]} - E [y_{[i+1]}|x_{[i+1]}]) |x_{[i+1]}] g_{[i]} W_{[i]} / f (x_{[i]})] + o_p (n^{-1/2}) \\ &= E [g_{[i]} W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}]) / f (x_{[i]})] + 0 + o_p (n^{-1/2}) \\ &= E [g_{[i]}] E [W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}]) / f (x_{[i]})] + o_p (n^{-1/2}) \\ &= E [W_{[i]} (y_{[i]} - E [y_{[i]}|x_{[i]}]) / f (x_{[i]})] + o_p (n^{-1/2}). \end{aligned} \quad (68)$$

The expression for $\widehat{\Omega}$ given in the theorem thus follows. ■

Acknowledgements

This research was supported in part by the National Science Foundation through grants SES-9905010 (Lewbel) and SES-0214068 (Schnnach). The authors would like to thank Peter Robinson and two anonymous referees for helpful comments, Eric Maurin for providing data, Zhihong Chen for research assistance, and Daniel McFadden for suggesting this class of estimators.

References

- Abadie, A. and G. Imbens, 2002, Simple and Bias-Corrected Matching Estimators for Average Treatment Effects, *Econometrica*, forthcoming.
- Barbe, P., 1994, Joint approximation of processes based on spacings and order statistics, *Stochastic Processes and their Applications* 53, 339–349.
- Bloch, D. A. and J. L. Gastwirth, 1968, On a Simple Estimate of the Reciprocal of the Density Function, *Annals of Mathematical Statistics* 39, 1083–1085.
- Cheng, M.-Y., J. Fan, and J. S. Marron, 1997, On automatic boundary corrections, *Annals of Statistics* 25, 1691–1708.
- Cogneau, D. and E. Maurin, 2002, Parental Income and School Attendance in a Low-Income Country: a semi-parametric analysis, Unpublished Manuscript.
- Hall, P. and J. L. Horowitz, 1990, Bandwidth Selection in Semiparametric Estimation of Censored Linear Regression Models, *Econometric Theory* 6, 123–150.
- Hall, P. and A. Yatchew, 2005, Unified approach to testing functional hypotheses in semiparametric contexts, *Journal of Econometrics*, forthcoming.
- Härdle, W. and T. M. Stoker, 1989, Investigating Smooth Multiple Regression by the Method of Average Derivatives, *Journal of the American Statistical Association*, 84, 986–995.
- Härdle, W., J. Hart, J. S. Marron and A. B. Tsybakov, 1992, Bandwidth Choice for Average Derivative Estimation, *Journal of the American Statistical Association*, 87, 218–226.
- Hausman, J. A. and W. Newey, 1995, Nonparametric Estimation of Exact Consumers Surplus and Dead-weight Loss, *Econometrica* 63, 1445–1476.
- Hirano, K., G. Imbens, and G. Ridder, 2003, Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score, *Econometrica* 71, 1161–1189.

- Hong, Y. and H. White, 2005, Asymptotic Distribution Theory for Nonparametric Entropy Measures of Serial Dependence, *Econometrica* 73, 837–901.
- Klein, R. and R. H. Spady, 1993, An efficient Semiparametric Estimator for Binary Response Models, *Econometrica* 61, 387–421.
- Lewbel, A., 1997, Semiparametric Estimation of Location and Other Discrete Choice Moments, *Econometric Theory* 13, 32–51.
- Lewbel, A., 1998, Semiparametric Latent Variable Model Estimation With Endogenous or Mismeasured Regressors, *Econometrica* 66, 105–121.
- Lewbel, A., 2000, Semiparametric Qualitative Response Model Estimation With Unknown Heteroskedasticity or Instrumental Variables, *Journal of Econometrics* 97, 145–177.
- Lewbel, A., 2002, Endogeneous Selection or Treatment Model Estimation, Unpublished Manuscript.
- Lewbel, A., O. Linton, and D. L. McFadden, 2001, Estimating Features of a Distribution From Binomial Data, unpublished manuscript.
- Magnac, T. and E. Maurin, 2003, Identification and Information in Monotone Binary Models, unpublished manuscript.
- McFadden, D. L., 1999, Computing Willingness-to-Pay in Random Utility Models, in: J. Moore, R. Riezman, and J. Melvin, eds., *Trade Theory, and Econometrics: Essays in Honour of John S. Chipman* (Routledge, New York) 253–274.
- Mack, Y. P. and H. G. Müller, 1988, Convolution Type Estimators for Nonparametric Regression, *Statistics and Probability Letters* 7, 229–239.
- Newey, W. K., 1994, The Asymptotic Variance of Semiparametric Estimators, *Econometrica* 62, 1349–1382.
- Newey, W. K., 1997, Convergence rates and asymptotic normality for series estimators, *Journal of econometrics* 79, 147–168.

- Newey, W. K. and P. A. Ruud, 1994, Density Weighted Linear Least Squares, University of California at Berkeley working paper.
- Priestley, M. B. and M. T. Chao, 1972, Nonparametric Function Fitting, *Journal of the Royal Statistical Society, Series B* 34, 385–392.
- Robinson, P. M., 1987, Asymptotically Efficient Estimation in the Presence of Heteroskedasticity of Unknown Form, *Econometrica* 55, 875–891.
- Staiger D., and J. H. Stock, 1997, Instrumental Variables Regression with Weak Instruments, *Econometrica* 65, 557–586.
- Wand, M. P. and M. C. Jones, 1995, *Kernel Smoothing* (CRC Press, Boston).
- Weiss, L., 1958, Limiting distributions of homogenous functions of sample spacings, *Annals of Mathematical Statistics* 29, 310–312.
- Yatchew, A., 1997, An Elementary Estimator of the Partial Linear Model, *Economics Letters* 57, 135–43.

TABLE 1: MONTE CARLO RESULTS

	mean	stddev	lower	median	upper	rmse	mae	mdae
true density $\hat{\theta}_0$.9959	.2834	.7967	.9820	1.177	.2835	.2264	.1919
ordered data 1 $\hat{\theta}_1$	1.000	.2663	.8209	1.001	1.175	.2663	.2112	.1770
ordered data 2 $\hat{\theta}_2$	1.003	.2470	.8374	1.001	1.164	.2470	.1954	.1634
ordered data 3 $\hat{\theta}_3$	1.006	.2429	.8441	1.002	1.166	.2430	.1921	.1608
Silverman b kernel $\hat{\theta}_4$.9925	.2338	.8334	.9891	1.151	.2339	.1866	.1588
$b/2$ kernel $\hat{\theta}_5$.9679	.2262	.8170	.9647	1.121	.2285	.1820	.1528
$2b$ kernel $\hat{\theta}_6$	1.029	.2546	.8526	1.025	1.199	.2562	.2048	.1744

TABLE 2: SCHOOL ATTENDANCE ESTIMATES

	Kernel OLS	Ordered OLS	Kernel 2SLS	Ordered 2SLS
Constant	-2.758 (.474)	-2.251 (.757)	-12.394 (2.195)	-8.394 (3.324)
Boy	.030 (.041)	-.007 (.067)	.033 (.047)	-.005 (.070)
Log Income	.173 (.035)	.144 (.056)	.882 (.161)	.597 (.244)
Mother's Education	.194 (.026)	.036 (.039)	-.089 (.061)	-.068 (.094)

Note: Standard errors are in parentheses.