# ESTIMATING FEATURES OF A DISTRIBUTION FROM BINOMIAL DATA

*Arthur Lewbel*
*Oliver Linton*
*Daniel McFadden*

# ESTIMATING FEATURES OF A DISTRIBUTION FROM BINOMIAL DATA[¤]

Arthur Lewbel[y]

Boston College

Oliver Linton[z]

London School of Economics

Daniel McFadden

University of California, Berkeley[x]

December 4, 2001

## Abstract

A statistical problem that arises in several ...elds is that of estimating the features of an unknown distribution, which may be conditioned on covariates, using a sample of binomial observations on whether draws from this distribution exceed threshold levels set by experimental design. One application is destructive duration analysis, where the process is censored at an observation test time. Another is referendum contingent valuation in resource economics, where one is interested in features of the distribution of values placed by consumers on a public good such as endangered species. Sample consumers are asked whether they would vote for a referendum that would provide the good at a cost speci...ed by experimental design. This paper provides estimators for moments and quantiles of the unknown distribution in this problem.

# 1   Introduction

A statistical problem that arises in several ...elds is that of estimating the features of an unknown distribution, which may be conditioned on covariates, using a sample of binomial observations on whether draws from this distribution exceed threshold levels set by experimental design. Three applications illustrate the problem:

Bioassay - Find the distribution of survival times until the onset of an abnormality in laboratory animals exposed to an environmental hazard. The animals are sacri...ced at times determined by experimental design, and tested for the abnormality. An observation consists of a vector of covariates, a test time, and an indicator for the test result.

Destructive Testing - Find the distribution of speeds at which air bags fail to protect passengers in automobile crashes. At speeds selected by experimental design, drive cars into a barrier and determine whether a dummy occupant is injured. An observation consists of covariates, a test speed, and an indicator for injury.

Survey research with Shadow E¤ects - Find the distribution of a household economic variable such as wealth. Subjects are asked if their economic variable exceeds a test value chosen by design. An observation consists of covariates, a test value, an indicator for the response. Follow up queries are shadowed by the framing e¤ect of the ...rst bid. This shadowing e¤ect is common in unfolding bracket survey questions on economic variables, and on stated willingness to pay (WTP) for economic goods.[1]

Given a set of covariates, when the experimental design is randomized with a strictly positive test value density and mild regularity conditions, we propose consistent estimators for conditional (on covariates) moments of the unknown distribution. We also provide root n consistent estimators for the case where the unknown distribution depends on covariates through a single index location shift. In addition, we provide estimators of conditional quantiles of the unknown distribution.

# 2   Model Speci...cation

The goal is estimation of conditional moments or quantiles of a latent, unobserved random scalar $W$, conditioned on a vector of observed covariates $X$. The conditional cumulative distribution function

---

[1]McFadden (1994) provides references and experminatal evidence that responses to follow up test values can be biased. There are additional issues of the impact of framing of questions on survey responses, particularly anchoring to test values, including the initial test value; see Green et al. (1998) and Hurd et al. (1998). The data generation process may then be a convolution of the target distribution and a distribution of psychometric errors. This paper will ignore these issues and treat the data generation process as if it is the target distribution. The di¢cult problem of deconvoluting a target distribution in the presence of psychometric errors is left for future research.

of $W$, denoted $G(w \mid x)$, is unknown but assumed to be smooth.

A test value $v$ is set by a randomized experimental design or natural experiment. The value $v$ is a realization of a random variable $V$, drawn from either a known or unknown conditional density $h(v \mid x)$ (we consider both cases). It is assumed that $W$ is conditionally independent of $V$, conditioning on $X$ (consistent with experimental design).

De…ne $Y$ to equal one in the event that $W$ exceeds $V$, and zero otherwise, so $Y = I(W > V)$ where $I(\phi)$ is the indicator function. The observed data consist of a random sample of realizations of covariates $X$, test values $V$, and outcomes $Y$. The framework is similar to random censored regressions (with censoring point $v$), except that for random censoring we would observe $w$ for observations having $w > v$, whereas in the present context we only observe $y = I(w > v)$.

Given a function $r(w; x)$, the goal is estimation of the conditional moment $\mu_r(x) = E[r(W; X) \mid X = x]$ for any chosen $x$ in the support of $X$. Of particular interest are the moments based on $r(W; X) = W^k$ for integers $k$: In addition to moments we may also be interested in quantiles. Let $w_q(x)$ denote the $q$'th quantile of $W$ given $x$:

If the conditional distribution of $W$ given $X = x$ is …nitely parameterized, then those parameters can generally be e¢ciently estimated by maximum likelihood (corresponding to ordinary binary choice model estimation, e.g., logit or probit models), thereby yielding e¢cient estimates for conditional moments $\mu_r(x)$ and quantiles $w_q(x)$ de…ned in terms of those parameters.

Assuming that the conditional distribution of $W$ given $X$ is not …nitely parameterized, we propose semiparametric and nonparametric estimators for these moments and quantiles. The semiparametric estimators assume that the conditional mean of $W$ is …nitely parameterized. The nonparametric estimators only require smoothness assumptions, but su¤er from the usual curse of dimensionality. We provide limit normal distributions for these estimators. The semiparametric estimators all converge at the rate that would be obtained if draws $w$ were observed.

In the example of willingness-to-pay models, $W_i$ would be an individual $i$'s unknown willingness to pay for a resource, $V_i$ would be a bid that was posed to the individual, and $X_i$ would be observable characteristics of the individual (such as age, income level, geographic location, and political party a¢liation). Objects of interest might include the average willingness-to-pay for the resource among individuals in certain locations and income levels, or for voting models, the median willingness-to-pay among subsets of likely voters.

The next section provides results that will form the basis for the proposed estimators. Later sections provide limiting distributions.

# 3  Identi...cation

Make the following assumptions.

**Assumption A.1.** The covariate vector $X$ has compact support $X \subseteq R^d$. The latent scalar $W$ has an unknown, twice continuously di¤erentiable conditional c.d.f. $G(w \mid x)$; with a compact support $[\circledR_0(X); \circledR_1(X)]$. The test variable $V$ is continuously distributed with a known or unknown positive probability density function $h(v \mid x)$ having compact support $[\pm_0(X); \pm_1(X)]$ such that $\pm_0(X) \cdot \circledR_0(X)$ and $\pm_1(X) , \circledR_1(X)$: The variables $W$ and $V$ are conditionally independent, given $X$. Let $Z = (X; V; Y)$.

De...ne $m(v; x)$ by

$$m(v; x) = E[Y \mid V = v; X = x]$$

and let $m^{i\,1}$ be the inverse of the function $m$ with respect to its ...rst element (which exists on the support of $W$ given assumption A.1), so if $t = m(v; x)$ then $v = m^{i\,1}(t \mid x)$ for $v \; 2 \; [\circledR_0(X); \circledR_1(X)]$.

**Assumption A.2.** The function $r(w; x)$, chosen by the researcher, is regular, meaning that it is continuous in $(w; x)$ for all $w$ and $x$ on their supports, and for each $x$ is twice continuously di¤erentiable in $w$. Let $\cdot$ be a known constant that is in the support of $W$. The moment $\mu_r(x)$ exists, where $\mu_r(x)$ is de...ned by

$$\mu_r(x) = E[r(W; X) \mid X = x]:$$

De...ne $r^{0}(w; x) = @r(w; x){=}@w$ and $s_r(z)$ by

$$s_r(z) = r(\cdot ; x) + \frac{r^{0}(v; x)[y \; i \; 1(v < \cdot )]}{h(v \mid x)}:$$

For any regular function $r$, Theorem 1 below provides an expression for the conditional mean $\mu_r(x)$. Also provided is the $q'$th conditional quantile of $W$ given $x$, denoted $w_q(x)$:

**Theorem 1.** Let Assumptions A.1 and A.2 hold. Then

$$\mu_r(x) = E[s_r(Z) \mid X = x]:$$
$$w_q(x) = m^{i\,1}(1 \; i \; q \mid x)$$

Proof of Theorem 1. First observe that, given the conditional independence of $W$ and $V$,

$$m(v;x) = E[Y|V = v;X = x] = 1 - G(v|x).$$

Next, by definition, $\mu_r(x) = \int_{\alpha_0(x)}^{\alpha_1(x)} r(v;x)[\partial G(v|x)/\partial v]dv$. Integration by parts yields

$$\mu_r(x) = r(v;x)[G(v|x) - 1(v \geq \cdot)]\big|_{v=\alpha_0(x)}^{\alpha_1(x)} - \int_{\alpha_0(x)}^{\alpha_1(x)} r'(v;x)[G(v|x) - 1(v \geq \cdot)]dv.$$

Therefore, collecting terms we find that

$$\mu_r(x) = r(\cdot;x) + \int_{\alpha_0(x)}^{\alpha_1(x)} r'(v;x)[E(Y|V = v;X = x) - 1(v < \cdot)]dv$$

$$= \int_{\alpha_0(x)}^{\alpha_1(x)} E[s_r(Z)|V = v;X = x)h(v|x)dv$$

$$= E[s_r(Z)|X = x];$$

where the last equality uses the assumptions regarding the supports of $W$ and $V$, and the law of iterated expectation. The conditional quantile expression follows from $G(v|x) = 1 - m(v;x)$. ∎

Theorem 1 provides the basis for the nonparametric moment estimators described in the next section, and for some semiparametric and quantile estimators. Essentially, based on Theorem 1, $\mu_r(x)$ may be estimated as the fitted values of either a parametric or nonparametric regression of $s_r(z)$ on x.

Corollary 1 below will be used to obtain faster converging moment and quantile estimators, based on stronger assumptions.

Assumption A.3. The latent $W$ satisfies $W = g(X;\mu_0) - \varepsilon$, where g is a known function, $\mu_0 \in \Xi$ is a vector of parameters, and $\varepsilon$ is a disturbance that is distributed independently of $V;X$, with unknown, twice continuously differentiable c.d.f. $G^\varepsilon(\varepsilon)$ and compact support $[a_0;a_1]$ that contains zero. Define $U = g(X;\mu_0) - V$. Let $\tilde{a}(U)$ denote the unconditional probability density function of $U$. The support of $U$ contains the interval $[a_0;a_1]$.

Define $s_r^\pi(x;u;y)$ by

$$s_r^\pi(x;u;y) = r[g(x;\mu_0);x] + \frac{r'[g(x;\mu_0) - u;x][y - 1(u > 0)]}{\tilde{a}(u)}.$$

5

**Corollary 1.** Let Assumptions A.1, A.2, and A.3 hold. Then

$$
\begin{aligned}
\tilde{A}(u) &= E[h(g(X;\mu_0) - U) \mid U = u] \\
G^{\pi}(u) &= E(Y \mid U = u) \\
{}^1_r(x) &= E[s^{\pi}_r(x;U;Y)] \\
w_q(x) &= g(x;\mu_0) - G^{\pi - 1}(1 - q)
\end{aligned}
$$

**Proof of Corollary 1.** Having $\tilde{A}(u) = E[h(g(X;\mu_0) - u)]$ follows from the de…nitions of $U$, $\tilde{A}$, and $h$: Also from de…nitions, $Y = I(" < U)$ which implies that $G^{\pi}(u) = E(Y \mid U = u)$. Next, following the same steps as in Theorem 1 we have

$$
\begin{aligned}
{}^1_r(x) &= \int_{a_0}^{a_1} r[g(x;\mu_0) - u; x][\partial G^{\pi}(u){=}\partial u]du \\
&= r[g(x;\mu_0);x] + \int_{a_0}^{a_1} r'[g(x;\mu_0) - u; x][G^{\pi}(u) - 1(u > 0)]du \\
&= \int_{a_0}^{a_1} E[s^{\pi}_r(x;U;Y) \mid U = u]\tilde{A}(u)du = E[s^{\pi}_r(x;U;Y)]
\end{aligned}
$$

Finally, the quantile expression follows from $G(W \mid X = x) = 1 - G^{\pi}[g(X;\mu_0) - W]$. ∎

The advantage of Corollary 1 over Theorem 1 for estimation is that in Corollary 1, ${}^1_r(x)$ and $\tilde{A}(u)$ are expressed as unconditional expectations and so can be estimated using ordinary sample averages (given an estimate of $\mu$). Similarly, using Corollary 1 estimation of the quantiles $w_q(x)$ given $\mu$ only requires estimation of the one dimensional regression $G^{\pi}(u) = E(Y \mid U = u)$, instead of the high dimensional $m(v;x)$.

# 4 Estimators

Assume that a random sample $Z_i = (X_i; V_i; Y_i)$ for $i = 1; \ldots ; n$ is observed, where $V_i$ is a realization of $V$, $Y_i$ is a realization of $Y$, and $X_i$ is a realization of $X$. Theorem 1 and Corollary 1 suggest a number of possible estimators for ${}^1_r(x)$: To describe these estimators, let $\hat{E}$ denote an estimated expectation. An unconditional estimated expectation just denotes the sample average, while a conditional estimated expectation denotes a nonparametric regression.

## 4.1 Nonparametric Estimators

Let Assumptions A.1 and A.2 hold.

If the experimental design, and hence the density function $h$, is known, then $s_r(z_i)$ can be constructed for each observation $i$, and $\mu_r(x)$ may then be consistently estimating by nonparametrically regressing $s_r(z)$ on $x$. This ...rst estimator is

$$\hat{\mu}_{1r}(x) = \hat{E}[s_r(z) \mid X = x]$$

Note that $\hat{\mu}_{1r}(x)$ depends on the design density $h$. One could replace $h(v \mid x)$ with an estimate $\hat{h}(v \mid x)$ (using, e.g., kernel density estimation) in the de...nition of $s_r(z)$. Call the result $\hat{s}_r(z)$. An estimator of $\mu_r(x)$ that can be used when $h$ is unknown is then $\hat{\mu}_{1r}^\pi(x) = \hat{E}[\hat{s}_r(z) \mid X = x]$.

An estimator that does not entail knowing or estimating the density $h$ is the following. Recall that $m(v; x) = E[Y \mid V = v; X = x]$: Let $\hat{m}(v; x)$ be a consistent estimator of $m$, that is, a nonparametric regression of $y$ on $x; v$, so

$$\hat{m}(v; x) = \hat{E}[Y \mid V = v; X = x]$$

Let $a_0$ and $a_1$ be known or estimated constants such that $a_0 \cdot \circledR_0(x)$ and $a_1 \; \circledR_1(x)$. Then, based on the proof of Theorem 1, a consistent estimator of $\mu_r(x)$ is given by

$$\hat{\mu}_{2r}(x) = r(\cdot \; ; x) + \int_{a_0}^{a_1} r^0(v; x)[\hat{m}(v; x) \; \textbf{1}(v < \cdot)]dv;$$

where the integral may be evaluated numerically. We give some more details later about the construction of the estimators.

## 4.2 Conditional Mean Estimation

This section considers estimation of the conditional mean of the unobserved $W$; when this conditional mean is ...nitely parameterized. This structure permits estimation at rate root $n$, instead of the slow convergence of the estimators described in the previous section.

For this section, continue to let Assumptions A.1 and A.2 hold, but assume also that $E(W \mid X = x) = \mu(x; \bar{\theta}_0)$ for all $x \in X$, where $\mu(x; \bar{\theta}_0)$ is a known function and $\bar{\theta}_0 \in B$ is an unknown vector of parameters. The function $\mu(x; \bar{\theta})$ equals $\mu_r(x)$ for $r(w; x) = w$. Whatever other moments may be considered, it is likely that the conditional mean function $\mu$ would be of interest as well.

Assume the identi...cation condition $\Pr[E(W \mid X = x) \notin \mu(X; \bar{\theta})] > 0$ for all $\bar{\theta} \neq \bar{\theta}_0; \bar{\theta} \in B$. It follows from Theorem 1 that

$$\mu(x; \bar{\theta}_0) = E\left[ \frac{Y \; \textbf{1}(V < 0)}{h(V \mid X)} \mid X = x \right]$$

This suggests the conditional mean estimator $\mu(x; \hat{\theta})$; where $\hat{\theta}$ is de...ned by

$$\hat{\theta} = \arg\min_{\bar{\theta} \in B} \frac{1}{n} \sum_{i=1}^{n} \left( \frac{Y_i \; \textbf{1}(V_i < 0)}{h(V_i \mid X_i)} \; \mu(X_i; \bar{\theta}) \right)^2 \hat{\omega}(Z_i);$$

7

where $b(z)$ is a known or estimated positive weight function chosen for e¢ciency. The estimator $b$ is an ordinary nonlinear weighted least squares, and so $b$, and therefore also $^1(x; b)$ is root $n$ consistent and asymptotically normal with a standard limiting distribution, under standard regularity conditions. No nonparametric plug in functions are required. This estimator might not be e¢cient, since it violates the principle of ancillarity due to its dependence on the design density $h$:

If $h$ is not known, one could replace $h(V j X)$ with an estimate $\hat{h}(V j X)$ in the de...nition of $b$. The resulting estimator would then take the form of an ordinary two step estimator with a nonparametric ...rst step (the estimation of $h$) which, with regularity, will be root $n$ consistent and asymptotically normal. This estimator is equivalent to the estimator for general binary choice models proposed by Lewbel (2000), though Lewbel provides other extensions, such as to estimation with endogenous regressors.

## 4.3   Semiparametric Estimators

This section discusses rate root $n$ estimation of arbitrary conditional moments based on Corollary 1. For these estimators we let Assumption A.3 hold, in addition to Assumptions A.1 and A.2. It will be convenient to ...rst consider the case where $\mu_0$ in Assumption 3 is known, implying that the conditional mean of $W$ is known up to an arbitrary location (since " is not required to have mean zero). A special case of known $\mu_0$ is when $x$ is empty, i.e., estimation of unconditional moments of $W$, since in that case we can without loss of generality take $g$ to equal zero.

### 4.3.1   Estimation With Known $\mu$

Assume that $\mu_0$ is known. Considering ...rst the case where the design density $h$ is also known, for a given $u$ de...ne the sample average $\tilde{A}(u)$ by

$$\tilde{A}(u) = \frac{1}{n} \sum_{i=1}^{X} h[g(X_i; \mu_0) ; u]:$$

Then, based on Corollary 1, we have the estimator

$$b_{3r}(x) = r[g(x; \mu_0); x] + \frac{1}{n} \sum_{i=1}^{X} \frac{r^0[g(x; \mu_0) ; U_i; x][Y_i ; 1(U_i > 0)]}{\tilde{A}(U_i)}:$$

This estimator is computationally extremely simple, since it entails only sample averages. Special cases of the estimator $b_{3r}(x)$ were proposed by McFadden (1994) and by Lewbel (1997).

Let $\hat{A}(u)$ be an estimator of $\tilde{A}(u)$ that does not depend on $h$. For example $\hat{A}(u)$ could be a (one dimensional) kernel density estimator of the density of $U$, based on the data $\hat{U}_i$ and evaluated at $u$. We then have the estimator

$$b_{4r}(x) = r[g(x; \mu_0); x] + \frac{1}{n} \sum_{i=1}^{X} \frac{r^0[g(x; \mu_0) ; U_i; x][Y_i ; 1(U_i > 0)]}{\hat{A}(U_i)};$$

which may be used when h is unknown.

Another approach to estimating $\iota_r(x)$ uses ordered observations in place of a preliminary estimate of $\tilde{A}$. Augment the observed $(U; Y)$ pairs with the artificial observations $(\circledR_0; 0)$ and $(\circledR_1; 1)$. Recode each observation $(U; Y)$ having $U < \circledR_0$ as $(\circledR_0; 0)$ and each observation having $U > \circledR_1$ as $(\circledR_1; 1)$. Then, index the observations so that the $U$'s, including the artificial ones, are in non-decreasing order, and denote them by $U_{n0} \cdot \; : \cdot \; U_{n;n+1}$. The probability of ties in the interior of the support is zero. Let $Y_{ni}$ denote the observed $Y$ associated with $U_{ni}$: The proposed ordered data estimator is then

$$\flat_{5r}(x) = r[g(x; \mu_0); x] + \sum_{i=1}^{X} r^0[g(x; \mu_0) \text{ ¡ } U_{ni}; x][Y_{ni} \text{ ¡ } 1(U_{ni} > 0)]\frac{U_{n;i+1} \text{ ¡ } U_{n;i \text{ ¡ } 1}}{2}$$

This estimator can be interpreted as being the same as $\flat_{4r}(x)$, with the density $\tilde{A}(u)$ estimated by differencing the empirical distribution of $U$; although this function, $2n=(U_{n;i+1} \text{ ¡ } U_{n;i \text{ ¡ } 1})$; is not a consistent estimator of $\tilde{A}(U_i)$: The estimator $\flat_{5r}(x)$ has the advantages of just equalling a sample average and of not requiring knowledge of h, however, it can be shown to be less efficient than $\flat_{4r}(x)$ with an optimally estimated $\tilde{A}(U_i)$:

## 4.3.2  Estimation with Unknown μ

Assumptions A.1, A.2, and A.3 imply that $E(W \text{ j } X = x) = \circledR + g(x; \mu)$ for some arbitrary location constant $\circledR$ (since no location constraint is imposed upon "). Therefore, based on the estimator of $\iota(x; \bar{\;})$ described in section 4.2, we may obtain, a root n consistent, asymptotically normal estimate $\hat{\beta}$ of μ by the simple least squares criterion

$$\hat{\beta} = \arg\min_{\mu} \left[\min_{\circledR} \frac{1}{n} \sum_{i=1}^{X} \left(\frac{Y_i \text{ ¡ } 1(V_i < 0)}{h(V_i \text{ j } X_i)} \text{ ¡ } \circledR \text{ ¡ } g(X_i; \mu)\right)^2\right]$$

An estimator $\hat{h}$ may be used in place of h if h is unknown.

Assumptions A.1, A.2, and A.3 make the latent error " independent of $X$, and therefore the binary choice estimator of Klein and Spady (1993) will provide a semiparametrically efficient estimator of μ (note that the Klein and Spady estimator does not identify a location constant $\circledR$, but that is not required, since no location constraint is imposed upon "). We will provide a numerically simpler estimator that is asymptotically equivalent to Klein and Spady.

Let $\hat{\beta}$ denote the chosen root N consistent, asymptotically normal estimator for $\mu_0$. Replacing $\mu_0$ with any $\mu \in \pounds$ we may rewrite the estimators of the previous section as $\flat_{\text{¸}r}(x; \mu)$ for ¸ = 3; 4; or 5. Note that in addition to directly appearing in the equations for $\flat_{\text{¸}r}$, μ also appears in the definition of $U_i = g(X_i; \mu) \text{ ¡ } V_i$: We later derive the root N consistent, asymptotically normal limiting distribution for each estimator $\flat_{\text{¸}r}(x; \hat{\beta})$. The estimators are not differentiable in $U_i$; which complicates the

9

derivation of their limiting distribution (e.g., Theorem 6.1 of Newey and McFadden (1994) is not directly applicable due to this nondi¤erentiability).

# 5 Quantile Estimators

In addition to moments, one may also desire estimates of conditional quantiles of $W$. Let Assumptions A.1 and A.2 hold and de...ne $G(v \mid x) = 1 - m(v; x)$: Then, based on Theorem 1, an estimate of the $q$'th quantile of $W$ given $X = x$ is just

$$\hat{w}_q(x) = \hat{G}^{-1}(q \mid x)$$

The rate of convergence of this estimate will be slow, because of the high dimension of $m$.

If Assumption A.3 holds in addition to A.1 and A.2, then faster convergence is possible. Given Corollary 1 we have $U = g(X; \mu_0) - V$, $G^{\pi}(u) = E(Y \mid U = u)$, and $w_q(x) = g(x; \mu_0) - G^{\pi-1}(1-q)$. Therefore, let $\hat{U}_i = g(X_i; \hat{\mu}) - V_i$ and estimate the conditional quantile $w_q(x)$ by

$$\hat{G}^{\pi}(u) = \hat{E}(Y \mid \hat{U} = u)$$
$$\hat{w}_q(x) = g(x; \hat{\mu}) - \hat{G}^{\pi-1}(1-q)$$

where the function $\hat{G}^{\pi}$ is obtained by nonparametrically regressing $Y$ on $\hat{U}$, and is then numerically inverted to obtain $\hat{G}^{\pi-1}$. This estimator $\hat{w}_q(x)$ will converge at a faster rate than the nonparametric quantile estimator $\hat{w}_q(x)$: With su¢cient regularity, $\hat{w}_q(x)$ is asymptotically normal and converges at the same rate as a one dimensional nonparametric regression estimator, i.e., the same as the best rate that could be obtained if realizations of the latent $w$ were observed.

# 6 Estimation Details and Distribution Theory

In this section we provide a bit more detail about the computation of the estimators $\hat{b}_{jr}(x)$ and their distribution theory.

## 6.1 Nonparametric Estimators

We ...rst consider a fairly general class of nonparametric estimators and then specialize to kernels. Speci...cally, we consider a class of linear estimators of $m$; that is, let

$$\hat{m}(v; x) = \sum_{i=1}^{X} w_{ni}(v; x) Y_i; \tag{1}$$

where $w_{ni}(v;x)$ are some smoothing weights that depend only on $\{(X_1;V_1);\ldots;(X_n;V_n)\}$ and satisfy certain conditions as in Stone (1982). This includes a large class of commonly used estimation schemes such as kernels, local polynomial, nearest neighbor, series, smoothing splines, etc., see Härdle and Linton (1994) for further discussion. It does exclude local median or local quantile estimators whether they be based on kernels or nearest neighbors; it also excludes the popular neural networks class of estimators. Finally, this framework so far excludes methods that have selected the smoothing parameter based on the data. However, it should be possible to extend the treatment to a class of asymptotically linear smoothers, which includes almost any smoothing method that can be asymptotically normal.

We shall suppose also that the smoothing weights in (1) satisfy exactly $\sum_{i=1}^n w_{ni}(v;x) = 1$; this is the case for many linear estimators. In this case, we can write

$$\hat{m}(v;x) - m(v;x) = \sum_{i=1}^n w_{ni}(v;x)\varepsilon_i + \sum_{i=1}^n w_{ni}(v;x)\{m(V_i;X_i) - m(v;x)\}; \tag{2}$$

where the error term $\varepsilon_i = Y_i - m(V_i;X_i)$ is independent across $i$ and satisfies $E(\varepsilon_i|V_i;X_i) = 0$; define also $\sigma_i^2 = var(\varepsilon_i|V_i;X_i)$. The first term on the right hand side of (2) is conditional mean zero and determines the limiting variance, while the second term determines the bias. Under additional conditions, we can approximate the bias term $\sum_{i=1}^n w_{ni}(v;x)\{m(V_i;X_i) - m(v;x)\}$ by $\beta_n\bar{b}(v;x)$ for some bounded and continuous function $\bar{b}(v;x)$ and deterministic sequence $\beta_n \to 0$ as $n \to \infty$. We can also replace the weights $w_{ni}(v;x)$ in $\sum_{i=1}^n w_{ni}(v;x)\varepsilon_i$ by some approximation $\tilde{w}_{ni}(v;x)$ that depends only on $(X_i;V_i)$. We arrive at the expansion

$$\hat{m}(v;x) - m(v;x) = \sum_{i=1}^n \tilde{w}_{ni}(v;x)\varepsilon_i + \beta_n\bar{b}(v;x) + R_n(v;x); \tag{3}$$

where $R_n(v;x)$ is a remainder term that contains the various approximation errors described above. We next state conditions under which $\hat{m}(v;x)$ is asymptotically normal.

**Theorem 2.** Suppose that: (i) $0 < \underline{\sigma}^2 \le \sigma_i^2 \le \bar{\sigma}^2 < \infty$; (ii) $R_n(v;x)=\min\{\rho_n;\beta_n\} \to^p 0$; where $\rho_n = 1/\sqrt{\sum_{i=1}^n \tilde{w}_{ni}^2(v;x)} \to^p 0$; and (iii) $\max_{1\le i\le n} \tilde{w}_{ni}^2(v;x)=\sum_{i=1}^n \tilde{w}_{ni}^2(v;x) \to^p 0$: Then

$$\frac{\hat{m}(v;x) - m(v;x) - \beta_n\bar{b}(v;x)}{\sqrt{\sum_{i=1}^n \tilde{w}_{ni}^2(v;x)\sigma_i^2}} \to^d N(0;1):$$

This result is a standard application of the Lindeberg-Feller central limit theorem. The magnitude of the bias term, $\beta_n$; depends on the method used and on the smoothness of $m$ (and perhaps also on the smoothness of the covariate density). The magnitude of $\rho_n$ depends on the estimation method and on the covariate density in general. The optimal rate in the central limit theorem is achieved when $\rho_n$ and $\beta_n$ are the same magnitude.

By appropriately rede...ning y and m, Theorem 2 can be immediately applied to yield the limiting distribution for $\mathbb{b}_{1r}(x)$.

Now consider estimating $^1{}_r(x)$ using $\mathbb{b}_{2r}(x)$, which is equivalent to

$$\mathbb{b}_{2r}(x) = r(\cdot \ ; x) + \int_{\circledR_0(x)}^{\circledR_1(x)} r^0(v; x)[\text{\ss}(v; x) \ ¡ \ 1(v < \cdot)]dv$$

where $\text{\ss}(v; x)$ is de...ned in (1). This estimator is in the class of marginal integration/partial mean estimators sometimes used for estimating additive nonparametric regression models, see Linton and Nielsen (1995), Newey (1994), and Tjøstheim and Auestad (1994), except that the integrating measure $\varsigma$; where $d_\varsigma(v) = \ ¡ \ r^0(v; x)1(\circledR_0(x) \cdot v \cdot \circledR_1(x))dv$; is not necessarily a probability measure, i.e., it may not be positive or integrate to one. The distribution theory for the class of marginal integration estimators has been worked out for speci...c smoothing methods like kernels or nearest neighbors. We give a derivation at the higher level of generality given by the de...nition (1). Then,

$$\mathbb{b}_{2r}(x) \ ¡ \ ^1{}_r(x) \ = \ \int_{\circledR_0(x)}^{\circledR_1(x)} r^0(v; x)[\text{\ss}(v; x) \ ¡ \ m(v; x)]dv$$

$$= \ \sum_{i=1}^{n} \overline{w}_{ni}(x)"_i + \pm_n \overline{\phantom{x}}(x) + \overline{R}_n(x);$$

where: $\overline{w}_{ni}(x) = \int w_{ni}(v; x)d_\varsigma(v)$; $\overline{\phantom{x}}(x) = \int \bar{\phantom{x}}(v; x)d_\varsigma(v)$; and $\overline{R}_n(x) = \int R_n(v; x)d_\varsigma(v)$: We next state conditions under which $\mathbb{b}_{2r}(x)$ is asymptotically normal.

**Theorem 3.** Suppose that condition (i) from Theorem 2 is true, and that: (i) $\overline{R}_n(x) = \min f\circledR_n; \pm_n g \ ! \ ^p \ 0;$ where $\circledR_n = 1 = \sqrt{\sum_{i=1}^{n} \overline{w}_{ni}^2(x)} \ ! \ ^p \ 0;$ and (ii) $\max_{1 \cdot i \cdot n} \overline{w}_{ni}^2(x) = \sum_{i=1}^{n} \overline{w}_{ni}^2(x) \ ! \ ^p \ 0.$ Then,

$$\frac{\mathbb{b}_{2r}(x) \ ¡ \ ^1{}_r(x) \ ¡ \ \pm_n \overline{\phantom{x}}(x)}{\sqrt{\sum_{i=1}^{n} \overline{w}_{ni}^2(x) \frac{3}{4}_i^2}} \ ¡ \ ^d \ N(0; 1): \tag{4}$$

Note that the magnitude of the bias, $\pm_n$; is the same for $\mathbb{b}_{2r}(x)$ as for $\text{\ss}(v; x)$: However, the magnitude of the asymptotic variance of $\mathbb{b}_{2r}(x)$, which is $\circledR_n^2$; can be expected to be of smaller magnitude than $\circledR_n^2$ [i.e., the asymptotic variance of $\text{\ss}(v; x)$] by virtue of the integration. This specially a¤ects the veri...cation of condition 3(i) because we must make the remainder term in Theorem 3 of smaller order than those in Theorem 2. In the next subsection we verify the conditions of Theorem 3 for a kernel estimator that falls in the class de...ned by (1). The optimal rate will balance $\circledR_n$ with $\pm_n$:

### 6.1.1 Veri...cation of Conditions for Kernels

The Nadaraya-Watson kernel estimator has weights

$$w_{ni}(v; x) = \frac{k \left(\frac{v_i \ V_i}{\mathfrak{z}}\right) K \left(\frac{x_i \ X_i}{\mathfrak{z}}\right)}{\sum_{i=1}^{n} k \left(\frac{v_i \ V_i}{b}\right) K \left(\frac{x_i \ X_i}{b}\right)}$$

in (1), where $k$ is a kernel function and $K(t) = \prod_{j=1}^{d} k(t_j)$. De...ne $\mu_2(k) = \int t^2 k(t)dt$: Let $\nabla$; $\nabla^2$ denote the ...rst and second derivative operators.

**Theorem 4**. Suppose that assumptions B1 and B2 in the appendix hold and that the bandwidth sequence $b = b(n)$ satis...es $b \to 0$ and $nb^{d+2} = \log n \to 1$: Then, (4) holds with $\pm_n = b^2$;

$$\bar{}(v;x) = \frac{\mu_2(k)}{2}\mathrm{tr}(\nabla^2 m(v;x) + \nabla m(v;x)\nabla s(v;x));\tag{5}$$

where $s(v;x) = \log f_{V;X}(v;x)$; while

$$nb^d \sum_{i=1}^{n} \overline{w}_{ni}^2(x)\frac{3}{4}_i^2 \to \left\| kKk^2 \right\| \int_{\circledR_0(x)}^{\circledR_1(x)} \frac{3}{4}^2(v;x) \left( \frac{r^{\emptyset}(v;x)}{f_{V;X}(v;x)} \right)^2 f_{V;X}(v;x)dv \; \acute{} \; \bar{\top}(x):$$

Thus $\mathbf{b}_{2r}(x)$ is asymptotically normal with mean $\mu_r(x) + b^{2}\bar{}(x)$ and variance $n^i\, {}^1b^{i\, d}\bar{\top}(x)$:

## 6.2   Semiparametric Estimators

### 6.2.1   Estimation of $\bar{}_0$

Consider brie‡y the conditional mean estimator $\mu(X_i; \mathbf{b})$ where

$$\mathbf{b} = \arg\min \frac{1}{n} \sum_{i=1}^{n} \frac{Y_i i 1(V_i < 0)}{\hat{h}(V_i j X_i)} i \mu(X_i; \bar{})^2 \mathbf{b}(Z_i);$$

If $\hat{h}(V_i j X_i)$ is replaced with a known design density $h(V_i j X_i)$ in the above, then the limiting distribution is given by ordinary weighted nonlinear least squares. The root $n$ limiting distribution of $\mathbf{b}$ using an asymptotically trimmed kernel estimator of $\hat{h}(V_i j X_i)$ is given by Lewbel (2000) (for the case where $\mu$ is linear). If $V$ is independent of $X$; then $\hat{h}(V_i j X_i)$ can be replaced with

$$\hat{h}(V_i) = \frac{1}{nb} \sum_{j=1}^{n} k\left( \frac{V_i i V_j}{b} \right):$$

and the resulting estimator $\mathbf{b}$ is a special case of the general theory covered by Andrews (1994, Theorem 1) and Newey and McFadden (1994, section 8).

### 6.2.2   Estimation of $\mu_0$

Consider ...rst the estimator $\hat{\mu}$ de...ned by

$$\hat{\mu} = \arg\min_{\mu 2 \pounds} \min_{\circledR} \frac{1}{n} \sum_{i=1}^{n} \frac{Y_i i 1(V_i < 0)}{\hat{h}(V_i j X_i)} i \circledR i g(X_i; \mu)^2 \mathbf{b}(Z_i);$$

The results of the previous section can be immediately applied taking $\bar{} = (\circledR; \mu)$ since $\mu(X; \bar{}) = \circledR i g(X; \mu)$:

To obtain a semiparametrically e¢cient estimator of $\mu_0$ we apply one-step estimation to Klein and Spady (1993), starting from this initial root-n consistent $\hat{\beta}$. De…ne $\bar{\imath}_i(\mu) = E[Y \mid g(X; \mu) + V = g(X_i; \mu) + V_i] = G^{\alpha}[g(X_i; \mu) + V_i]$, and

$$\hat{\mathbb{p}}_i(\mu) = \frac{\sum_{j=1}^{n} Y_j k \left[\frac{g(X_i;\mu)+V_i \mathbf{i} \; g(X_j;\mu)\mathbf{i} \; V_j}{b}\right]^{3}}{\sum_{j=1}^{n} k \left[\frac{g(X_i;\mu)+V_i \mathbf{i} \; g(X_j;\mu)\mathbf{i} \; V_j}{b}\right]^{3}}$$

for every $\mu$: Let

$$\mathbb{Q}(\mu) = \frac{1}{n} \sum_{i=1}^{X} Y_i \ln[\hat{\mathbb{p}}_i(\mu)] + (1 \mathbf{i} \; Y_i) \ln[1 \mathbf{i} \; \hat{\mathbb{p}}_i(\mu)]:$$

The semiparametrically e¢cient Klein and Spady estimator is $\bar{\mu} = \arg\sup_{\mu 2 \mathcal{E}} \mathbb{Q}(\mu)$, which satis…es the …rst order conditions $\mathbb{Q}^0(\bar{\mu}) = 0$. Given the well known problems with computing the Klein and Spady estimator we instead propose the one-step estimator

$$\hat{\beta} = \hat{\beta} \mathbf{i} \; [\hat{H}(\hat{\beta})]^{\mathbf{i} \; 1} \mathbb{Q}^0(\hat{\beta}); \qquad (6)$$

where

$$\mathbb{Q}^0(\hat{\beta}) = \frac{1}{n} \sum_{i=1}^{X} \left( \frac{Y_i}{\hat{\mathbb{p}}_i(\hat{\beta})} \mathbf{i} \; \frac{1 \mathbf{i} \; Y_i}{1 \mathbf{i} \; \hat{\mathbb{p}}_i(\hat{\beta})} \right) \hat{\mathbb{p}}_i^0(\hat{\beta})$$

and

$$\hat{H}(\hat{\beta}) = \frac{1}{n} \sum_{i=1}^{X} \left( \frac{\mathbf{i} \; 1}{\hat{\mathbb{p}}_i(\hat{\beta})} + \frac{\mathbf{i} \; 1}{1 \mathbf{i} \; \hat{\mathbb{p}}_i(\hat{\beta})} \right) \hat{\mathbb{p}}_i^0(\hat{\beta}) \hat{\mathbb{p}}_i^0(\hat{\beta})^0;$$

and in which $\hat{\mathbb{p}}_i^0(\mu) = @\hat{\mathbb{p}}_i(\mu)=@\mu$. The initial condition $\hat{\beta}$ is any root-N consistent estimator of $\mu$. Two-step estimation in semiparametric models have been examined in some detail in the monograph Bickel, Klaassen, Ritov and Wellner (1993) and the references therein. Under some regularity conditions, it can be shown that $\hat{\beta}$ has the same asymptotic distribution as $\bar{\mu}$:

### 6.2.3 Estimation of $^1$ with estimated $\mu$

Here we state the asymptotic properties of the conditional moment estimators based on Corollary 1. De…ne

$$\hat{\mathbb{B}}_{3r}^{\alpha}(x;\hat{\beta}) = r[g(x;\hat{\beta}); x] + \frac{1}{n} \sum_{i=1}^{X} \frac{r^0[g(x;\hat{\beta}) \mathbf{i} \; \hat{\mathbb{b}}_i; x][Y_i \mathbf{i} \; 1(\hat{\mathbb{b}}_i > 0)]}{\hat{A}(\hat{\mathbb{b}}_i)}$$

$$\hat{\mathbb{B}}_{4r}^{\alpha}(x;\hat{\beta}) = r[g(x;\hat{\beta}); x] + \frac{1}{n} \sum_{i=1}^{X} \frac{r^0[g(x;\hat{\beta}) \mathbf{i} \; \hat{\mathbb{b}}_i; x][Y_i \mathbf{i} \; 1(\hat{\mathbb{b}}_i > 0)]}{\tilde{A}(\hat{\mathbb{b}}_i)};$$

where $\hat{\mathbb{b}}_i = g(X_i; \hat{\beta}) \mathbf{i} \; V_i$ and

$$\hat{A}(\hat{\mathbb{b}}_i) = \frac{1}{n} \sum_{j=1}^{X} h[g(X_j; \hat{\beta}) \mathbf{i} \; \hat{\mathbb{b}}_i] \quad ; \quad \tilde{A}(\hat{\mathbb{b}}_i) = \frac{1}{nb} \sum_{j=1}^{X} k \left( \frac{\hat{\mathbb{b}}_i \mathbf{i} \; \hat{\mathbb{b}}_j}{b} \right):$$

We shall suppose that

$$\sqrt{n}(\hat{\mu} - \mu_0) = \frac{1}{\sqrt{n}} \sum_{i=1}^{n} m(Z_i; \mu_0) + o_p(1)$$

for some function m that is mean zero and has ...nite variance. The estimators $\hat{b}_{3r}(x)$ and $\hat{b}_{4r}(x)$ are the special cases of $\hat{b}_{3r}^{\pi}(x; \hat{\mu})$ and $\hat{b}_{4r}^{\pi}(x; \hat{\mu})$ in which $\mu$ is known, and so correspond to the case of m being identically zero.

For each $\mu \in \pounds$ and $x \in X$ let

$$f_0(Z_i; \mu) = \frac{r^0[g(x; \mu) - U_i(\mu); x][Y_i - 1(U_i(\mu) > 0)]}{\tilde{A}(U_i)}$$

$$f_1(Z_i; \mu) = r[g(x; \mu); x] + \frac{r^0[g(x; \mu) - U_i(\mu); x][Y_i - 1(U_i(\mu) > 0)]}{\tilde{A}(U_i)}$$

$$\tilde{A}_{iF} = \left. \left( \frac{@}{@\mu} E[f_1(Z_i; \mu)] \right) \right|_{\mu=\mu_0};$$

where $U_i(\mu) = g(X_i; \mu) - V_i$: The quantities $f_0; f_1$ and $\tilde{A}_{iF}$ depend on x but we have suppressed this notationally. Note also that $E f_1(Z_i; \mu_0) = \textbf{¹}_r(x)$: Finally, let

$$\circ_i = \frac{@g}{@\mu}(X_i; \mu_0) - E\left[ \frac{@g}{@\mu}(X_i; \mu_0) \right]:$$

**Theorem 5.** Suppose that Assumptions C1-C3 in the Appendix hold. Then, as $n \to \infty$;

$$\sqrt{n}[\hat{b}_{3r}^{\pi}(x; \hat{\mu}) - \textbf{¹}_r(x; \mu_0)] = \frac{1}{\sqrt{n}} \sum_{j=1}^{n} \acute{}_j + o_p(1); \tag{7}$$

where $\acute{}_j = \acute{}_{1j} + \acute{}_{2j} + \acute{}_{3j}$; with:

$$\acute{}_{1j} = f_1(Z_j; \mu_0) - E f_1(Z_j; \mu_0)$$

$$\acute{}_{2j} = \left\{ \tilde{A}_{iF} - E\left[ f_0(Z_i; \mu_0) \frac{\tilde{A}^0(U_i)}{\tilde{A}(U_i)} \circ_i \right] \right\} m(Z_j; \mu_0)$$

$$\acute{}_{3j} = - E\left[ f_0(Z_i; \mu_0) \frac{h[g(X_j; \mu_0) - U_i] - \tilde{A}(U_i)}{\tilde{A}(U_i)} jX_j \right]:$$

This implies that $\sqrt{n}[\hat{b}_{3r}^{\pi}(x; \hat{\mu}) - \textbf{¹}_r(x; \mu_0)]$ is asymptotically normal with mean zero and variance $\S_{\acute{}} = \text{var}(\acute{}_j)$ by the Lindeberg-Feller central limit theorem: The three terms $\acute{}_{1j}; \acute{}_{2j};$ and $\acute{}_{3j}$ are all mean zero and have ...nite variance. They are generally mutually correlated. When $\mu_0$ is known, the term $\acute{}_{2j} = 0$ and this term is missing from the asymptotic expansion.

The result can be extended to a functional central limit theorem in x because the $o_p(1)$ term in (7) is uniform in $x \in X$ and the stochastic process $\acute{}_j(x)$ is tight in x due to the smoothness properties of r:

15

We next give the distribution theory for the semiparametric estimator $\hat{b}^{\pi}_{4r}(x;\hat{\beta})$. Let

$$\circ^{\pi}_i = \frac{@g}{@\mu^0}(X_i;\mu_0) \; i \; E\left[\frac{@g}{@\mu^0}(X_i;\mu_0)jU_i\right];$$

where $\acute{A}_u(U_i) = E[@g(X_i;\mu_0)=@\mu jU_i]$.

**Theorem 6.** Suppose that assumptions B1,B2 and C1-C4 in the Appendix hold. Then

$$\stackrel{P}{\overline{n}}[\hat{b}^{\pi}_{4r}(x;\hat{\beta}) \; i \; ^1_r(x)] = \stackrel{P}{\overline{n}}\frac{1}{=}\sum_{i=1}^{X^n} \; ^{\pi}_i + o_p(1);$$

where $^{\pi}_j = \; ^{\pi}_{1j} + \; ^{\pi}_{2j} + \; ^{\pi}_{3j};$ with $^{\pi}_{1j} = \; ^{\prime}_{1j};$ while

$$^{\pi}_{2j} = \; i \; F \; i \; E\left[f_0(Z_i;\mu_0)\left[\frac{\tilde{A}^0(U_i)}{\tilde{A}(U_i)}\circ^{\pi}_i \; i \; A^0_u(U_i)\right]\right] m(Z_j;\mu_0)$$

$$^{\pi}_{3j} = \; i \; \frac{r^0[g(x;\mu_0) \; i \; U_j;x]}{\tilde{A}(U_j)}(Y_j \; i \; E[Y_jjU_j]):$$

This implies that $\stackrel{P}{\overline{n}}[\hat{b}^{\pi}_{4r}(x;\hat{\beta}) \; i \; ^1_r(x;\mu_0)]$ is asymptotically normal with mean zero and variance $\S^{\pi} = var(^{\pi}_j)$. The three terms $^{\pi}_{1j};\;^{\pi}_{2j};$ and $^{\pi}_{3j}$ are all mean zero and have ...nite variance. They are generally correlated. When $\mu_0$ is known, the term $^{\pi}_{2j} = 0$ and this term is missing from the asymptotic expansion.

**Remarks**

1. Standard errors can be constructed by substituting population quantities by estimated ones. That is, let $\hat{\S} = n^{i\,1}\sum_{j=1}^{n}\hat{b}^2_j;$ where $\hat{b}_j = \hat{b}_{1j} + \hat{b}_{2j} + \hat{b}_{3j}$ and for example

$$\hat{b}_{1j} = \frac{r^0[g(x;\hat{\beta}) \; i \; \hat{\beta}_i;x][Y_i \; i \; 1(\hat{\beta}_i > 0)]}{\hat{A}(\hat{\beta}_i)} \; i \; \frac{1}{n}\sum_{i=1}^{X^n}\frac{r^0[g(x;\hat{\beta}) \; i \; \hat{\beta}_i;x][Y_i \; i \; 1(\hat{\beta}_i > 0)]}{\hat{A}(\hat{\beta}_i)}:$$

To construct $\hat{b}_{2j}$ we replace $E[f_0(Z_i;\mu_0)(\tilde{A}^0(U_i)=\tilde{A}(U_i))\circ_i]$ in $^{\prime}_{2j}$ by

$$\frac{1}{n}\sum_{i=1}^{X^n}\frac{r^0[g(x;\hat{\beta}) \; i \; \hat{\beta}_i;x][Y_i \; i \; 1(\hat{\beta}_i > 0)]}{\hat{A}(\hat{\beta}_i)}\frac{\hat{A}^0(\hat{\beta}_i)}{\hat{A}(\hat{\beta}_i)}\left[@\frac{@g}{@\mu}(X_i;\hat{\beta}) \; i \; \frac{1}{n}\sum_{j=1}^{X^n}\frac{@g}{@\mu}(X_j;\hat{\beta})\right]:$$

For the quantity $i\;F$ we must use numerical derivatives, i.e.,

$$\hat{\beta}_F = \frac{1}{n}\sum_{i=1}^{X^n}\frac{\hat{f}_1(Z_i;\hat{\beta} + \pm e_k) \; i \; \hat{f}_1(Z_i;\hat{\beta})}{\pm};$$

where $e_k$ is the elementary vector in direction $k$ and $\pm$ is a small number, while

$$\hat{f}_1(Z_i;\hat{\beta}) = r[g(x;\hat{\beta});x] + \frac{r^0[g(x;\hat{\beta}) \; i \; U_i(\hat{\beta});x][Y_i \; i \; 1(U_i(\hat{\beta}) > 0)]}{\hat{A}(U_i(\hat{\beta}))}:$$

16

For the conditional expectation in $\acute{}_{3j}$ we should use a kernel regression smoother on the estimated quantities.

2. Regarding e¢ciency, it is not possible to provide a ranking of the two estimators $\mathfrak{b}_{3r}^{¤}(x;\hat{\mu})$ and $\mathfrak{b}_{4r}^{¤}(x;\hat{\mu})$ uniformly throughout the 'parameter space'. However, one step in that direction might be to use a semiparametrically e¢cient estimator of $\mu_0$: It may be possible to develop an e¢ciency bound for estimation of the function $^1_r(:)$ by following the calculations of Bickel, Klaassen, Ritov and Wellner (1993, Chapter 5). Since there are no additional restrictions on $^1_r$, the plug-in estimator with e¢cient $\hat{\mu}$ should be e¢cient.

### 6.2.4 Quantile estimators

The distribution theory is trivial. The estimator $\hat{w}_q(x) = \hat{\mathbb{G}}^{i1}(q \mid x)$ has the distribution theory for standard conditional quantile estimators. The distribution theory for $\hat{w}_q(x) = g(x;\hat{\mu}) _i \hat{\mathbb{G}}^{¤i1}(1 _i q)$ is the same as the distribution theory for $\hat{w}_q(x) = g(x;\mu_0) _i \mathbb{G}^{¤i1}(1 _i q)$; where

$$\mathbb{G}^{¤}(u) = \hat{\mathbb{P}}(Y \mid U = u);$$

which is again basically a standard one-dimensional conditional quantile estimator. This is because $\hat{\mu}$ converges at rate root-n, so the estimation error in $\hat{\mu}$ is asymptotically irrelevant given the slower convergence rate of quantiles.

# 7 Conclusions

We have provided some estimators of conditional moments and quantiles of the latent $W$: We have for convenience assumed throughout that the support of $V$ (which must contain the support of $W$) is bounded. Most of the results in this paper should extend to the in...nite support case, although some of the estimators may then require asymptotic trimming to deal with issues arising from division by a density estimate when the true density is not bounded away from zero.

The precision of these estimators depends in part on the density h. When designing experiments one may wish to choose h to maximize e¢ciency based on the variance estimators.

# 8 Appendix

## 8.1 Regularity Conditions

We ...rst state some regularity conditions that are needed for the nonparametric estimation of h:

**Assumption B.1.** k is a symmetric probability density with bounded support, and is Lipschitz continuous on its support, i.e.,

$$jk(t) \; i \; k(s)j \cdot \; cjt \; i \; sj$$

for some constant c.

**Assumption B.2.** The variables $(V; X)$ are continuously distributed with Lebesgue density $f_{V;X}(v; x)$ that satis...es $\inf_{\circledR_0(x) \cdot \; v \cdot \; \circledR_1(x)} f_{V;X}(v; x) > 0$: Furthermore, m and $f_{V;X}$ are twice continuously di¤erentiable for all v with $\circledR_0(x) \cdot \; v \cdot \; \circledR_1(x)$. The set $[\circledR_0(x); \circledR_1(x)] \; £ \; fxg$ is strictly contained in the support of $(V; X)$:

We also need some conditions on the estimator and on the regression functions and densities.

**Assumption C.1.** Suppose that

$$\frac{p}{n}(\beta \; i \; \mu_0) = \frac{1}{p_n} \sum_{i=1}^{X} m(Z_i; \mu_0) + o_p(1)$$

for some function m such that $E[m(Z_i; \mu_0)] = 0$ and $- = E[m(Z_i; \mu_0)m(Z_i; \mu_0)^0] < 1$:

**Assumption C.2.** The function g is twice continuously di¤erentiable in μ and

$$\sup_{k\mu_i \; \mu_0k \cdot \; \pm_n} \left| \frac{@g}{@\mu}(x; \mu) \right| \cdot \; d_1(x) \quad ; \quad \sup_{k\mu_i \; \mu_0k \cdot \; \pm_n} \left| \frac{@^2g}{@\mu@\mu^0}(x; \mu) \right| \cdot \; d_2(x)$$

with $E d_1(X_i) < 1$ and $E d_2(X_i) < 1$:

**Assumption C.3.** The density function h is continuous and is strictly positive on its support and is twice continuously di¤erentiable.

**Assumption C.4.** The kernel k is twice continuously di¤erentiable on its support, and therefore $\sup_t jk^{00}(t)j < 1$ : The bandwidth b satis...es $b \; ! \; 0$ and $nb^6 \; ! \; 1$ :

## 8.2 Distribution Theory for Nonparametric Estimators

**Proof of Theorem 4.** Under Assumptions B1 and B2 the expansion (3) holds with $\pm_n = b^2$ and $^-(v; x)$ is as stated in (5), the weights

$$w_{ni}(v; x) = \frac{1}{f_{V;X}(v; x)} \frac{1}{nb^{d+1}} K\left(\frac{v \; i \; V_i}{b}\right) K\left(\frac{x \; i \; X_i}{b}\right);$$

while the remainder term satis...es

$$\sup_{\circledR_0(x) \cdot \; v \cdot \; \circledR_1(x)} jR_n(v; x)j = O_p\left(\sqrt{\frac{\log n}{nb^{d+1}}}\right) + o_p(b^2):$$

See for example Masry (1996a, 1996b).

Provided $nb^{d+2}/\log n \to 1$; condition (i) of Theorem 3 is satis...ed because $\mathcal{T}_n = O_p(1/\sqrt{nb^d})$ as we now show. We have

$$
E[\overline{w}_{ni}(x)] = \frac{1}{nb^d} K\left(\frac{x_i - X_i}{b}\right) \frac{r^0(V_i; x)}{f_{V;x}(V_i; x)}
$$

$$
= \frac{1}{nb^d} K\left(\frac{x_i - X_i}{b}\right)\left[ \int_{\circledR_0(x)}^{\circledR_1(x)} \frac{r^0(v; x)}{f_{V;x}(v; x)} \frac{1}{b} K\left(\frac{v_i - V_i}{b}\right) dv - \frac{r^0(V_i; x)}{f_{V;x}(V_i; x)} \right]
$$

$$
= \frac{1}{nb^d} K\left(\frac{x_i - X_i}{b}\right) \int_{\circledR_0(x)}^{\circledR_1(x)} \left[ \frac{r^0(v; x)}{f_{V;x}(v; x)} - \frac{r^0(V_i; x)}{f_{V;x}(V_i; x)} \right] \frac{1}{b} K\left(\frac{v_i - V_i}{b}\right) dv \quad \text{for large } n
$$

$$
= O(b^2)
$$

by a change of variables and dominated convergence argument that is in wide use in nonparametrics (see, e.g., Newey and McFadden 1994 section 8). It works in this case because the set $[\circledR_0(x); \circledR_1(x)]$ is contained in the support of $V$ and the conditions on $K$ etc. Therefore, the asymptotic variance of $\mathcal{B}_{2r}(x)$ is

$$
\sum_{i=1}^{\aleph} \overline{w}_{ni}^2(x)\frac{3}{4}_i^2 = \frac{1}{n^2 b^{2d}} \sum_{i=1}^{\aleph} K^2\left(\frac{x_i - X_i}{b}\right)\left[ \frac{r^0(V_i; x)}{f_{V;x}(V_i; x)} \right]^2 \frac{3}{4}_i^2 = O_p(n^{-1} b^{-d});
$$

as follows from Markov's inequality. Therefore, $\mathcal{T}_n = O_p(1/\sqrt{nb^d})$ as required. In fact, $\sum_{i=1}^n \overline{w}_{ni}^2(x)\frac{3}{4}_i^2$ satis...es a law of large numbers and is approximately

$$
\frac{1}{nb^d} E\left[ \frac{1}{b^{dx}} K^2\left(\frac{x_i - X_i}{b}\right)\left[ \frac{r^0(V_i; x)}{f_{V;x}(V_i; x)} \right]^2 \frac{3}{4}^2(V_i; X_i) \right]
$$

$$
= \frac{1}{nb^d} kKk^2 \int \frac{3}{4}^2(v; x) \left[ \frac{r^0(v; x)}{f_{V;x}(v; x)} \right]^2 f_{V;x}(v; x) dv;
$$

where $\frac{3}{4}^2(V_i; X_i) = \frac{3}{4}_i^2$; by a change of variables and dominated convergence. Furthermore, condition (ii) of Theorem 3 is satis...ed by the arguments used in Gozalo and Linton (1999, Lemma CLT). $\blacksquare$

## 8.3  Distribution Theory for Semiparametric Quantities

Let $E_i$ denote expectation conditional on $Z_i$:

**Proof of Theorem 5.** Recall that

$$
\mathcal{B}_{3r}^{\alpha}(x; \mathcal{B}) = r[g(x; \mathcal{B}); x] + \frac{1}{n} \sum_{i=1}^{\aleph} \frac{r^0[g(x; \mathcal{B}) - \mathcal{B}_i; x][Y_i - 1(\mathcal{B}_i > 0)]}{\mathcal{A}(\mathcal{B}_i)};
$$

where $\mathcal{B}_i = g(X_i; \mathcal{B}) - V_i$ and

$$
\mathcal{A}(\mathcal{B}_i) = \frac{1}{n} \sum_{j=1}^{\aleph} h[g(X_j; \mathcal{B}) - \mathcal{B}_i]:
$$

19

By a geometric series expansion we can write

$$B_{3r}^n(x;\hat{\beta}) = \frac{1}{n}\sum_{i=1}^n f_1(Z_i;\hat{\beta}) - \frac{1}{n}\sum_{i=1}^n f_2(Z_i;\mu_0)[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]$$

$$- \frac{1}{n}\sum_{i=1}^n [f_2(Z_i;\hat{\beta}) - f_2(Z_i;\mu_0)][\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]$$

$$+ \frac{1}{n}\sum_{i=1}^n \frac{r^0[g(x;\hat{\beta}) - \hat{\theta}_i;x][Y_i - 1(\hat{\theta}_i > 0)]}{\tilde{A}^2(U_i)\hat{A}(\hat{\theta}_i)}[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]^2;$$

where

$$f_2(Z_i;\mu) = \frac{r^0[g(x;\mu) - U_i(\mu);x][Y_i - 1(U_i(\mu) > 0)]}{\tilde{A}^2(U_i)}.$$

**Leading Terms.** We make use of Lemmas 1 and 2 given below. Lemma 1 implies that

$$\sqrt{n}\frac{1}{n}\sum_{i=1}^n [f_1(Z_i;\hat{\beta}) - Ef_1(Z_i;\mu_0)] = \sqrt{n}\frac{1}{n}\sum_{i=1}^n f_{i\ F}m(Z_i;\mu_0) + [f_1(Z_i;\mu_0) - Ef_1(Z_i;\mu_0)]g + o_p(1): \quad (8)$$

Furthermore, by Lemma 2

$$\left|\frac{1}{n}\sum_{i=1}^n f_2(Z_i;\mu_0)[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i) - \frac{1}{n}\sum_{j=1}^n L(Z_i;Z_j)]\right|$$

$$\cdot \frac{1}{n}\sum_{i=1}^n |f_2(Z_i;\mu_0)| \le \max_{1 \le i \le n}\left|\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i) - \frac{1}{n}\sum_{j=1}^n L(Z_i;Z_j)\right|$$

$$= o_p(n^{-1=2});$$

where $L(Z_i;Z_j) = \gg_j(U_i) + \jmath(Z_i)m(Z_j;\mu_0)$; and

$$\gg_j(u) = h[g(X_j;\mu_0) - u] - E(h[g(X_j;\mu_0) - u])$$

$$\jmath(Z_i) = \tilde{A}^0(U_i)\left[\frac{@g}{@\mu}(X_i;\mu_0) - E\left[\frac{@g}{@\mu}(X_i;\mu_0)\right]\right]:$$

Then

$$\frac{1}{n^2}\sum_{i=1}^n\sum_{j=1}^n f_2(Z_i;\mu_0)L(Z_i;Z_j) = \sum_{i=1}^n\sum_{j=1}^n \varsigma_n(Z_i;Z_j)$$

$$= \frac{1}{n}\sum_{j=1}^n !(Z_j) + o_p(n^{-1=2});$$

where

$$!(Z_j) = E[f_2(Z_i;\mu_0)\jmath(Z_i)]m(Z_j;\mu_0) + E[f_2(Z_i;\mu_0)\gg_j(U_i)|Z_j]$$

20

by standard U-statistic theory. We have

$$E\left[f_2(Z_i;\mu_0)\varkappa_j(U_i)|Z_j\right] = E\left[f_2(Z_i;\mu_0)h[g(X_j;\mu_0)-U_i]|Z_j\right] - E\left[f_2(Z_i;\mu_0)\tilde{A}(U_i)\right]$$

$$= E\left[(f_1(Z_i;\mu_0)-r[g(x;\mu_0);x])\frac{h[g(X_j;\mu_0)-U_i]}{\tilde{A}(U_i)}|Z_j\right]$$

$$- E\left[(f_1(Z_i;\mu_0)-r[g(x;\mu_0);x])\right]$$

$$= E\left[\frac{r'[g(x;\mu)-U_i;x][Y_i-1(U_i>0)]}{\tilde{A}(U_i)}\frac{h[g(X_j;\mu_0)-U_i]-\tilde{A}(U_i)}{\tilde{A}(U_i)}|Z_j\right]$$

$$E\left[f_2(Z_i;\mu_0)-\lambda(Z_i)\right] = E\left[f_2(Z_i;\mu_0)\tilde{A}(U_i)\frac{\tilde{A}'(U_i)}{\tilde{A}(U_i)}\left\{\frac{\partial g}{\partial \mu}(X_i;\mu_0)-E\left[\frac{\partial g}{\partial \mu}(X_i;\mu_0)\right]\right\}\right]$$

$$= E\left[\frac{r'[g(x;\mu_0)-U_i;x][Y_i-1(U_i>0)]\tilde{A}'(U_i)}{\tilde{A}(U_i)}\frac{\tilde{A}'(U_i)}{\tilde{A}(U_i)}\circ_i\right];$$

so that the leading terms are as stated.

Remainders. By the Cauchy-Schwarz inequality

$$\left|\frac{1}{n}\sum_{i=1}^{n}[f_2(Z_i;\hat\theta)-f_2(Z_i;\mu_0)][\hat{A}(\hat\theta_i)-\tilde{A}(U_i)]\right|$$

$$\cdot \left\{\frac{1}{n}\sum_{i=1}^{n}[f_2(Z_i;\hat\theta)-f_2(Z_i;\mu_0)]^2\right\}^{1=2}\left\{\frac{1}{n}\sum_{i=1}^{n}[\hat{A}(\hat\theta_i)-\tilde{A}(U_i)]^2\right\}^{1=2}$$

$$= O_p(n^{-1})$$

from another application of Lemmas 1 and 2:

We have assumed that $\inf_{u\in U}\tilde{A}(u)>0$; which implies that

$$\min_{1\le i\le n}\tilde{A}(\hat\theta_i) = \inf_{u\in U}\tilde{A}(u) + O_p(n^{-1=2})$$

is bounded away from zero with probability tending to one. Therefore,

$$\left|\frac{1}{n}\sum_{i=1}^{n}\frac{r'[g(x;\hat\theta)-\hat\theta_i;x][Y_i-1(\hat\theta_i>0)]}{\tilde{A}^2(\hat\theta_i)\hat{A}(\hat\theta_i)}[\hat{A}(\hat\theta_i)-\tilde{A}(\hat\theta_i)]^2\right|$$

$$\cdot \frac{\sup_{u\in U}[\hat{A}(u)-\tilde{A}(u)]^2 + O_p(n^{-1=2})}{\inf_{u\in U}\tilde{A}^2(u)\hat{A}(u) + O_p(n^{-1=2})}\frac{1}{n}\sum_{i=1}^{n}|r'[g(x;\hat\theta)-\hat\theta_i;x]|\cdot(|Y_i|+1)$$

$$= O_p(n^{-1}):$$

21

In conclusion,

$$\sqrt{n}[\hat{\beta}^{\pi}_{3r}(x;\hat{\theta}) - {}^{1}_{r}(x;\mu_0)] = \frac{1}{\sqrt{n}}\sum_{i=1}^{n} \acute{}_i + o_p(1);$$

as required. The asymptotic distribution of $\sqrt{n}[\hat{\beta}^{\pi}_{3r}(x;\hat{\theta}) - {}^{1}_{r}(x;\mu_0)]$ follows from the central limit theorem for independent random variables.

∎

**Proof of Theorem 6.** By a geometric series expansion we can write

$$\hat{\beta}^{\pi}_{4r}(x;\hat{\theta}) = \frac{1}{n}\sum_{i=1}^{n} f_1(Z_i;\hat{\theta}) - \frac{1}{n}\sum_{i=1}^{n} f_2(Z_i;\mu_0)[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]$$

$$- \frac{1}{n}\sum_{i=1}^{n} [f_2(Z_i;\hat{\theta}) - f_2(Z_i;\mu_0)] \pounds [\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]$$

$$+ \frac{1}{n}\sum_{i=1}^{n} \frac{r^0[g(x;\hat{\theta}) - \hat{\theta}_i;x][Y_i - 1(\hat{\theta}_i > 0)]}{\tilde{A}^2(U_i)\hat{A}(\hat{\theta}_i)}[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]^2:$$

**Leading Terms.** We make use of Lemma 3 given below. The term $n^{-1}\sum_{i=1}^{n} f_1(Z_i;\hat{\theta})$ has already been analyzed above. By Lemma 3 we have with probability tending to one

$$\left[\frac{1}{n}\sum_{i=1}^{n} f_2(Z_i;\mu_0) 4[\hat{A}(\hat{\theta}_i) - \tilde{A}(U_i)]\right]^2 - \frac{1}{n}\sum_{j=1}^{n} L^{\pi}(Z_i;Z_j) 5 \cdot \left(\frac{\hat{A}}{nb^3}\frac{1}{n}\sum_{i=1}^{n} jf_2(Z_i;\mu_0)jd(X_i)\right)$$

$$= O_p(n^{-1}b^{-3}) \tag{9}$$

for some function $d(:)$ and random variables $L^{\pi}(Z_i;Z_j) = b^{-1}k((U_i - U_j)=b) - \tilde{A}(U_i) + {}^{\pi}(Z_i)\pounds m(Z_j;\mu_0);$ where

$$-{}^{\pi}(Z_i) = \tilde{A}^0(U_i)\left[\frac{@g}{@\mu^0}(X_i;\mu_0) - E\left[\frac{@g}{@\mu^0}(X_i;\mu_0)jU_i\right]\right] - \tilde{A}(U_i)\hat{A}^0_u(U_i):$$

Under our bandwidth conditions, the right hand side of (9) is $o_p(n^{-1=2}):$ Furthermore,

$$\frac{1}{n}\sum_{i=1}^{n} f_2(Z_i;\mu_0)\frac{1}{n}\sum_{j=1}^{n} L^{\pi}(Z_i;Z_j) = \sum_{i=1}^{n}\sum_{j=1}^{n} \acute{}_n(Z_i;Z_j)$$

where

$$\acute{}_n(Z_i;Z_j) = \frac{1}{n^2}f_2(Z_i;\mu_0)\left[\frac{1}{b}k\left(\frac{U_i - U_j}{b}\right) - \tilde{A}(U_i) + {}^{\pi}(Z_i)\pounds m(Z_j;\mu_0)\right]:$$

Note that $E_i\acute{}_n(Z_i;Z_j) = 0$ but $E_j\acute{}_n(Z_i;Z_j) \neq 0:$ We write

$$T_{n4} = \frac{1}{n}\sum_{j=1}^{n} \tilde{A}(Z_j) + \sum_{\substack{i=1 \\ i\neq j}}^{n}\sum_{j=1}^{n} \acute{}_n(Z_i;Z_j);$$

22

where

$$\tilde{A}(Z_j) = n^2 E_j \,'_n(Z_i; Z_j)$$

$$!_n(Z_i; Z_j) = \,'_n(Z_i; Z_j) \; ¡ \; E_j \,'_n(Z_i; Z_j);$$

so that $!_n(Z_i; Z_j)$ is a degenerate kernel satisfying $E_i!_n(Z_i; Z_j) = E_j!_n(Z_i; Z_j) = 0$: We next compute $\tilde{A}_n(Z_j)$; using integration by parts, a change of variable, and dominated convergence we have

$$\tilde{A}(Z_j) = (f_2(Z_j; \mu_0) \; ¡ \; E[f_2(Z_j; \mu_0)jU_j]) \tilde{A}(U_j) + E[f_2(Z_i; \mu_0) ¡ \,^\pi(Z_i)] ¢ m(Z_j; \mu_0) + O_p(b^2):$$

Finally,

$$(f_2(Z_i; \mu_0) \; ¡ \; E[f_2(Z_i; \mu_0)jU_i]) \tilde{A}(U_i) = \frac{r^0[g(x; \mu_0) \; ¡ \; U_i; x]}{\tilde{A}(U_i)} (Y_i \; ¡ \; E[Y_ijU_i])$$

$$E[f_2(Z_i; \mu_0) ¡ \,^\pi(Z_i)] = E\left[ \frac{r^0[g(x; \mu_0) \; ¡ \; U_i; x]}{\tilde{A}(U_i)}[Y_i \; ¡ \; 1(U_i > 0)] £ \left( \frac{\tilde{A}^0(U_i)}{\tilde{A}(U_i)} \,^\pi_i \; ¡ \; Á^0_u(U_i) \right) \right]:$$

**Remainder Terms.** First,

$$\left| \frac{1}{n} \sum_{i=1}^{n} [f_2(Z_i; \hat{\beta}) \; ¡ \; f_2(Z_i; \mu_0)][\hat{A}(\hat{\theta}_i) \; ¡ \; \tilde{A}(U_i)] \right|$$

$$\cdot \left( \frac{1}{n} \sum_{i=1}^{n} [f_2(Z_i; \hat{\beta}) \; ¡ \; f_2(Z_i; \mu_0)]^2 \right)^{1=2} \left( \frac{1}{n} \sum_{i=1}^{n} [\hat{A}(\hat{\theta}_i) \; ¡ \; \tilde{A}(U_i)]^2 \right)^{1=2}$$

$$= o_p(n^{¡1=2}):$$

Second

$$\left| \frac{1}{n} \sum_{i=1}^{n} \frac{r^0[g(x; \hat{\beta}) \; ¡ \; \hat{\theta}_i; x][Y_i \; ¡ \; 1(\hat{\theta}_i > 0)]}{\tilde{A}^2(U_i)\hat{A}(\hat{\theta}_i)}[\hat{A}(\hat{\theta}_i) \; ¡ \; \tilde{A}(U_i)]^2 \right|$$

$$\cdot \frac{\sup_{u2U}[\hat{A}(u) \; ¡ \; \tilde{A}(u)]^2(1 + o_p(1))}{\inf_{u2U} \tilde{A}^2(u)\hat{A}(u) + o_p(1)} \frac{1}{n} \sum_{i=1}^{n} jr^0[g(x; \hat{\beta}) \; ¡ \; \hat{\theta}_i; x]j ¢ (jY_ij + 1)$$

$$= o_p(n^{¡1=2}): \qquad \blacksquare$$

## 8.4 Subsidiary Results

De...ne

$$F_n(\mu) = \frac{1}{n} \sum_{i=1}^{n} f(Z_i; \mu)$$

23

for some function $f$; and let $F(\mu) = EF_n(\mu)$ and $\Lambda_F = \partial F(\mu_0)/\partial\mu$:

**Lemma 1**. Assume:

(i) For some vector $m$

$$\sqrt{n}(\hat\beta - \mu_0) = \frac{1}{\sqrt{n}}\sum_{i=1}^{\infty} m(Z_i;\mu_0) + o_p(1)$$

where $E[m(Z_i;\mu_0)] = 0$ and $\Omega = E[m(Z_i;\mu_0)m(Z_i;\mu_0)'] < \infty$ :

(ii) There exists a finite matrix $\Lambda_F$ of full (column) rank such that

$$\lim_{\|\mu - \mu_0\|\to 0} \frac{\|F(\mu) - \Lambda_F(\mu - \mu_0)\|}{\|\mu - \mu_0\|} = 0:$$

(iii) For every sequence of positive numbers $\{\pm_n\}$ such that $\pm_n \to 0$;

$$\sup_{\|\mu - \mu_0\|\cdot \pm_n} \left\{\sqrt{n}[F_n(\mu) - F(\mu)] - \sqrt{n}[F_n(\mu_0) - F(\mu_0)]\right\} = o_p(1):$$

Then

$$\sqrt{n}[F_n(\hat\beta) - F(\mu_0)] \Rightarrow N(0;V);$$

where

$$V = \text{var}[\Lambda_F m(Z_i;\mu_0) + f(Z_i;\mu_0)]$$

$$= \Lambda_F \Omega \Lambda_F' + \text{var}[f(Z_i;\mu_0)] + 2\Lambda_F Em(Z_i;\mu_0)f(Z_i;\mu_0):$$

See below for a discussion on the verification of (iii).

**Proof**: Since $\hat\beta$ is root-n consistent, there exists a sequence $\pm_n \to 0$ such that

$$\Pr[\|\sqrt{n}(\hat\beta - \mu_0)\| > \pm_n] \to 0$$

as $n \to \infty$ : We can therefore suppose that $\|\sqrt{n}(\hat\beta - \mu_0)\| \cdot \pm_n$ with probability tending to one. We have

$$\sqrt{n}[F_n(\hat\beta) - F(\mu_0)] = \sqrt{n}[F(\hat\beta) - F(\mu_0)] + \sqrt{n}[F_n(\hat\beta) - F(\hat\beta)]$$

$$= \Lambda_F \sqrt{n}(\hat\beta - \mu_0) + \sqrt{n}[F_n(\mu_0) - F(\mu_0)] + o(\|\sqrt{n}(\hat\beta - \mu_0)\|)$$
$$+ \sqrt{n}\{[F_n(\hat\beta) - F(\hat\beta)] - [F_n(\mu_0) - F(\mu_0)]\}$$

$$= \Lambda_F \sqrt{n}(\hat\beta - \mu_0) + \sqrt{n}[F_n(\mu_0) - F(\mu_0)] + o_p(1) \text{ [by (ii) and (iii)]}$$

$$= \frac{1}{\sqrt{n}}\sum_{i=1}^{\infty} \{\Lambda_F m(Z_i;\mu_0) + [f(Z_i;\mu_0) - Ef(Z_i;\mu_0)]\} + o_p(1);$$

24

and the result now follows from standard CLT arguments. ∎

Lemma 2. As $n \to \infty$

$$\max_{1 \le i \le n} \left| \hat{A}(\hat{\theta}_i) - \tilde{A}(U_i) - \frac{1}{n}\sum_{j=1}^{n} L(Z_i; Z_j) \right| = o_p(n^{-1/2}); \qquad (10)$$

where $L(Z_i; Z_j) = \varkappa_j(U_i) + \iota(Z_i) m(Z_j; \mu_0)$ and

$$\varkappa_j(u) = h[g(X_j; \mu_0) - u] - E(h[g(X_j; \mu_0) - u])$$

$$\iota(Z_i) = \tilde{A}'(U_i) \left\{ \frac{@g}{@\mu}(X_i; \mu_0) - E\left[ \frac{@g}{@\mu}(X_i; \mu_0) \right] \right\}:$$

**Proof.** We have for any $u$;

$$\hat{A}(u) - \tilde{A}(u) = \frac{1}{n}\sum_{j=1}^{n} h[g(X_j; \hat{\theta}) - u] - E(h[g(X_j; \mu_0) - u])$$

$$= \frac{1}{n}\sum_{j=1}^{n} h[g(X_j; \mu_0) - u] - E(h[g(X_j; \mu_0) - u])$$

$$+ \frac{1}{n}\sum_{j=1}^{n} h'[g(X_j; \mu_0) - u]\frac{@g}{@\mu}(X_j; \mu_0)(\hat{\theta} - \mu_0) + R_n(u);$$

where

$$R_n(u) = \frac{1}{2n}\sum_{j=1}^{n} h''[g(X_j; \bar{\mu}) - u](\hat{\theta} - \mu_0)'\frac{@g}{@\mu}(X_j; \bar{\mu})\frac{@g}{@\mu'}(X_j; \bar{\mu})(\hat{\theta} - \mu_0)$$

$$+ \frac{1}{2n}\sum_{j=1}^{n} h'[g(X_j; \bar{\mu}) - u](\hat{\theta} - \mu_0)'\frac{@^2g}{@\mu@\mu'}(X_j; \bar{\mu})(\hat{\theta} - \mu_0);$$

where $\bar{\mu}$ are intermediate values between $\hat{\theta}$ and $\mu_0$. With probability tending to one for a sequence $\pm_n \to 0$ we have by the Cauchy Schwarz inequality

$$|R_n(u)| \le ||\hat{\theta} - \mu_0||^2 \frac{1}{2n}\sum_{j=1}^{n} \sup_{\|\mu - \mu_0\| \le \pm_n} |h''[g(X_j; \mu) - u]| \sup_{\|\mu - \mu_0\| \le \pm_n} \left\|\frac{@g}{@\mu}(X_j; \mu)\right\|$$

$$+ ||\hat{\theta} - \mu_0||^2 \frac{1}{2n}\sum_{j=1}^{n} \sup_{\|\mu - \mu_0\| \le \pm_n} |h'[g(X_j; \mu) - u]| \sup_{\|\mu - \mu_0\| \le \pm_n} \left\|\frac{@^2g}{@\mu@\mu'}(X_j; \mu)\right\|$$

$$\le ||\hat{\theta} - \mu_0||^2 \left( \sup_t |h''(t)| \frac{1}{2n}\sum_{j=1}^{n} d_1(X_j) + \sup_t |h'(t)| \frac{1}{2n}\sum_{j=1}^{n} d_2(X_j) \right) \qquad (11)$$

$$= O_p(n^{-1}):$$

Since the right hand side of (11) does not depend on $u$; this order is uniform in $u$: Furthermore because $h'$ is bounded and continuous, by a standard uniform law of large numbers

$$\sup_{u \in U} \left| \frac{1}{n}\sum_{j=1}^{n} h'[g(X_j; \mu_0) - u]\frac{@g}{@\mu}(X_j; \mu_0) - E\left[ h'[g(X_j; \mu_0) - u]\frac{@g}{@\mu}(X_j; \mu_0) \right] \right| = o_p(1);$$

25

where $U$ is the support of $U_i = g(X_i; \mu_0) - V_i$. Therefore,

$$\sup_{u \in U} \left| \hat{A}(u) - \tilde{A}(u) - \frac{1}{n}\sum_{j=1}^{n} \varkappa_j(u) - J(u)(\hat{\beta} - \mu_0) \right| = o_p(n^{-1/2}); \qquad (12)$$

where $\varkappa_j(u) = h[g(X_j; \mu_0) - u] - E(h[g(X_j; \mu_0) - u])$ are i.i.d. with mean zero and finite variance, and

$$J(u) = E\left[ h'[g(X_j; \mu_0) - u]\frac{\partial g}{\partial \mu}(X_j; \mu_0) \right].$$

Because $\sqrt{n}(\hat{\beta} - \mu_0) = O_p(1)$ the supremum over $u \in U$ is the same as a maximum over $\hat{U}_i$.

Furthermore, by a second order Taylor series expansion

$$
\begin{aligned}
\frac{1}{n}\sum_{j=1}^{n} \varkappa_j(\hat{U}_i) + J(\hat{U}_i)(\hat{\beta} - \mu_0) &= \frac{1}{n}\sum_{j=1}^{n} \varkappa_j(U_i) + \frac{1}{n}\sum_{j=1}^{n} \frac{\partial \varkappa_j}{\partial u}(U_i)\frac{\partial g}{\partial \mu}(X_j; \mu_0)(\hat{\beta} - \mu_0) \\
&\quad + J(U_i)(\hat{\beta} - \mu_0) + o_p(n^{-1/2}) \\
&= \frac{1}{n}\sum_{j=1}^{n} \varkappa_j(U_i) + E_i\left[ \frac{\partial \varkappa_j}{\partial u}(U_i)\frac{\partial g}{\partial \mu}(X_j; \mu_0) \right](\hat{\beta} - \mu_0) \\
&\quad + J(U_i)(\hat{\beta} - \mu_0) + o_p(n^{-1/2});
\end{aligned}
$$

where

$$\frac{\partial \varkappa_j}{\partial u}(u) = -h'[g(X_j; \mu_0) - u] + E(h'[g(X_j; \mu_0) - u]);$$

and the error term is bounded in the same way as above using the continuous second derivatives of $h, g$. That is, $\max_{1 \le i \le n} |J(\hat{U}_i) - J(U_i)| = O_p(1)$ and

$$\sup_{||\mu - \mu_0|| \le \eta_n} \max_{1 \le i \le n} \left| \frac{\partial^2 \varkappa_j}{\partial u^2}(U_i(\mu)) \right| \cdot d(Z_i)$$

with $E d(Z_i) < 1$. Note that

$$
\begin{aligned}
J(U_i) + E_i\left[ \frac{\partial \varkappa_j}{\partial u}(U_i)\frac{\partial g}{\partial \mu}(X_j; \mu_0) \right] &= E_i\left[ h'[g(X_j; \mu_0) - U_i]\frac{\partial g}{\partial \mu}(X_j; \mu_0) \right] \\
&\quad + E_i\left[ -h'[g(X_j; \mu_0) - U_i] + E_i(h'[g(X_j; \mu_0) - U_i])g\frac{\partial g}{\partial \mu}(X_j; \mu_0) \right] \\
&= E_i(h'[g(X_j; \mu_0) - U_i]) E\left( \frac{\partial g}{\partial \mu}(X_j; \mu_0) \right);
\end{aligned}
$$

so that

$$\hat{A}(\hat{U}_i) - \tilde{A}(\hat{U}_i) = \frac{1}{n}\sum_{j=1}^{n} \varkappa_j(U_i) + \left( E_i[h'[g(X_j; \mu_0) - U_i]] \cdot E\left( \frac{\partial g}{\partial \mu}(X_j; \mu_0) \right) \right) \cdot (\hat{\beta} - \mu_0) + o_p(n^{-1/2}). \qquad (13)$$

26

Finally,

$$\bar{A}(\hat{\theta}_i) - \bar{A}(U_i) = -E_i(h'[g(X_j;\mu_0) - U_i]) \cdot \frac{\partial g}{\partial \mu}(X_i;\mu_0) \cdot (\hat{\beta} - \mu_0) + o_p(n^{-1/2}): \qquad (14)$$

These results are uniform under some additional conditions. Combining (13) and (14) we obtain the result (10).

∎

Lemma 3. We have with probability tending to one

$$\left| \bar{A}(\hat{\theta}_i) - \bar{A}(U_i) - \frac{1}{n}\sum_{j=1}^{X} L^{\pi}(Z_i;Z_j) \right| \cdot \frac{k}{nb^3} d(X_i)$$

for some function $d$; where

$$L^{\pi}(Z_i;Z_j) = \frac{1}{b}k\left(\frac{U_i - U_j}{b}\right) - \bar{A}(U_i) + i^{\pi}(Z_i) \cdot m(Z_j;\mu_0)$$

$$i^{\pi}(Z_i) = \bar{A}'(U_i)\left[\frac{\partial g}{\partial \mu'}(X_i;\mu_0) - E\left[\frac{\partial g}{\partial \mu'}(X_i;\mu_0)|U_i\right]\right] - \bar{A}(U_i)\bar{A}_u'(U_i):$$

Proof. Making a second order Taylor series expansion we have $\bar{A}(\hat{\theta}_i) - \bar{A}(U_i) = T_{ni} + R_{ni}$; where

$$T_{ni} = \bar{A}(U_i) - \bar{A}(U_i) + \frac{1}{nb^2}\sum_{j=1}^{X} k'\left(\frac{U_i - U_j}{b}\right)\left[\frac{\partial g}{\partial \mu'}(X_i;\mu_0) - \frac{\partial g}{\partial \mu'}(X_j;\mu_0)\right](\hat{\beta} - \mu_0)$$

$$R_{ni} = \frac{1}{nb^3}\sum_{j=1}^{X} k''\left(\frac{U_i^{\pi} - U_j^{\pi}}{b}\right)\left[\frac{\partial g}{\partial \mu}(X_i;\mu_0) - \frac{\partial g}{\partial \mu}(X_j;\mu_0)\right](\hat{\beta} - \mu_0)(\hat{\beta} - \mu_0)'\left[\frac{\partial g}{\partial \mu}(X_i;\mu_0) - \frac{\partial g}{\partial \mu}(X_j;\mu_0)\right]'$$

$$+\frac{1}{nb^2}\sum_{j=1}^{X} k'\left(\frac{U_i - U_j}{b}\right)(\hat{\beta} - \mu_0)'\left[\frac{\partial^2 g}{\partial\mu\partial\mu'}(X_i;\mu^{\pi}) - \frac{\partial^2 g}{\partial\mu\partial\mu'}(X_j;\mu^{\pi})\right](\hat{\beta} - \mu_0);$$

where $\mu^{\pi}$ are intermediate values between $\hat{\beta}$ and $\mu_0$; and $U_i^{\pi} = U_i(\mu^{\pi})$: We have with probability tending to one

$$|R_{ni}| \cdot b^{-3}\sup_u|k''(u)| \cdot ||\hat{\beta} - \mu_0||^2 \cdot \left\{\left|\frac{\partial g}{\partial\mu}(X_i;\mu_0)\right|^2 + \frac{1}{n}\sum_{j=1}^{X}\left|\frac{\partial g}{\partial\mu}(X_j;\mu_0)\right|^2\right\}$$

$$+b^{-1}||\hat{\beta} - \mu_0||^2 \cdot \frac{1}{nb}\sum_{j=1}^{X} k'\left(\frac{U_i - U_j}{b}\right)(d_1(X_i) + d_2(X_j))$$

by the Cauchy-Schwarz inequality: By the uniform convergence results of Masry (1996a,1996b):

$$\max_{1 \cdot i \cdot n} \frac{1}{nb}\sum_{j=1}^{X} k'\left(\frac{U_i - U_j}{b}\right) = O_p(1)$$

$$\max_{1 \cdot i \cdot n} \frac{1}{nb}\sum_{j=1}^{X} k'\left(\frac{U_i - U_j}{b}\right)d_2(X_j) = O_p(1);$$

27

so that for suitable constants and dominating functions

$$|R_{ni}| \leq \frac{k_1}{nb^3}(d_3(X_i) + k_2) + \frac{k_3}{nb}(d_1(X_i) + k_4)$$

with probability tending to one. This gives the result. Furthermore, if the functions $d_j$ are bounded this translates into a uniform result

$$\max_{1 \leq i \leq n} |R_{ni}| = O_p(n^{-1}b^{-3}).$$

Provided $nb^6 \to \infty$; this term is $o_p(n^{-1/2})$. With additional smoothness conditions on $k$ this condition can be substantially weakened.

Furthermore, by Masry (1996a, 1996b)

$$\max_{1 \leq i \leq n} \left| \frac{1}{nb^2} \sum_{j=1}^{n} k'\left(\frac{U_i - U_j}{b}\right) - E\left[\frac{1}{b^2}k'\left(\frac{U_i - U_j}{b}\right)\Big|U_i\right] \right| = O_p\left(\sqrt{\frac{\log n}{nb^3}}\right)$$

$$\max_{1 \leq i \leq n} \left| \frac{1}{nb^2} \sum_{j=1}^{n} k'\left(\frac{U_i - U_j}{b}\right)\frac{\partial g}{\partial \mu}(X_j; \mu_0) - E\left[\frac{1}{b^2}k'\left(\frac{U_i - U_j}{b}\right)\hat{A}_u(U_j)\Big|U_i\right] \right| = O_p\left(\sqrt{\frac{\log n}{nb^3}}\right).$$

Also,

$$E\left[\frac{1}{b^2}k'\left(\frac{U_i - U_j}{b}\right)\hat{A}_u(U_j)\Big|U_i\right] - [\hat{A}_u(U_i)\tilde{A}(U_i)]'$$

$$= \int \frac{1}{b^2}k'\left(\frac{U_i - u}{b}\right)\hat{A}_u(u)\tilde{A}(u)du - [\hat{A}_u(U_i)\tilde{A}(U_i)]'$$

$$= \int \frac{1}{b}k\left(\frac{U_i - u}{b}\right)[\hat{A}_u(u)\tilde{A}(u)]'du - [\hat{A}_u(U_i)\tilde{A}(U_i)]'$$

$$= \int k(t) \left([\hat{A}_u(U_i + tb)\tilde{A}(U_i + tb)]' - [\hat{A}_u(U_i)\tilde{A}(U_i)]'\right) dt$$

$$= O_p(b^2)$$

by integration by parts, change of variables and dominated convergence using the symmetry of $k$. This order is uniform in $i$ by virtue of the boundedness and continuity of the relevant functions. Therefore,

$$\max_{1 \leq i \leq n} \left| T_{ni} - \frac{1}{n}\sum_{j=1}^{n} L^*(Z_i; Z_j) \right| = o_p(n^{-1/2}).$$

Finally, we have

$$\max_{1 \leq i \leq n} \left| \frac{1}{n}\sum_{j=1}^{n} L^*(Z_i; Z_j) \right| = O_p\left(b^2 + \sqrt{\frac{\log n}{nb}}\right)$$

by standard results for kernel estimates. ∎

28

## 8.4.1 Stochastic Equicontinuity Results

We now show that condition (iii) of Lemma 1 is satis...ed. Let $£_n(c) = f\mu: \sqrt{n}^{\bar{}}|\mu_i \;\mu^{0}_{\bar{}}\cdot\;cg$: Since $\sqrt{n}(\hat{\beta}_i\mu^0) = O_p(1)$; for all $2 > 0$ there exists a $c_2$ and an integer $n_0$ such that for all $n \;, n_0$; $\Pr[\hat{\beta} 2 £_n(c_2)] \; 1_i \; 2$: De...ne the stochastic process

$$°_n(\mu) = \frac{1}{\sqrt{n}}\sum_{i=1}^{n} f(Z_i;\mu)_i\; E[f(Z_i;\mu)]; \quad \mu \; 2 \; £;$$

where

$$f(Z_i;\mu) = r[g(x;\mu);x] + \frac{r^0[g(x;\mu)_i\; U_i(\mu);x][Y_i\;_i\;1(U_i(\mu) > 0)]}{\tilde{A}(U_i)}$$

and de...ne the pseudo-metric

$$\frac{1}{2}(\mu;\mu^0) = E^3\;[f(Z_i;\mu)_i\;f(Z_i;\mu^0)]^{2^{'}};$$

on $£$: Under this metric, the parameter space $£$ is totally bounded. We are only interested in the behaviour of this process as $\mu$ varies in the small set $£_n$: By writing $\mu = \mu^0 + °n_i^{1=2}$; we shall make a reparameterization to $°_n(°)$; where $° \; 2 \;_i(c)\; \frac{1}{2}\;R^p$: We establish the following result:

$$\sup_{°2_i} j°_n(°)_i\;°_n(0)j = o_p(1) \tag{15}$$

To prove (15) it is su¢cient to show a pointwise law of large numbers, e.g., $°_n(°)_i\;°_n(0) = o_p(1)$ for any $° \; 2 \;_i$; and stochastic equicontinuity of the process $°_n$ at $° = 0$. The pointwise result is immediate because the random variables are sums of i.i.d. random variables with ...nite absolute moment and zero mean; the probability limit of $°_n(°)$ is the same for all $° \; 2 \;_i$ by the smoothness of the expected value in $°$. To complete the proof of (15) we shall use the following lemma, proved below, which states that $°_n$ is stochastically equicontinuous in $\mu$. The di¢culty in establishing the required equicontinuity arises solely because the function $g$ inside $U$ is nonlinear in $\mu$:

**Lemma SE.** Under the above assumptions, the process $°_n(°)$ is stochastically equicontinuous, i.e., for all $2 > 0$ and $´ > 0$; there exists $\pm > 0$ such that

$$\limsup_{n!\;1} \Pr\left[\sup_{\frac{1}{2}(t_1;t_2)<\pm} j°_n(t_1)_i\;°_n(t_2)j > ´\right] < 2:$$

**Proof of Lemma SE.** By a second order Taylor series expansion of $g(Z_i;\mu)$ around $g(Z_i;\mu^0)$:

$$g(Z_i;\mu^0 + °n_i^{1=2}) = g(Z_i;\mu^0) + \frac{1}{\sqrt{n}}\sum_{k=1}^{p}\frac{@g}{@\mu_k}(Z_i;\mu^0)°_k + \frac{1}{n}\sum_{k=1}^{p}\sum_{r=1}^{p}\frac{@^2g}{@\mu_k@\mu_r}(Z_i;\bar{\mu})°_k°_r \tag{16}$$

29

for some intermediate points $\bar{\mu}$: De...ne the linear approximation to $g(Z_i; \mu^0 + {}^{\circ}n^{i\ 1=2})$;

$$T(Z_i; {}^{\circ}) = g(Z_i; \mu^0) + \sum_{k=1} \frac{@g}{@\mu_k}(Z_i; \mu^0){}^{\circ}{}_k$$

for any ${}^{\circ}$: By assumption C2, for all $k; r$; $\sup_{\mu 2 \mathsterling} j@^2 g(Z_i; \mu)=@\mu_k@\mu_r j^2 \cdot d(Z_i)$ with $Ed(Z_i) < 1$: Therefore, for all $\pm > 0$ there exists an $" > 0$ such that

$$\Pr\left["\frac{1}{n}\max_{i;k;r}\sup_{\mu 2 \mathsterling_n}\left\vert\frac{@^2 g}{@\mu_k@\mu_r}(Z_i; \mu)\right\vert > "\right\#] \cdot n\sum_{k;r}\Pr\left["\frac{1}{n}\sup_{\mu 2 \mathsterling_n}\left\vert\frac{@^2 g}{@\mu_k@\mu_r}(Z_i; \mu)\right\vert > "\right\#]$$

$$\cdot \frac{\sum_{k;r} E[d(Z_i)]}{"2}$$

$$\cdot \pm$$

by the Bonferroni and Chebychev inequalities. Therefore, with probability tending to one

$$\max_{1\cdot i \cdot n}\left\vert\frac{1}{n}\sum_{k=1}\sum_{r=1}\frac{@^2 g}{@\mu_k@\mu_r}(Z_i; \bar{\mu}){}^{\circ}{}_k {}^{\circ}{}_r\right\vert \cdot \frac{\frac{1}{4}}{\sqrt{n}}$$

for some $\frac{1}{4} < 1$: De...ne the stochastic process

$$^{\circ}{}_{n1}({}^{\circ}; \frac{1}{4}) = \frac{1}{\sqrt{n}}\sum_{i=1}\bar{f}(Z_i; \mu_0 + {}^{\circ}n^{i\ 1=2}; \frac{1}{4}n^{i\ 1=2}) \text{ i } E\bar{f}(Z_i; \mu_0 + {}^{\circ}n^{i\ 1=2}; \frac{1}{4}n^{i\ 1=2})$$

on ${}^{\circ} 2 \text{ i and } \frac{1}{4} 2 \text{ ¦} = [0; \frac{1}{4}]$; where

$$\bar{f}(Z_i; \mu_0 + {}^{\circ}n^{i\ 1=2}; \frac{1}{4}n^{i\ 1=2})$$

$$= r[g(x; \mu_0 + {}^{\circ}n^{i\ 1=2}); x] + \frac{r^0[g(x; \mu_0 + {}^{\circ}n^{i\ 1=2}) \text{ i } U_i(\mu_0 + {}^{\circ}n^{i\ 1=2}); x]}{\tilde{A}(U_i)}[Y_i \text{ i } 1(T(Z_i; {}^{\circ}n^{i\ 1=2}) + \frac{\frac{1}{4}}{\sqrt{n}} > 0)]$$

It su¢ces to show that $^{\circ}{}_{n1}({}^{\circ}; \frac{1}{4})$ is stochastically equicontinuous in ${}^{\circ}; \frac{1}{4}$; and the deterministic centering term is of smaller order. The latter argument is a standard Taylor expansion. The argument for $^{\circ}{}_{n1}({}^{\circ}; \frac{1}{4})$ is very similar to that contained in Sherman (1993) because we basically have a linear index structure in this part. One can apply Lemma 2.12 in Pakes and Pollard (1989). ■

# References

[1] Andrews, D.W.K. (1994), Asymptotics for Semiparametric Econometric Models by Stochastic Equicontinuity. Econometrica 62, 43-72.

[2] Bickel, P.J., C.A.J. Klaassen, J. Ritov, and J. Wellner (1993), E¢cient and Adaptive Estimation for Semiparametric Models. Springer: Berlin.

[3] Gozalo, P., and O.B. Linton (1999), Local Nonlinear Least Squares: Using Parametric Information in Nonparametric Regression. Journal of Econometrics, 99, 63-106.

[4] Hardle, W., and O.B. Linton (1994), "Applied Nonparametric Methods," The Handbook of Econometrics, vol. IV, eds. D.F. McFadden and R.F. Engle III. North Holland.

[5] Klein, R. and R. H. Spady (1993), "An e¢cient Semiparametric Estimator for Binary Response Models," Econometrica 61, 387-421.

[6] Lewbel, A. (1997), "Semiparametric Estimation of Location and Other Discrete Choice Moments," Econometric Theory, 13, 32-51.

[7] Lewbel, A. (2000), "Semiparametric Qualitative Response Model Estimation With Unknown Heteroscedasticity or Instrumental Variables," Journal of Econometrics 97, 145-177.

[8] Linton, O. and J.P. Nielsen (1995), "A kernel method of estimating structured nonparametric regression based on marginal integration," Biometrika, 82, 93-100.

[9] Masry, E. (1996a), "Multivariate local polynomial regression for time series: Uniform strong consistency and rates," J. Time Ser. Anal. 17, 571-599.

[10] Masry, E., (1996b), "Multivariate regression estimation: Local polynomial ...tting for time series. Stochastic Processes and their Applications 65, 81-101.

[11] McFadden, D. (1994), "Contingent Valuation and Social Choice," American Journal of Agricultural Economics, 76, 4.

[12] Newey, W. K. (1994), "The Asymptotic Variance of Semiparametric Estimators," Econometrica, 62, 1349–1382.

[13] Newey, W. K. and D. McFadden (1994), "Large Sample Estimation and Hypothesis Testing," in Handbook of Econometrics, vol. iv, ed. by R. F. Engle and D. L. McFadden, pp. 2111-2245, Amsterdam: Elsevier.

[14] Pakes, A. and D. Pollard. (1989), "Simulation and the Asymptotics of Optimization Estimators," Econometrica, 57, 1027-57.

[15] Sherman, R. P. (1993), "The Limiting Distribution of the Maximum Rank Correlation Estimator," Econometrica, 61, 123-37.

[16] **Silverman, B.** (1986): Density estimation for statistics and data analysis. London, Chapman and Hall.

[17] **Stone, C.J. (1982).** Optimal global rates of convergence for nonparametric regression. Annals of Statistics 8, 1040-1053.

[18] **Tjostheim, D. and B. H. Auestad** (1994), "Nonparametric Identi...cation of Nonlinear Time Series: Projections," Journal of the American Statistical Association, 89, 1398-1409.