# ON INTERCEPT ESTIMATION IN THE
# SAMPLE SELECTION MODEL[*]

Marcia M A Schafgans
Department of Economics, London School of Economics

and

Victoria Zinde-Walsh
Department of Economics, McGill University

Contents:

**Abstract**

We provide a proof of the consistency and asymptotic normality of the estimator suggested by Heckman (1990) for the intercept of a semiparametrically estimated sample selection model. The estimator is based on "identification at infinity" which leads to non-standard convergence rate. Andrews and Schafgans (1998) derived asymptotic results for a smoothed version of the estimator. We examine the optimal bandwidth selection for the estimators and derive asymptotic MSE rates under a wide class of distributional assumptions. We also provide some comparisons of the estimators and practical guidelines.

**Keywords:** Asymptotic normality; sample selection model; semiparametric estimation.

**JEL Nos.:** C14, C34, C35.

# 1   Introduction

Semiparametric estimation of sample selection models has attracted considerable interest in the last decade. More recently the estimation of the intercept of the semiparametric estimated sample selection model has received due attention, see Heckman (1990) and Andrews and Schafgans (1998).

The discussion around the estimation of the intercept arose, since, with the exception of Gallant and Nychka (1978), all semiparametric estimation approaches to the sample selection model precluded the estimation of the intercept; the intercept was absorbed in the nonparametric sample selection bias correction term. The semi-nonparametric estimator proposed by Gallant and Nychka (1978), however, has a drawback in that although it is consistent, its asymptotic distribution is unknown.

The importance of this intercept is evident, e.g., when using the sample selection model in the evaluation of social programs. Estimation of the intercept allows one to evaluate the net bene…t of a social program, by allowing one to compare the actual outcome of participants with the expected outcome had they chosen not to participate. Empirically, the estimation of the intercept of semiparametrically estimated sample selection models has proven desirable in the estimation of wages. Its estimation allows for a decomposition of the wage-gap between socio-economic groups (e.g., male–female) in order to assess the extent of "discrimination" (Schafgans (1998a) and allows for a discussion of its evolution over time (Buchinsky (1998)).

In Andrews and Schafgans (1998), the …rst consistent and asymptotically normal estimator was derived for the intercept, $^1_0$. Their estimator was based on a suggestion by Heckman (1990) to estimate $^1_0$ using only those observations for which the probability of selection in the truncated or censored sample is close to one and in the limit as $n \, ! \, 1$ is one. The justi…cation of this approach is that the conditional mean of the errors in the outcome equation for the observations having probability of selection close to one is close to zero. Due to the di¢culty in deriving the asymptotic distribution of the Heckman (1990) estimator, arising from the non-di¤erentiability of the indicator function, Andrews and Schafgans (1998) introduced a smooth monotone [0,1]-valued function, s(¢). Since we will make reference to the Andrews and Schafgans (1998) paper frequently, we will call it AS henceforth.

In this paper, we derive the consistency and asymptotic distribution of the Heck-

man estimator itself. This provides the empirical researcher with the advantage of not having to specify the smoothing function introduced by Andrews and Schafgans. We investigate a wide class of distributional assumptions for the model and derive "optimal" bandwidth parameters and corresponding asymptotic rates for mean squared error (MSE) for the two estimators. Since the solution for the optimal bandwidth may not be practical, we provide simpler bounds on the optimal bandwidth parameter; using a bound may imply preference for either AS or the Heckman estimator.

The remainder of this paper is organized as follows: Section 2 introduces the sample selection model considered and the estimators of Heckman (1990) and AS. The asymptotic normality result for the Heckman estimator is given in Section 3. Section 4 discusses the optimal selection of the bandwidth and the rate of the asymptotic mean squared error. Section 5 concludes. Various appendices follow. Appendix A contains the proof of the asymptotic normality result given in Section 3. Appendix B derives the asymptotic bias and variance for the two estimators under a class of general distributional assumptions and the optimal bandwidth choices given in Section 4.

## 2   Intercept Estimation

The sample selection model that we consider can be written as:

$$
\begin{aligned}
Y_i^* &= {}^1{}_0 + Z_i^0 \mu_0 + U_i \ ; \\
D_i &= 1(X_i^{0-}{}_0 > {}^{"}_i) \ ; \quad \text{and} \\
Y_i &= Y_i^* D_i \quad \text{for} \quad i = 1; ::::; n \ ;
\end{aligned}
\tag{1}
$$

where $(Y_i; D_i; Z_i; X_i)$ are observed random variables. The ...rst equation is the outcome equation and the second equation is the participation equation. For convenience, we set

$$
W_i = X_i^{0-}{}_0 : \tag{2}
$$

The literature on semiparametric estimation of sample selection models gives several root-n consistent and asymptotically normal estimators for the selection parameters, $^-{}_0$ (up to some unknown scale), and the slope parameters of the outcome equation, $\mu_0$. For instance, one could consider: Ichimura (1993), Han (1987), Newey (1988), Robinson (1988), Powell (1989), Powell, Stock, and Stoker (1989), Ichimura

and Lee (1990), Andrews (1991), and Klein and Spady (1993). The existing litera-
ture and AS can deal both with censored samples (as in the model given in (1)), or
truncated samples. In the latter case, that is where $Y_i$ is observed only if $D_i = 1$,
$\beta_0$ and $\gamma_0$ need to be estimated simultaneously using, e.g., Ichimura and Lee (1990).
Regarding the selection parameters $\gamma_0$; furthermore, it should be noted that only the
slope parameters are required in the context of estimating the intercept $\mu_0$. The loss
of identi...cation of the intercept in the selection equation, e.g., when using Ichimura
(1993) or Ichimura and Lee (1990) is innocuous therefore.

A consistent and asymptotically normal estimator for the intercept, $\mu_0$, which
uses these preliminary estimators, was provided by AS. Their estimator, call it the
AS estimator, is given by

$$\hat{\mu}_s = \frac{\sum_{i=1}^{n}(Y_i - Z_i^0\hat{\beta})D_i s(X_i^0\mathbf{b} - \gamma_n)}{\sum_{i=1}^{n}D_i s(X_i^0\mathbf{b} - \gamma_n)};\tag{3}$$

where $s(\cdot)$ is a non-decreasing $[0,1]$-valued function that has three derivatives bounded
over R and for which $s(x) = 0$ for $x \cdot 0$ and $s(x) = 1$ for $x \geq b$ for some $0 < b < 1$
(AS, Assumption 3). The preliminary estimators $(\hat{\beta}; \mathbf{b})$ are root-n consistent estima-
tors of $(\mu_0; \gamma_0)$. The parameter $\gamma_n$ is called the bandwidth or smoothing parameter,
where the bandwidth parameter is chosen such that $\gamma_n \to \infty$ as $n \to \infty$.

The Heckman (1990) estimator, on which the AS estimator was based, is given by

$$\hat{\mu}_I = \frac{\sum_{i=1}^{n}(Y_i - Z_i^0\hat{\beta})D_i 1(X_i^0\mathbf{b} > \gamma_n)}{\sum_{i=1}^{n}D_i 1(X_i^0\mathbf{b} > \gamma_n)}\tag{4}$$

Comparing the two formulae (3) and (4), it is clear that the AS estimator $\hat{\mu}_s$
di¤ers from Heckman's (1990) $\hat{\mu}_I$ only in that it replaces the indicator function $1(\cdot)$
with a smooth function $s(\cdot)$.

Heckman's estimator $\hat{\mu}_I$ is essentially a sample average of the random variables
$U_i + \mu_0$ over a fraction of all observations, since $Y_i - Z_i^0\hat{\beta} \to_p U_i + \mu_0$ as $n \to \infty$ for
all $i \geq 1$. The e¤ective sample size is equal to the number of observations used for
the estimation of $\mu_0$: Since AS introduced a weighting scheme for these observations,
viz., the smooth function $s(\cdot)$; the estimator $\hat{\mu}_s$ is a weighted sample average of the
random variables $U_i + \mu_0$; where observations with $X_i^0\mathbf{b}$ greater than $\gamma_n$ and with $X_i^0\mathbf{b}$
close to the threshold $\gamma_n$ are weighted less than those further away.

3

Estimation using the AS or Heckman estimator involves two choices, that of the bandwidth parameter $°_n$ and that of a function $s(¢)$ (or $1(¢)$). It is clear that the choice of $°_n$ has the most important consequences for the properties of the estimator while the impact of the function $s(¢)$ is small in comparison. This is con...rmed in the analysis of Section 4; nevertheless there are cases when the choice of $s(¢)$ (or $1(¢)$) a¤ects the asymptotic rate of the MSE; results are presented in Section 4.

First, we turn to our asymptotic normality result for the Heckman estimator.

# 3 Asymptotic normality of the Heckman estimator

Here we prove the conjecture made by Andrews and Schafgans that the Heckman estimator also is asymptotically normal. In the unrealistic case where the true $\mu_0$ and $\bar{}_0$ are known, Andrews and Schafgans already showed that the Heckman estimator, $^\wedge_{I;0}$, is asymptotically normal (i.e., in Andrews and Schafgans (1998) the indicator replaces the $s(¢)$ function when the true $\mu_0$ and $\bar{}_0$ are known).[1]

For our purposes, all we need to show now is that

$$\frac{\textbf{p}\overline{nED_i1(W_i > °_n)}}{¾}\, i\, {}^\wedge_I\, i\, {}^\wedge_{I;0}\, {}^¢\, i\, \overset{\textbf{p}}{!}\, 0; \tag{5}$$

where $¾^2 = Var(U_i)$. Essentially, the proof requires us to deal directly with the non-di¤erentiability of the indicator function.

There are di¤erent ways of dealing with asymptotics for non-di¤erentiable functions. Typically assumptions regarding the probability density function are required. This is due to the fact that the expectation of the Dirac $\pm$-function, which is the generalized derivative of the indicator function, equals the value of the p.d.f. at zero. In our case we need to consider the non-di¤erentiable function $®_{in}$ given by

$$®_n(^\Delta; W_i; X_i)\ \text{with}\ ®_n(^\bar{}; W; X) = 1(X^{0\bar{}} > °_n)\, i\, 1(W > °_n); \tag{6}$$

where $W = X^{0}\bar{}_0$. Transform $X_i$ via a linear transformation into the random vector partitioned as $(W_i; \,_i)$: For our purposes, it will be convenient to let $\,_i = X_{i(i\,1)}$; where $X_{i(i\,1)}$ is $X_i$ with the exclusion of its ...rst component.

---

[1] The estimator $^\wedge_{I;0}$ is identical to $^\wedge_I$ with the preliminary estimators replaced by their true values.

We add the following Assumption A to the Assumptions 1–7 of AS.[2]

**Assumption A:** (a) For some $A > 0$ the conditional probability density function $p.d.f._{W_i|J_i}$ exists for all $W > A$ and declines monotonically. The marginal probability density function $p.d.f._W$ is such that for some $d > 0$; $dn^{1-\delta} \to 1$

$$\frac{p.d.f._W(\bar{c}_n - d)}{\Pr(W_i > \bar{c}_n)^{3/4}} = O(1).$$

(b) For any $W_i > A$ the conditional moment $E_{|W=W_i}(\overset{\circ}{X}_{i(i-1)}{}^{\overset{\circ}{3}})$ exists.

Similar to Assumptions 4 and 7 of AS, Assumption A(a) relates to the upper tail behaviour of the selection index $W_i = X_i^{\beta-}{}_0$. It is satisfied if $W_i$ has a Weibull, Pareto or "combined" upper tail. If Assumption 4 of AS is satisfied with $\gg = 0$, then this condition can be replaced by:

$$\frac{p.d.f._W(\bar{c}_n)}{\Pr(W_i > \bar{c}_n)^{3/4}} = O(1).$$

In both cases, the condition is less strong than requiring a bounded hazard function on $W_i$; since $p.d.f._W(\bar{c}_n - d)^{1/4} = o(1)$.

The second part of Assumption A requires the existence of the conditional moment, but does not place any restrictions on its behaviour as a function of $W_i$. This assumption is satisfied, for example, if the unbounded components of $X_i$ have a joint normal or spherical distribution with $W_i$:

The following theorem summarizes our result for the Heckman estimator which satisfies Assumption $3^0$ of AS.

**Theorem 1:** Under Assumptions 1, 2, 4–7 of AS and Assumption A

(a) $$\frac{\sqrt{nE D_i 1(W_i > \bar{c}_n)}}{\sqrt[3]{4}}\left(\hat{\gamma}_1 - \gamma_{10} - \frac{E U_i D_i 1(W_i > \bar{c}_n)}{E D_i 1(W_i > \bar{c}_n)}\right) \overset{d}{\to} N(0; 1)$$

(b) $$\frac{\sqrt{nE D_i 1(W_i > \bar{c}_n)}}{\sqrt[3]{4}}(\hat{\gamma}_1 - \gamma_{10}) \overset{d}{\to} N(0; 1) \text{ iff Assumption 8 of AS holds.}$$

---

[2] Essentially, Assumptions 1 and 2 of AS require existence of moments and independence between $(U_i; \varepsilon_i)$ and $(Z_i; X_i)$; Assumption 3 and $3^0$ deal with the shape of $s(\phi)$; Assumption 4 characterizes the upper tail of $W_i$ in terms of a parameter $0 \cdot \gg < 1/3$ with "fatter" tails if $\gg = 0$; Assumption 5 is root-n consistency and asymptotic normality of $(\hat{\beta}; \hat{b})$; Assumption 6 is $\bar{c}_n \to 1$; where its speed is restricted by Assumption 7 in terms of the tail of $W_i > \bar{c}_n$:

**Proof:** See Appendix A. $\quad 2$

To test hypotheses and construct con…dence intervals for functions of $(\mathbf{1}_0; \mu_0; \bar{\phantom{.}}_0)$, we need a joint asymptotic normality result for $(^{\wedge}_1; \hat{\beta}; \hat{b})$. This result, similar to that in AS, is given by

**Theorem 2:** Under Assumptions 1, 2, 4–8 of AS and Assumption A

$$
\begin{pmatrix} O \\ B \\ @ \end{pmatrix} \frac{\sqrt{nED_i 1(W_i > °_n)}}{\sqrt{n} \bar{\phantom{.}}_{-i} \sum_{1=2} \begin{matrix} \hat{\beta}_i \mu_0 \\ \hat{b}_i \bar{\phantom{.}}_0 \end{matrix}} (^{\wedge}_i i\ \mathbf{1}_0) \begin{pmatrix} \phantom{.} \\ C \\ A \end{pmatrix}_i \overset{d}{\to} N(0; I):
$$

**Proof:** See Appendix A. $\quad 2$

In the following section, we compare the performance of the estimators and provide guidelines for selection of the bandwidth parameter and function $s(\mathbb{¢})$ (or $1(\mathbb{¢})$).

# 4  Bandwidth selection and comparison of the estimators

Here we use the asymptotic MSE as a criterion for bandwidth selection and choice of the estimator. Two characteristics of the model are of importance for these choices: the tail behaviour of the selection index, $W_i$, and the tail behaviour of the function $!(W)$ de…ned below that determines the asymptotic bias of the estimator. Speci…cally, let

$$
!(W) = E_{jW=W} U_i 1("_i > W): \tag{7}
$$

The asymptotic bias (abias) of the estimator $\hat{b}_s$ (or $\hat{b}_l$ for $s(\mathbb{¢}) = 1(\mathbb{¢})$) is given by[3]

$$
abias(\hat{b}_s) = \frac{i\ E[!(W_i)s(W_i\ i\ °_n)]}{[ED_i s(W_i\ i\ °_n)]}: \tag{8}
$$

Under the bivariate normality assumption of $(U_i; "_i)$, $!(W)$ is equal to $\frac{¾_{"u}}{¾_{"}} Á(\frac{W}{¾_{"}})$; where $Á(\mathbb{¢})$ denotes the standard normal density function.

---

[3] Since $E(U_i s(W_i > °_n)) = 0$ by independence of $U_i$ and $W_i$, this follows as $E(U_i D_i s(W_i > °_n)) = i\ E(U_i(1\ i\ D_i)s(W_i > °_n)) = i\ E(U_i 1("_i > W_i)s(W_i > °_n))$: By law of iterated expectations, this equals $i\ E(!(W_i)s(W_i > °_n))$.

There are circumstances when $!(W)$ may be zero for large enough $W$, e.g. if the distribution of $"_i$ has ...nite support, in which case a comparison of the asymptotic variance will determine the optimal estimator. When $!(W)$ di¤ers from zero for large $W$ we need to establish the importance of the asymptotic bias relative to the asymptotic variance to select our estimator, which of course will depend on the behaviour of this conditional expectation. From Theorem 2 in AS and Theorem 1 in this paper, the asymptotic variance (avar) of the estimator $\mathbf{b}_s$ (or $\mathbf{b}_l$ for $s(\mathbb{C}) = 1(\mathbb{C})$) is given by

$$\mathrm{avar}(\mathbf{b}_s) = \frac{\frac{3}{4}^2 E\,[D_i s^2(W_i \,\mathbf{i}\,\,^\circ{}_n)]}{n\,[E\,D_i s(W_i \,\mathbf{i}\,\,^\circ{}_n)]^2}: \tag{9}$$

Proposition 1 shows that there may be a trade-o¤ between asymptotic variance and asymptotic bias depending on the choice of function $s(\mathbb{C})$ (or $1(\mathbb{C})$). Here and below, "$a(x)\,\frac{1}{4}\,b(x)$" is de...ned to mean that $a(x) = b(x)(1 + o(1))$ as $x\,!\,\,1$.

**Proposition 1**: Under Assumptions 1–7 of AS and Assumption A, for a given sequence $^\circ{}_n$

(a) $\mathrm{avar}(\mathbf{b}_s)\,_\mathbf{o}\,\mathrm{avar}(\mathbf{b}_l)$

(b) $\mathrm{jabias}(\mathbf{b}_s)j\,\cdot\,\mathrm{jabias}(\mathbf{b}_l)j$; if $E1(W_i > \,^\circ{}_n) = E1(W_i > \,^\circ{}_n + b)\,\frac{1}{4}\,1$ ($» = 0$ in Assumption 4 of AS ) and $!(W_i)\,_\mathbf{o}\,0\,8W_i > \,^\circ{}_n$ (or $!(W_i)\,\cdot\,0\,8W_i > \,^\circ{}_n$):

**Proof**: See Appendix B.  $\square$

If there is no asymptotic bias, naturally the Heckman estimator is preferred based on asymptotic MSE (and variance).

We can characterize the tail behaviour of the selection index $W_i$ as "fat-tailed" if Assumption 4 of AS is satis...ed with $» = 0$; if Assumption 4 of AS is satis...ed only with $» > 0$, we say that $W_i$ is "thin-tailed". Examples of fat tails of $W_i$ are Pareto upper tails (i.e., $1\,\mathbf{i}\,F(W)\,\frac{1}{4}\,c_W\,W^{\,\mathbf{i}\,\cdot};\,_\mathbf{o} > 0$) or Weibull ($_\mathbf{o}; c$) upper tails (i.e., $1\,\mathbf{i}\,F(W)\,\frac{1}{4}\,c_W\,\exp(\mathbf{i}\,_\mathbf{o}W^c);\,_\mathbf{o} > 0$) with $c\,\cdot\,1$: For Pareto and fat Weibull tails of $W$ the condition in (b) is satis...ed.

Proposition 2 shows that for fat-tailed distributions of $W_i$ the choice of the function $s(\mathbb{C})$ or $1(\mathbb{C})$ does not a¤ect the asymptotic variance; the asymptotic bias if una¤ected if additionally $!(W_i)$ does not go to zero too fast.

**Proposition 2:** Under Assumptions 1–7 of AS and Assumption A,

(a) If $E1(W_i > {}^\circ_n) = E1(W_i > {}^\circ_n + b) \frac{1}{4} 1$ ( $» = 0$ in Assumption 4 of AS) $avar(\hat{b}_s) \frac{1}{4} avar(\hat{b}_l)$:

(b) If additionally to (a), $jE!(W)1(W > {}^\circ_n + b)j = jE!(W)1(W > {}^\circ_n)j \frac{1}{4} 1$ and $!(W_i) \, 0 \, 8W_i > {}^\circ_n$ (or $!(W_i) \cdot 0 \, 8W_i > {}^\circ_n$); then $abias(\hat{b}_s) \frac{1}{4} abias(\hat{b}_l)$:

**Proof:** See Appendix B. 2

Unlike the assumptions encompassed in our Proposition 2, it is frequently assumed that $W_i$ has thinner upper tails, e.g., the normal (Lee (1982)). We next examine the AS and Heckman estimators for a class of models with tails of the selection index $1_i F(W) \frac{1}{4} c_W W^{\circledR} e^{i \cdot W^c}$ and $!(W) \frac{1}{4} c_! W^{\&} e^{i \, 1 W^d}$ where the parameters ${}^\circledR; \, ; c; \&; 1; d$ are such that the functions $1_i F(W)$ and $!(W) ! 0$ as $W ! 1$: This class of models includes the Pareto, Weibull (with $c \cdot$ or $c > 1$) as well as "combined" tails. If $U; "$ and $W$ are jointly normally distributed $d = c = 2; {}^\circledR = i 1; \, = \frac{1}{2 \frac{3}{4} \frac{2}{W}}; 1 = \frac{1}{2 \frac{3}{4} \frac{2}{i}};$ and $\& = 0.$[4] In order to facilitate the derivation of the asymptotic mean squared error for this class of distributions, we restrict our attention to $s(\cent)$ functions satisfying the following assumption

**Assumption S:** Let $s(\cent)$ be a function satisfying Assumption 3 of AS. For some $q$ its derivatives at zero are such that

$$
s^{(i)}(0) = \begin{cases} \gtrless 0 & i < q \\ a_q \ 6 \ 0 & i = q \\ \text{exists} & i = q + 1: \end{cases}
$$

Note that any function that satis...es AS for which the lowest order of non-zero derivative is $q \cdot 2$ satis...es Assumption S as well; it is only functions with two (or more) zero derivatives at 0 that require this additional assumption.

The following proposition provides expressions for the asymptotic variance and asymptotic bias. To simplify the expressions in Proposition 3 we omit the constant

---

[4] In AS it was mistakenly claimed that the normal distribution has a Weibull tail with $c = 2$; in fact its tail is $W^{i \, 1} \exp(i \, \frac{1}{2 \frac{3}{4} \frac{2}{W}} W^2)(1 + o(W^{i \, 1}))$:

factors, if they are present the expressions below for avar acquire $c_W^{i\,1}$ and for abias $c_i$ as a factor. Furthermore, we omit the subscript $n$ on $°$.

**Proposition** 3: Under Assumptions 1, 2, 5–7 of AS, Assumptions A and S

(a) If $1 \ge F(W) \frac{1}{4} W^{®} \exp(i \, _{,}W^c)$

$$\mathrm{avar}(^{*}_s) = \begin{cases} \mathrm{avar}(^{*}_l) = \frac{3}{4}^2 n^{i\ 1°i\ ®} \exp(_{,}°^c) & \text{if } c \le 1 \\ \frac{3}{4}^2 n^{i\ 1} \frac{2q}{q} °i\ ® \exp(_{,}°^c) & \text{if } c > 1; \end{cases}$$

where $\mathrm{avar}(^{*}_l)$ when $c > 1$ obtains for $q = 0$ ( $\frac{i_0^{\text{¢}}}{0} = 1$ ).

(b) If additionally to (a) $!(W) \frac{1}{4} W^{\&} \exp(i\ ^1 W^d)$

$$\mathrm{abias}(^{*}_s) = \begin{cases} \mathrm{abias}(^{*}_l) = i\ ^{°\&} \exp(i\ ^{1°d}) & \text{if } c \le 1; d \le 1; d < c \\ \mathrm{abias}(^{*}_l) = i\ \frac{}{_{,}+1} ^{°\&} \exp(i\ ^{1°c}) & \text{if } c \le 1; d \le 1; c = d^5 \\ \mathrm{abias}(^{*}_l) = i\ \frac{c}{1d} ^{°\&+c_i\ d} \exp(i\ ^{1°d}) & \text{if } c \le 1; d \le 1; d > c^6 \\ i\ \frac{a_q\,_,c}{(^1d)^{q+1}} ^{°\&+(c_i\ d)(1+q)} \exp(i\ ^{1°d}) & \text{if } c \le 1; d > 1^6 \\ i\ \frac{°\&}{_3} \exp(i\ ^{1°d}) & \text{if } c > 1; d < c; \\ i\ \frac{}{_3 \frac{}{_{,}+1}^{q+1}} ^{°\&} \exp(i\ ^{1°c}) & \text{if } c > 1; d = c \\ i\ \frac{c}{1d}^{q+1} ^{°\&+(c_i\ d)(1+q)} \exp(i\ ^{1°d}) & \text{if } c > 1; d > c; \end{cases}$$

where $\mathrm{abias}(^{*}_l)$ obtains for $q = 0$ ($0! = 1$) where it is not de...ned explicitly and $a_0 ´ 1$.

**Proof**: See Appendix B.    2

We see that under our assumptions on $s(x)$, the asymptotic MSE is a¤ected by the choice of function $s(\text{¢})$ via $q$ and the value of the derivative $s^{(q)}(0)$ only. When $\mathrm{avar}(^{*}_s)$ depends on $q$ it is an increasing function of $q$, while if $\mathrm{abias}(^{*}_s)$ depends on $q$, its absolute value declines with $q$.

As an example of Proposition 3, if $U; "$ and $W$ are jointly normally distributed ($c = d = 2$), the asymptotic bias and variance of $^{*}_s$ (including all relevant constant

---

[6] If $c = d = 0$, the constant $\frac{}{_{,}+1}$ becomes $\frac{®}{®+\&}$:

[6] If $c = 0$, then the expression $_{,}c$ in the constant becomes $i\ ^{®}$:

factors) equals

$$\text{abias}(\hat{\lambda}_s) = i \frac{\rho \frac{\frac{3}{4} \cdot U}{2\frac{1}{4}\frac{3}{4} \cdot}}{\rho \frac{\frac{3}{4}}{2\frac{1}{4}}} \left(\frac{\frac{3}{4}\frac{2}{\cdot\cdot}}{\frac{3}{4}\frac{2}{\cdot\cdot} + \frac{3}{4}\frac{2}{W}}\right)^{q+1} \exp(i \frac{1}{2\frac{3}{4}\frac{2}{\cdot\cdot}} \circ^2)$$

$$\text{avar}(\hat{\lambda}_s) = \frac{\frac{3}{4}^2}{n} \frac{\frac{3}{4}}{\frac{3}{4}_W} \circ \mu_{2q} \frac{\P}{q} \exp(\frac{1}{2\frac{3}{4}\frac{2}{W}} \circ^2):$$

The asymptotic bias and variance of $\hat{\lambda}_l$ under the joint normality assumption obtains for $q = 0$:

If all the parameters determining the tail behaviour of $W$ and the function $!(W)$ were known a solution that would provide an optimal $\circ^{\pi} = \text{argmin}(MSE)$ as a function of $n$, $q$ (the $s(\mathbb{C})$ function) and all those parameters could be obtained (at least via a numerical algorithm) from the formulae in Proposition 3. If the asymptotic bias is not present (in which case one would choose the Heckman estimator based on Proposition 1) the bandwidth parameter arising from reducing MSE (or equivalently avar) as $\circ \ !\ 1$ can be presented as $\circ^{\pi} = (\frac{\mu}{2} \ln n)^{1=c}$ if $c \ne 0$ and $\circ^{\pi} = n^{i \ \mu=®}$ if $c = 0$ with $\mu$ close to zero and would result in a MSE proportional to $n^{i \ 1+\mu}$. Proposition 4 deals with situations where an asymptotic bias is present and may be severe. It characterizes the bandwidth parameter $\circ^{\pi}$ and the best possible rate for MSE depending on the relation between the rate of decline in the tail of $W$ and the function $!(W)$; we also provide simple bounds on $\circ^{\pi}$ which bring MSE close to achieving the best possible rate.

**Proposition 4**: Under Assumptions 1, 2, 5–7 of AS and Assumptions A and S, if $1 \ i \ F(W) \ \tfrac{1}{4} \ c_W W^{®} e^{i \ _sW^c}$ and $!(W) \ \tfrac{1}{4} \ c_! W^{\&} e^{i \ _!W^d}$ as $W \ !\ 1$:

(a) There exists a sequence $\circ^{\pi}_n$ unique up to $o(\circ^{\pi i \ v})$ for some $v > 0$ that minimizes the asymptotic $MSE(\hat{\imath}_s)$ (or $\hat{\imath}_l$):

(b) The optimal asymptotic $MSE^{\pi}$ can be represented as a product of a polynomial component $n^{i \ \zeta}; \zeta_{\ s} \ 0$ and a logarithmic component $O((\ln n)^{\circ})$; where $\zeta$ depends only on the parameters which characterize the leading term in the tail of $W$, i.e. $®$ for a Pareto and $_s; c$ for a Weibull or combined tail, and parameters of the leading term of $!(W)$:

(c) There exist bounds $\circ_H$ and $\circ_L$ such that $\circ_L < \circ^{\pi} < \circ_H$; where $\circ_H$ and $\circ_L$ are

functions of the coe₵cients of the leading terms in tail of $W$ and $!(W)$ only and $MSE(°_H)$, $MSE(°_L)$ decline at a rate with the polynomial component $n^{i \; ¿(H)}$; $n^{i \; ¿(L)}$ with one (or both) of $¿(H)$ and $¿(L)$ either equal to $¿$; or arbitrarily close to $¿$:

(d) When $0 \cdot d < c$, that is $!(W)$ goes to zero exponentially slower than the tail of $W$ distribution, $¿ = 0$ and only a logarithmic rate of decline (at best) can be obtained for $MSE$. When, conversely, $0 \cdot c < d$, $¿ = 1$:

**Proof**: See Appendix B. □

Appendix B also provides the speci...c form $¿$; $¿(L)$; and $¿(H)$ take for all cases considered in Proposition 3.

As an example of Proposition 4, if $U$; " and $W$ are jointly normally distributed ($c = d = 2$), $°^¤$ is bounded by $°_L = (\mu_L \ln n)^{1=2}$; $\mu_L < 2¾_¤^2¾_W^2 = (¾_¤^2 + 2¾_W^2)$, and $°_H = (\mu_H \ln n)^{1=2}$; $\mu_H > 2¾_¤^2¾_W^2 = (¾_¤^2 + 2¾_W^2)$. The optimal asymptotic MSE has the polynomial component $n^{i \; ¿}$ with $¿ = 1 \; ¡ \; ¾_¤^2 = (¾_¤^2 + 2¾_W^2)$. We note that $¿$ can get arbitrarily close to 1 if $¾_W^2 \; \grave{A} \; ¾_¤^2$: (In general, if $c = d$, $¿$ can be made arbitrarily close to 1 given $_¸ ¿ {}^1$). The $MSE(°_H) \; \frac{1}{4} \; avar(\mathbf{b}_s)$ with $¿_H = 1 \; ¡ \; \mu_H = (2¾_W^2)$, and $MSE(°_L) \; \frac{1}{4} \; abias(\mathbf{b}_s)$ with $¿_L = \mu_L = ¾_¤^2$: Both $¿_L$ and $¿_H$ are arbitrarily close to $¿$ when $\mu_L$; $\mu_H$ are close to $2¾_¤^2¾_W^2 = (¾_¤^2 + 2¾_W^2)$:

After characterizing the optimal bandwidth parameter, it remains to determine the "optimal" choice of function $s(¢)$ (or $1(¢)$) in situations where an asymptotic bias is present. Following Proposition 4, if all the parameters of the tail distribution of $W$ and in the function $!(W)$ are known solving the ...rst-order condition for $°$ and then substituting into the MSE and minimizing over $q$ (where it appears) as well would give us the "optimal" estimator. When the bounds $°_H$, or $°_L$, are used instead of $°^¤$ they imply dominance of MSE by asymptotic variance or abias, correspondingly. This in turn implies preference for the function $1(¢)$ or function $s(¢)$ with large value of $q$ (where it matters) correspondingly.

The use of the bounds $°_H$, or $°_L$, as the desired bandwidth (bringing MSE close or equal to its best possible rate) might be more practical since the bounds are functions of the leading terms of the tails of $W_i$ and $!(W_i)$ only. For fat tailed distributions, we can estimate the upper tail index of a distribution e.g., using Hill (1975) and

11

Danielsson and De Vries (1997) (see Huisman et al. (1997) for its estimation in small samples). Alternatively, a probability weighted moment estimator (or maximum likelihood) of the parameters from the generalized extreme value distribution can be considered (Hosking et al. (1985)).

# 5  Conclusions

The paper presents the asymptotic behaviour of the intercept in the sample selection model based on "identi...cation at in...nity," which was ...rst proposed by Heckman (1990). Technical problems in derivations arise from the non-di¤erentiability of the indicator function. This problem was circumvented by AS via introduction of a differentiable function to replace the indicator function. Here we deal with the problem by introducing an assumption on the p.d.f. that essentially permits to obtain the expectation of the "derivative" of the indicator function.

Next, the paper examines the selection of the bandwidth and choice of the estimator of the intercept (Heckman 1990 versus AS 1998) using as a criterion the asymptotic MSE. Two characteristics of the model are of importance for such a choice: the tail behaviour of the selection index, $W_i$, and the tail behaviour of a function $!(W)$ that determines the asymptotic bias of the estimator. A wide class of distributional assumptions for the model is investigated, speci...cally, $1 ¡ F(W) ¼ c_W W^® e^{i \cdot W^c}$ and $!(W) ¼ c_! W^& e^{i \cdot^1 W^d}$: This class of models includes Pareto, Weibull as well as "combined" tails.

We have shown that for fat-tailed distributions of $W_i$ the choice of the function $s(¢)$ or $1(¢)$ does not a¤ect the asymptotic variance; the asymptotic bias if una¤ected if additionally $!(W_i)$ does not go to zero too fast. In general, however, the asymptotic MSE may be a¤ected by the choice of function $s(¢)$ but then via $q$ and the value of the derivative $s^{(q)}(0)$ only, where $q$ is the order of the ...rst non-zero derivative of $s(¢)$ at zero.

If all the parameters determining the tail behaviour of $W$ and the function $!(W)$ were known a solution that would provide an optimal $°^¤ = argmin(MSE)$ as a function of $n$, $q$ (the $s(¢)$ function) and all those parameters could be obtained (at least via a numerical algorithm) from the formulae in Proposition 3. Similarly, the

12

"optimal" choice of function $s(\cdot)$ (or $1(\cdot)$) can be obtained. Since the solution for the optimal bandwidth (and choice of function $s(\cdot)$ (or $1(\cdot)$) may not be practical, we provide simpler bounds on the optimal bandwidth parameter; using a bound may imply preference for either AS or the Heckman estimator.

Asymptotically, we give preference to the Heckman estimator in cases where there is no asymptotic bias and reveal the equivalence of the two estimators under fat-tailed distributions of $W_i$ if additionally $\lambda(W_i)$ does not go to zero too fast. For thinner-tailed distributions the decision is less clear, nevertheless, we argue that the optimal selection of the bandwidth is of primary importance. For finite samples, the AS estimator might still have advantages over the Heckman estimator, in that the trade-off between bias and variance, like in nonparametric estimation problems, is better for smooth "kernels". Nevertheless, only observations at the margin are affected by the choice of the function $s(\cdot)$. In Schafgans (1998b) simulations are presented that reveal these findings clearly.

# Appendix A: Asymptotic normality result for the Heckman estimator

To prove asymptotic normality, all we need to show is that (5) holds. We start by deriving a few sufficient conditions for (5). As in the proof of Theorem A-1 of AS the left hand side of (5) can be written as $C(\frac{\hat{A}}{\hat{B}} - \frac{A}{B}) = C\frac{\hat{A}-A}{B}\frac{B}{\hat{B}} - C\frac{\hat{B}-B}{B}\frac{A}{B}\frac{B}{\hat{B}}$, where $C = \sqrt{nE D_i 1(W_i > \gamma_n)} = \frac{3}{4}$, $A = \sum_{i=1}^{n}(Y_i - Z_i'\mu_0)D_i 1(W_i > \gamma_n)$, $\hat{A} = \sum_{i=1}^{n}(Y_i - Z_i'\hat{\mu})D_i 1(X_i'\hat{b} > \gamma_n)$, $B = \sum_{i=1}^{n} D_i 1(W_i > \gamma_n)$, and $\hat{B} = \sum_{i=1}^{n} D_i 1(X_i'\hat{b} > \gamma_n)$. To show (5), therefore, it suffices to show that

$$
\begin{align*}
&\text{(i)} \quad \tfrac{\hat{B}}{B} \xrightarrow{p} 1 \\
&\text{(ii)} \quad C(\hat{A} - A) \xrightarrow{p} 0 \\
&\text{(iii)} \quad \tfrac{A}{B} = O_p(1) \\
&\text{(iv)} \quad C(\hat{B} - B) \xrightarrow{p} 0:
\end{align*}
\tag{A.1}
$$

From Assumption 7 and Lemma A-2 of AS it follows that $C \to 1$ which means that (iv) implies (i) in (A.1). From Lemmas A-1 and A-2 of AS one gets that $\frac{3}{4}C/B$ equals $n^{-1/2} \Pr(W_i > \gamma_n)^{-1/2}$ in probability and thus (ii) is implied by the following sufficient conditions.

$$\frac{\pm\frac{1}{n}\sum_{i=1}^{n}(U_i + {}^1{}_0)D_i{}^{\circledR}{}_{in}}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}} \,!^p\, 0; \tag{A.2}$$

$$\frac{i\;(\beta_i\;\mu_0)^0\pm\frac{1}{n}\sum_{i=1}^{n}Z_iD_i{}^{\circledR}{}_{in}}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}} \,!^p\, 0; \text{ and} \tag{A.3}$$

$$\frac{i\;(\beta_i\;\mu_0)^0\pm\frac{1}{n}\sum_{i=1}^{n}Z_iD_i1(W_i > {}^{\circ}{}_n)}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}} \,!^p\, 0; \tag{A.4}$$

where ${}^{\circledR}{}_{in}$ is given by (6).

Condition (iii) has been shown to hold in the proof of Theorem A-1 in AS, it is equivalent to

$$\frac{\frac{1}{n}\sum_{i=1}^{n}U_iD_i1(W_i > {}^{\circ}{}_n)}{\Pr(W_i > {}^{\circ}{}_n)} = O_p(1): \tag{A.5}$$

Finally, Condition (iv) would follow if we show that

$$\frac{\pm\frac{1}{n}\sum_{i=1}^{n}D_i{}^{\circledR}{}_{in}}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}} \,!^p\, 0: \tag{A.6}$$

Note that under Assumptions 1, 5 and 6 of AS and using Hölder's inequality the expression in (A.4) is bounded by

$$\frac{O_p(1)E\;kZ_ik\,1(W_i > {}^{\circ}{}_n)}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}} \cdot O_p(1)\,{}^i E\,kZ_ik^{3\,\mathbb{C}_{1=3}}(\Pr(W_i > {}^{\circ}{}_n))^{2=3i\;1=2}\,!\;0: \tag{A.7}$$

The proof of Theorem 1, therefore, requires us to show that (A.2), (A.3) and (A.6) which involve the discontinuous function ${}^{\circledR}$ hold.

We do this in three steps: ...rst, a technical lemma is given (Lemma 1), then a lemma is given which examines terms in the expressions of interest (Lemma 2) and ...nally we give the proofs of the theorems which combine the intermediate results.

Lemma 1: Under Assumptions 1, 2, 6, 7 of AS and Assumption A, there exists a su¢ciently slowly increasing $\{M_n\}$ such that $n^{i\;1=2}M_n < 1$ and

$$\sup_{k^{-}_i\;{}^{-}_0k<\frac{M_n}{\sqrt{n}}}\frac{E(\sum_k{}_{-i}{}^{\circledR}{}_n({}^{-};W_i;X_i))}{\Pr(W_i > {}^{\circ}{}_n)^{1=2}}\,!\;0 \text{ as } n\,!\;1 \text{ for } {}^{-}_i = Z_i=\sqrt{n},\; {}^{-}_i = U_i, \text{ or } {}^{-}_i = 1;\,8i:$$

Proof of Lemma 1:

14

It will be sufficient to show that for the sequence $\{n^{\frac{1}{3}}\}$

$$\sup_{\|\tau_i - \tau_0\| < \frac{M_n}{n}} \frac{\frac{1}{\sqrt{n}} E \left[ \|-_i\|^3 \cdot |\varphi_n(\tau; W_i; X_i)| \cdot 1(\|X_i\| > n^{\frac{1}{3}}) \right]}{\Pr(W_i > \circ_n)^{1/2}} \rightarrow 0; \qquad (A.8)$$

and

$$\sup_{\|\tau_i - \tau_0\| < \frac{M_n}{n}} \frac{\frac{1}{\sqrt{n}} E \left[ \|-_i\|^3 \cdot |\varphi_n(\tau; W_i; X_i)| \cdot 1(\|X_i\| \cdot n^{\frac{1}{3}}) \right]}{\Pr(W_i > \circ_n)^{1/2}} \rightarrow 0: \qquad (A.9)$$

Noting that $|\varphi_n(\tau; W_i; X_i)| \cdot 1$, we get for (A.8) in the case $-_i = U_i$ (using the independence condition of Assumption 2(b) of AS) and similarly for $-_i = 1$, under Assumptions 1 and 7 of AS

$$\cdot \frac{\frac{1}{\sqrt{n}} E \left( \|-_i\| |\varphi_n(\tau; W_i; X_i)| 1(\|X_i\| > n^{1/3}) \right)}{\Pr(W_i > \circ_n)^{1/2}} \cdot \frac{\frac{1}{\sqrt{n}} E \left( \|-_i\| 1(\|X_i\| > n^{1/3}) \right)}{\Pr(W_i > \circ_n)^{1/2}}$$
$$\cdot \frac{\frac{1}{\sqrt{n}} E\|-_i\| \Pr(\|X_i\| > n^{1/3})}{\Pr(W_i > \circ_n)^{1/2}} \cdot \frac{\frac{1}{\sqrt{n}} E\|-_i\| E\|X_i\|^3 = n}{\Pr(W_i > \circ_n)^{1/2}} \qquad (A.10)$$
$$= \frac{E\|-_i\| E\|X_i\|^3}{(n \Pr(W_i > \circ_n))^{1/2}} \rightarrow 0$$

using Jensen's and Markov's inequalities. In the case $-_i = Z_i = \frac{1}{\sqrt{n}}$; (A.8) converges to zero even faster, since similarly

$$\frac{E \left( \|Z_i\| |\varphi_n(\tau; W_i; X_i)| 1(\|X_i\| > n^{1/3}) \right)}{\Pr(W_i > \circ_n)^{1/2}} \cdot \frac{\left( E\|Z_i\|^3 \right)^{1/3} (E\|X_i\|^3)^{2/3}}{n^{1/6} (n \Pr(W_i > \circ_n))^{1/2}} \rightarrow 0: \qquad (A.11)$$

The left hand side of (A.9) for $-_i = Z_i = \frac{1}{\sqrt{n}}$ is bounded by

$$\frac{\left[ E\|Z_i\|^2 \right]^{\frac{1}{2}}}{(n \Pr(W_i > \circ_n))^{1/4}} @ \left( \sup_{\|\tau_i - \tau_0\| < \frac{M_n}{n}} \frac{\frac{1}{\sqrt{n}} E |\varphi_n(\tau; W_i; X_i)| 1(\|X_i\| < n^{1/3})}{\Pr(W_i > \circ_n)^{1/2}} \right)^{1/2} \qquad (A.12)$$

using Hölder's inequality. Given Assumptions 1 and 7 of AS, therefore, it remains only to show (A.9) for $-_i = 1$ and $U_i$.

Using the independence of $-_i$ and $X_i$ for $-_i = 1$ and $U_i$ (by Assumption 2(b) of AS) we can rewrite the denominator of (A.9) as

$$E (\|-_i\|) \cdot \sup_{\|\tau_i - \tau_0\| < \frac{M_n}{n}} \frac{1}{\sqrt{n}} E |\varphi_n(\tau; W_i; X_i)| 1(\|X_i\| < n^{1/3}): \qquad (A.13)$$

By examining the function $\varphi_n(\tau; W_i; X_i)$; (6), we can see that $|\varphi_n(\tau; W_i; X_i)|$ equals 1 if either $\circ_n < W_i < \circ_n + |X_i'(\tau_i - \tau_0)|$ for negative $X_i'(\tau_i - \tau_0)$ or if $\circ_n - |X_i'(\tau_i - \tau_0)| <$

15

$W_i < \circ_n$ for positive $X_i^0(\bar{}_i - \bar{}_0)$ and zero otherwise. In view of the restrictions on $k\bar{}_i - \bar{}_0 k$ and $kX_i k$ we have $jX_i^0(\bar{}_i - \bar{}_0)j < \frac{M_n}{n^{1=6}}$, which implies $W_i > \circ_n i \frac{M_n}{n^{1=6}}$. Hence,

$$\sup_{k\bar{}_i - \bar{}_0 k < \frac{M_n}{n}} \mathsf{P}\,\overline{n}E\,j\circledR_n(\bar{};W_i;X_i)j\,1(kX_i k < n^{1=3})$$

$$= \sup_{k\bar{}_i - \bar{}_0 k < \frac{M_n}{n}} \mathsf{P}\,\overline{n}E\,j\circledR_n(\bar{};W_i;X_i)j\,1(W_i > \circ_n i \frac{M_n}{n^{1=6}})1(kX_i k < n^{1=3}): \qquad (A.14)$$

Consider the linear transformation $L : X \,!\, \begin{smallmatrix} i w \cent \\ i \end{smallmatrix}$, where $i = X_{(i\,1)}$ is the vector of all components of $X$ except the …rst (NB the matrix of this transformation is

$L = \begin{smallmatrix} \bar{}_{01} & \cdots & \bar{}_{0k} \\ 0 & I_{k i 1} \end{smallmatrix}$ ; which is non-singular). Then $X = L^{i\,1} \begin{smallmatrix} i w \cent \\ i \end{smallmatrix}$ and $X^0(\bar{}_i$

$\bar{}_0) = W\,(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(1)} + i\,^0\,(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(i\,1)}$, where subscript (1) denotes the …rst component of a vector. We denote $(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(1)} = b^7$ and $(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(i\,1)} = B$. Re-examining, for negative $X_i^0(\bar{}_i - \bar{}_0)$; the condition $\circ_n < W_i < \circ_n + jX_i^0(\bar{}_i - \bar{}_0)j = \circ_n i\,W\,(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(1)} i\,i\,^0\,(L^{0i\,1}(\bar{}_i - \bar{}_0))_{(i\,1)}$ we realize that for negative $X_i^0(\bar{}_i - \bar{}_0)$ $j\circledR_n(\bar{};W_i;X_i)j$ equals 1 for $\circ_n < W_i < \frac{\circ_n i\,i\,^0 B}{1+b}$. Similarly for positive $X_i^0(\bar{}_i - \bar{}_0)$ $j\circledR_n(\bar{};W_i;X_i)j$ equals 1 for $\frac{\circ_n i\,i\,^0 B}{1+b} < W_i < \circ_n$. Thus (A.14) equals

$$\sup_{k\bar{}_i - \bar{}_0 k < \frac{M_n}{n}} \mathsf{P}\,\overline{n}E\,4 \int_{\circ_n}^{\frac{\circ_n i\,i\,^0 B}{1+b}} 1(X_i^0(\bar{}_i - \bar{}_0) < 0)1(kX_i k < n^{1=3})p:d:f:_{W_i|j_i}(W;_i)dW +$$

$$\int_{\frac{\circ_n i\,i\,^0 B}{1+b}}^{\circ_n} 1(X_i^0(\bar{}_i - \bar{}_0) > 0)1(W_i > \circ_n i \frac{M_n}{n^{1=6}})1(kX_i k < n^{1=3})p:d:f:_{W_i|j_i}(W;_i)dW \;\#:$$

$$(A.15)$$

The $p.d.f._{W_i|j_i}$ exists for large enough $n$, since the expression under the integrals is non-zero only if $W_i > \circ_n i \frac{M_n}{n^{1=6}}$ and thus as $\circ_n \,!\, 1$ for $M_n = O(n^{1=6})$ becomes greater than $A$; the conditional $p.d.f._{W_i|j_i}$ declines monotonically in $W$. Denote $\circ_n i \frac{M_n}{n^{1=6}}$ by $\overset{\circ}{n}_n$: The …rst integral on the right hand side of (A.15) is bounded by $p:d:f:_{W_i|j_i}(\overset{\circ}{n};_i)\cent \frac{\circ_n i\,i\,^0 B}{1+b} i\,\circ_n$, the second can be bounded by $p:d:f:_{W_i|j_i}(\overset{\circ}{n};_i)\cent \circ_n i \frac{\circ_n i\,i\,^0 B}{1+b}$ : Thus the sum is bounded by

$$\sup_{k\bar{}_i - \bar{}_0 k < \frac{M_n}{n}} 2\mathsf{P}\,\overline{n} \int p:d:f:_{W_i|j_i}(\overset{\circ}{n};_i)\left|\frac{i\,^0 B}{1+b}\right| dPr_i \qquad (A.16)$$

---

[7]Note: $jbj \cdot \overset{\circ}{\|}L^{i\,1}\overset{\circ}{\|} \cent k\bar{}_i - \bar{}_0 k \cdot \overset{\circ}{\|}L^{i\,1}\overset{\circ}{\|} \frac{M_n}{n}$. Assume that $M_n$ is such that $\overset{\circ}{\|}L^{i\,1}\overset{\circ}{\|} \frac{M_n}{n} \cdot G < 1$.

16

$$\cdot \quad \frac{2M_n}{1+b} \int \|\xi_i\| \, k \, p.d.f._W(\overset{\circ}{\omega}_n) dPr_{i|W=\overset{\circ}{\omega}_n} = \frac{2M_n}{1+b} p.d.f._W(\overset{\circ}{\omega}_n) E_{|W=\overset{\circ}{\omega}_n} \|\xi_i\| \, k ;$$

where $\xi_i = X_{(i\,1)}$ and $\|B\| = \overset{\circ}{\circ}(L^{i\,1}(\bar{}_i \, \bar{}_0))_{(i\,1)}\overset{\circ}{\circ} \; \leq \; 8 \overset{\circ}{\circ}(\bar{}_i \, \bar{}_0)_{(i\,1)}\overset{\circ}{\circ} \cdot \frac{M_n}{n^{1=2}}$ using the restriction on $\|\bar{}_i \, \bar{}_0\|$: The notation $dPr_i$ indicates that a Stiltjes integral with respect to the cumulative probability function $_i$ is taken, if a marginal density exists $dPr_i = p.d.f._{_i}(_i)d_i$. By Jensen's inequality, and in view of Assumption A(b) this is bounded almost surely by

$$\frac{2M_n}{1+b} \, \phi \, p.d.f._W(\overset{\circ}{\omega}_n) \, \phi \, E_{|W=\overset{\circ}{\omega}_n} \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3}{}^{1=3} : \qquad (A.17)$$

Let us consider $E_{|W=\overset{\circ}{\omega}_n} \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3}$. By Assumption 1 of AS, we know that the unconditional expectation $E \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3}$ is bounded; this implies that for the ...xed $A$ of Assumption A $E \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3} 1(W > A)$ is bounded as well.

$$E \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3} 1(W > A) = \iint \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3} 1(W > A) p.d.f._{W;\,_i|W}(_i) p.d.f._W(W) d_i \, dW$$

$$= \int J(W) p.d.f._W(W) 1(W > A) dW; \qquad (A.18)$$

where $J(W) = \int \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3} p.d.f._{W;\,_i|W}(_i) d_i$ is a function of $W$ only. Since the integral $\int J(W) p.d.f._W(W) dW \leq E \overset{\circ}{\circ}\|X_{(i\,1)}\|^{\overset{\circ}{\circ}3}$ exists, it implies that as $W \to 1$ $J(W) p.d.f._W(W) = o(W^{i\,1})$, or

$$J(W) = o\left(W^{i\,1}(p.d.f._W(W))^{i\,1}\right) : \qquad (A.19)$$

Equation (A.17), can then be rewritten as follows

$$\frac{2M_n}{1+b} \, \phi \, p.d.f._W(\hat{e}_n) \, \phi \, J(\hat{e}_n)^{1=3} = \frac{2M_n}{1+b} \, \phi \, p.d.f._W(\overset{\circ}{\omega}_n \, \frac{M_n}{n^{1=6}})^{2=3} o(\overset{\circ}{\omega}_n^{i\,1=3}): \qquad (A.20)$$

This implies that for $-_i = 1$ and $U_i$ (A.9) can be bounded as

$$\frac{E(\|-_i\|) \frac{2M_n}{1+b} \, \phi \, p.d.f._W(\overset{\circ}{\omega}_n \, \frac{M_n}{n^{1=6}})^{2=3} o(\overset{\circ}{\omega}_n^{i\,1=3})}{Pr(W_i > \overset{\circ}{\omega}_n)^{1=2}}$$

$$\cdot \quad E(\|-_i\|) \frac{2M_n}{1+b} \left(\frac{p.d.f._W(\overset{\circ}{\omega}_n \, \frac{M_n}{n^{1=6}})}{Pr(W_i > \overset{\circ}{\omega}_n)^{3=4}}\right)^{2=3} o(\overset{\circ}{\omega}_n^{i\,1=3}): \qquad (A.21)$$

---

[8] Here we use the fact that $L^{i\,1} = \begin{pmatrix} \frac{1}{_{01}} & \frac{\bar{}_{02}}{_{01}} \, \phi\phi\phi \, \frac{\bar{}_{0k}}{_{01}} \\ 0 & I_{k\,i\,1} \end{pmatrix}$.

17

Set $M_n = \min\{dn^{1-\delta}, {}^{\circ}{}^{1-\beta}\}$ with $d > 0$ as in Assumption A. Using Assumptions 1 and 6 of AS and Assumption A(a) the right-hand side of (A.21) converges to zero for $M_n$. $\blacksquare$

The next lemma will help to show that terms, involving $\mathfrak{R}_{in}$, which appear in (A.2), (A.3), and (A.6) have a zero probability limit.

**Lemma 2:** Under Assumptions 1, 2, 5–7 of AS, and Assumption A
$$\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_{in}j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} \xrightarrow{p} 0 \text{ as } n \to \infty \text{ for } \tau_i = Z_i, \tau_i = U_i, \text{ or } \tau_i = 1; \forall i:$$

**Proof of Lemma 2:**

We would like to show that for any $\varepsilon > 0$

$$\Pr\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_{in}j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} > \varepsilon\right) \to 0 \text{ as } n \to \infty: \tag{A.22}$$

From Lemma 1, let $M_n$ satisfy $M_n o({}^{\circ}{}_h^{1-\beta}) \to 0$ and $\frac{M_n}{n^{1-\delta}} \cdot d$. The left hand side of (A.22) is equal to

$$\Pr\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_{in}j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} > \varepsilon, \overset{\circ}{\mathbf{b}}_i - {}_0\overset{\circ}{} \cdot \frac{M_n}{\sqrt{n}}\right) + \Pr\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_{in}j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} > \varepsilon, \overset{\circ}{\mathbf{b}}_i - {}_0\overset{\circ}{} > \frac{M_n}{\sqrt{n}}\right)$$

$$\cdot \Pr\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_{in}j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} > \varepsilon, \overset{\circ}{\mathbf{b}}_i - {}_0\overset{\circ}{} \cdot \frac{M_n}{\sqrt{n}}\right) + \Pr\left(\overset{\circ}{\mathbf{b}}_i - {}_0\overset{\circ}{} > \frac{M_n}{\sqrt{n}}\right): \tag{A.23}$$

The second expression on the right hand side of (A.23) converges to zero by Assumption 5 of AS, since $\sqrt{n}(\mathbf{b}_i - {}_0) = O_p(1)$ implies $\sqrt{n}(\mathbf{b}_i - {}_0)/M_n = o_p(1)$. The first expression on the right hand side of (A.23) is bounded by

$$\sup_{\|k-i - {}_0\| \cdot k \cdot \frac{M_n}{\sqrt{n}}} \Pr\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_n(\tau; W_i; X_i)j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}} > \varepsilon\right) \tag{A.24}$$

$$\cdot \sup_{\|k-i - {}_0\| \cdot k \cdot \frac{M_n}{\sqrt{n}}} \frac{E\left(\frac{\frac{1}{n}\sum_{k-i}k_j\mathfrak{R}_n(\tau; W_i; X_i)j}{\Pr(W_i > {}^{\circ}{}_n)^{1-2}}\right)}{\varepsilon} = \sup_{\|k-i - {}_0\| \cdot k \cdot \frac{M_n}{\sqrt{n}}} \frac{\frac{1}{n}\sum E\,k_{-i}k_j\mathfrak{R}_n(\tau; W_i; X_i)j}{\varepsilon \Pr(W_i > {}^{\circ}{}_n)^{1-2}};$$

where the inequality is based on Markov's inequality. This term converges to zero for all $\varepsilon > 0$ by Lemma 1. $\blacksquare$

**Proof of Theorem 1:** For our proof it is sufficient to show (A.2), (A.3), and (A.6). By Assumptions 1 and 5 for $\mathbf{b}_i - \mu_0$ of AS, the left-hand sides of (A.2), (A.3) can be

18

bounded by:

$$\frac{\mathbb{P}\frac{1}{n}\sum_{i=1}^{n}|U_i||j^{\circledR}_{in}|}{\Pr(W_i > {}^{\circ}_n)^{1=2}} + |1_0| \frac{\mathbb{P}\frac{1}{n}\sum_{i=1}^{n}j^{\circledR}_{in}j}{\Pr(W_i > {}^{\circ}_n)^{1=2}};\qquad\text{(A.25)}$$

$$\frac{O_p(1)\frac{1}{n}\sum_{i=1}^{n}\|Z_i\|\,|j^{\circledR}_{in}j}{\Pr(W_i > {}^{\circ}_n)^{1=2}};\qquad\text{(A.26)}$$

respectively. From Lemma 2, (A.25) and (A.26) have zero probability limits.

Finally, to prove (A.6), we note that its left-hand side can be bounded by

$$\frac{\mathbb{P}\frac{1}{n}\sum_{i=1}^{n}j^{\circledR}_{in}j}{\Pr(W_i > {}^{\circ}_n)^{1=2}};\qquad\text{(A.27)}$$

which by Lemma 2 again converges in probability to zero. This completes the proof of our theorem. $\quad\square$

**Proof of Theorem 2:** By Cramer-Wold device, this result follows directly from (5) and Theorem A-4 in AS. In the latter, the result in Theorem 2 is shown for the case where $\mathbf{b}_l$ is replaced by $\mathbf{b}_{l;0}$. $\quad\square$

# Appendix B: Asymptotic variance and bias: Selection of Bandwidth and Estimator

In this Appendix we provide the asymptotic bias, variance, and mean squared error of the AS and Heckman estimator. Using Lemma A–2 of AS and Theorem 1 of this paper, we write the asymptotic bias and variance as

$$\text{abias}(\mathbf{b}_s)\ \tfrac{1}{4}\ \dot{\iota}\ \frac{E!\,(W_i)s(W_i\ \dot{\iota}\ {}^{\circ}_n)}{Es(W_i\ \dot{\iota}\ {}^{\circ}_n)};\quad \text{abias}(\mathbf{b}_l)\ \tfrac{1}{4}\ \dot{\iota}\ \frac{E!\,(W_i)1(W_i > {}^{\circ}_n)}{E1(W_i > {}^{\circ}_n)};\text{(B.1)}$$

$$\text{avar}(\mathbf{b}_s)\ \tfrac{1}{4}\ \tfrac{3}{4}^2 n^{\dot{\iota}\,1}\frac{Es^2(W_i\ \dot{\iota}\ {}^{\circ}_n)}{[Es(W_i\ \dot{\iota}\ {}^{\circ}_n)]^2};\quad \text{avar}(\mathbf{b}_l)\ \tfrac{1}{4}\ \tfrac{3}{4}^2 n^{\dot{\iota}\,1}\,(E1(W_i > {}^{\circ}_n))^{\dot{\iota}\,1}_{\dot{.}}\text{(B.2)}$$

In the following we let $!\,(W_i)\ {}_{\,\text{¿}}\ 0\ 8W_i > {}^{\circ}_n$ (similar proofs can be given when $!\,(W_i)\cdot\ 0\ 8W_i > {}^{\circ}_n$).

**Proof of Proposition 1:** For (a), we note that by Cauchy-Schwartz inequality

$$Es(W_i\ \dot{\iota}\ {}^{\circ}_n) = Es(W_i\ \dot{\iota}\ {}^{\circ}_n)1(W_i > {}^{\circ}_n)\cdot\ [Es^2(W_i\ \dot{\iota}\ {}^{\circ}_n)]^{1=2}[E1(W_i > {}^{\circ}_n)]^{1=2};\ \text{(B.3)}$$

This inequality combined with (B.2) gives the result that (a) holds for any sequence $°_n$.

Next we turn to (b). Using (B.1), we get

$$jabias(\mathbf{b}_s)j \; ¼ \; \frac{E! \, (W_i)s(W_i \, ¡ \, °_n)}{Es(W_i \, ¡ \, °_n)} \; \cdot \; \frac{E! \, (W_i)1(W_i > °_n)}{E1(W_i > °_n + b)}$$
$$¼ \; jabias(\mathbf{b}_l)j \frac{E1(W_i > °_n)}{E1(W_i > °_n + b)}: \tag{B.4}$$

Since $E1(W_i > °_n)=E1(W_i > °_n + b) \; ¼ \; 1$ (b) follows.   2

**Proof of Proposition 2**: From Proposition 1(a), we know

$$\frac{Es^2(W_i \, ¡ \, °_n)}{[Es(W_i \, ¡ \, °_n)]^2} \; , \; [E1(W_i > °_n)]^{i \, 1}: \tag{B.5}$$

Furthermore as $Es^2(W_i \, ¡ \, °_n) \cdot Es(W_i \, ¡ \, °_n)$ and $Es(W_i \, ¡ \, °_n) \; , \; E1(W_i > °_n + b)$;

$$\frac{Es^2(W_i \, ¡ \, °_n)}{[Es(W_i \, ¡ \, °_n)]^2} \cdot [E1(W_i > °_n + b)]^{i \, 1} \; ¼ \; [E1(W_i > °_n)]^{i \, 1}; \tag{B.6}$$

and (a) follows.

Under the same assumptions, we know from Proposition 1(b) that

$$\frac{E! \, (W_i)s(W_i \, ¡ \, °_n)}{Es(W_i \, ¡ \, °_n)} \cdot \frac{E! \, (W_i)1(W_i > °_n)}{E1(W_i > °_n + b)} \; ¼ \; \frac{E! \, (W_i)1(W_i > °_n)}{E1(W_i > °_n)}: \tag{B.7}$$

In addition, we have

$$\frac{E! \, (W_i)s(W_i \, ¡ \, °_n)}{Es(W_i \, ¡ \, °_n)} \; , \; \frac{E! \, (W_i)1(W_i > °_n + b)}{E1(W_i > °_n)} \; ¼ \; \frac{E! \, (W_i)1(W_i > °_n)}{E1(W_i > °_n)}: \tag{B.8}$$

Combining these inequalities we obtain (b).   2

For the remainder of this appendix we have omitted the subscript n on ° to simplify notation. For the proof of Proposition 3, we make use of the following technical lemma:

**Lemma 3**: As $° \, ! \, 1$,

(a) $\displaystyle\int_°^1 {}_,c(W \, ¡ \, °)^i W^v \exp(¡ \, {}_,W^c) dW = °^{i+v_i \, c+1_i \, ci} \exp(¡ \, {}_,°^c)({}_,c)^{i \, i}i!(1 + o(°^{i \, c}))$

20

(b) For $a>f$, $\int_0^1 a(W_i{}^\circ)^i W^v \exp(_i{}^1 W^a {}_i {}_{\circ}W^f)dW =$

$$\circ^{i+v_i a+1_i ai}\exp(_i{}^{1\circ a}{}_i{}_{\circ}^{\circ f})(^1 a)^{i}{}_i{}^i i!(1 + o(^{\circ_i a})):$$

**Proof:** The integral in (a), $I$, can be rewritten as

$$I = \int_\circ^{\mathbb{Z}} {}_\circ c \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^{i+j\circ i_i j} W^{j+v}\exp(_i{}_\circ W^c)dW$$

$$= \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^{i+j\circ i_i j} \int_\circ^{\mathbb{Z}} {}_\circ c W^{j+v}\exp(_i{}_\circ W^c)dW: \qquad (B.9)$$

By setting $W^c = z$, we get

$$I = \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^{i+j\circ i_i j} \mathbb{C}\int_{\circ c}^{\mathbb{Z}_1} z^{\frac{j+v_i c+1}{c}}\exp(_i{}_\circ z)dz: \qquad (B.10)$$

Combining 3.381#3 and 8.357 in Gradshteyn and Ryzhik (1994)[9] we obtain

$$I = \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^{i+j\circ i_i j} {}_\circ \mathbb{C}{}_i{}^{\frac{j+v+1}{c}}{}_i \left(\frac{j+v+1}{c}; {}_\circ{}^{\circ c}\right)$$

$$= \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^{i+j\circ i_i j} {}_\circ \mathbb{C}^{\circ j+v_i c+1}\exp(_i{}_\circ{}^{\circ c}) \mathbb{C}$$

$$\tilde{A}\left(\sum_{m=0}^1 \frac{(_i 1)^m{}_i (1{}_i \frac{j+v+1}{c} + m)}{\circ cm{}_\circ{}^m{}_i (1{}_i \frac{j+v+1}{c})} + o(^{\circ_i cL})\right) \qquad (B.11)$$

$$= (_i 1)^{i\circ i+v_i c+1}\exp(_i{}_\circ{}^{\circ c}) \sum_{j=0}^{\mu_i\P} \binom{}{j} (_i 1)^j \tilde{A}\left(\sum_{m=0}^1 \frac{P_m{}_{s=0}\cdot{}_s(\frac{i}{c})^s}{\circ cm{}_\circ{}^m} + o(^{\circ_i cL})\right);$$

where we have substituted $(_i 1)^m{}_i (1{}_i \frac{j+v+1}{c} + m)=_i (1{}_i \frac{j+v+1}{c}) = P^m{}_{s=0}\cdot{}_s(\frac{i}{c})^s$, with $\cdot_s = 1$ for $s = m$ and some known constant for $s = 0; ::::; m{}_i 1$: Using the fact that $P_{j=0}^{i_i\mathbb{C}} \binom{}{j}(_i 1)^j j^m = 0$ for $i{}_\circ m+1{}_\circ 1; 0^0 \,\acute{}\, 1$ (see 0.154#3 in Gradshteyn and Ryzhik

---

[9]There is a typographical error in 3.381#3 in G&R; the correct formula reads $\int_u^1 x^{v_i 1}e^{i{}^{1x}}dx = {}^1{}_i{}^v{}_i (v; {}^1 u):$

21

(1994)), we get

$$I = (\imath 1)^{\imath\circ i+v_{\imath} c+1}\exp(\imath\, \circ c)\sum_{j=0}^{\mu_{\imath}\P}(\imath 1)^j \frac{(\frac{\imath}{c})^i}{\circ ci\,\imath}(1 + o(^{\circ i\, c}))$$

$$= {}^{\circ i+v_{\imath} c+1\imath\, ci}\exp(\imath\, \circ c)(\,c)^{\imath\, \imath}i!(1 + o(^{\circ i\, c})); \tag{B.12}$$

where the second equality uses 0.154#4 in Gradshteyn and Ryzhik (1994).

Next consider the integral (b), $I^0$, for $a > f$. We notice that $I^0$ can be rewritten as follows,

$$I^0 = \exp(\imath\, \circ f)\int_{\circ}^{1} {}^{\imath}a(W \imath \circ)^i W^v \exp(\imath\, {}^1 W^a)^i \exp(\imath\, (W^f \imath \circ f))^{\mathrm{¢}} dW \tag{B.13}$$

$$= \exp(\imath\, \circ f)\int_{\circ}^{1} {}^{\imath}a(W \imath \circ)^i W^v \exp(\imath\, {}^1 W^a)\sum_{r=0}^{\tilde{A}} \frac{(\imath 1)^r {}_{\,}^r (W^f \imath \circ f)^r}{r!} dW;$$

where we have substituted the series representation for the last exponential function in the integral. By dominating convergence theorem, we can interchange the integral and summation, giving

$$I^0 = \exp(\imath\, \circ f)\sum_{r=0}^{} \frac{(\imath 1)^r {}_{\,}^r}{r!}\int_{\circ}^{1} {}^{\imath}a(W^f \imath \circ f)^r(W \imath \circ)^i W^v \exp(\imath\, {}^1 W^a)dW: \tag{B.14}$$

To complete the proof, we need to reapply the steps taken in (B.9)–(B.12). The ...nal result is obtained by setting $r = 0$ (all remaining terms converge to zero faster). A detailed proof of this can be obtained from the authors. 2

**Proof of Proposition 3:** In this proof we will not attempt to formally show all cases considered, but indicate the method used to derive the results, pointing primarily to the more complex derivations.

Using the results from Proposition 2, $avar(^1{}_{\imath}) \frac14 avar(^1{}_s)$ if $1\imath F(W) \frac14 W^{\circledR}e^{i\, W^c}$ with $c \cdot 1$, since Assumption 4 of AS holds with $» = 0$ and $E1(W > \circ)=E1(W > \circ + b) \frac14 1$. The variance in these cases can be easily derived using $\frac34^2=(nE1(W_i > \circ))$; see (B.2).[10]

When on the other hand $1 \imath F(W) \frac14 W^{\circledR}e^{i\, W^c}$ with $c > 1$ we need to derive $avar(^1{}_s)$ using the de...nition in (B.2) (Assumption 4 of AS does not hold with $» = 0$).

---

[10] For $W > \circ$ the p:d:f:(W) equals $_{\,}cW^{\circledR+c\imath\, 1}\exp(\imath\, W^c)(1 + o(^{\circ i\, c}))$, where $\circledR = \imath\,$ if $c = 0$ (Pareto tail case).

For this we need to apply Lemma 3(a). Notice that for any $\epsilon > 0$ we can write $E_s(W_i - \zeta)$ as

$$
\begin{aligned}
E_s(W_i - \zeta) \;=\; & \frac{a_q}{q!} \int_\zeta^1 (W_i - \zeta)^q \, cW^{\beta + c_i - 1} \exp(-\lambda W^c) dW \\
& -\; \frac{a_q}{q!} \int_{\zeta+\epsilon}^1 (W_i - \zeta)^q \, cW^{\beta + c_i - 1} \exp(-\lambda W^c) dW \\
& +\; \int_{\zeta+\epsilon}^{\zeta+\epsilon} \left[ s(W_i - \zeta) - \frac{a_q}{q!}(W_i - \zeta)^q \right] cW^{\beta + c_i - 1} \exp(-\lambda W^c) dW \\
& +\; \int_{\zeta+\epsilon}^1 s(W_i - \zeta) \, cW^{\beta + c_i - 1} \exp(-\lambda W^c) dW.
\end{aligned}
$$

By Lemma 3(a), the first integral is $\approx a_q \zeta^{q + \beta - cq} c^q (\lambda c)^{-q} \exp(-\lambda \zeta^c)$; the second is similarly $\approx - a_q (\zeta + \epsilon)^{q + \beta - cq} c^q (\lambda c)^{-q} \exp(-\lambda(\zeta + \epsilon)^c)$ which goes to zero at an exponentially faster rate than the first as long as $\epsilon \zeta^{c-1} \lambda \to \infty$ since $\exp(-\lambda(\zeta + \epsilon)^c) = \exp(-\lambda \zeta^c) \exp(-\lambda c \epsilon \zeta^{c-1}(1 + o(\frac{\epsilon}{\zeta})))$. Using Assumption S and Lemma 3(a), the absolute value of the third integral can be bounded by $O(1)\epsilon^{q+1} \zeta^\beta \exp(-\lambda \zeta^c)(1 + o(\zeta^{-c}))$ (we apply a $(q + 1)^{th}$ order Taylor expansion to $s(\cdot)$ around zero). If $\epsilon^{q+1} \lambda (\zeta^{c-1})^q \to 0$ this implies that the third integral goes to zero faster than the first. Finally, the fourth integral can be bounded by one where the function $1(\cdot)$ is substituted for $s(\cdot)$, thus by Lemma 3(a) it is bounded by $(\zeta + \epsilon)^\beta \exp(-\lambda(\zeta + \epsilon)^c)(1 + o(\zeta^{-c}))$ and (similarly to the second integral) goes to zero exponentially faster than the first one if $\epsilon \zeta^{c-1} \lambda \to \infty$. If the conditions are met, the first integral dominates $E_s(W_i - \zeta)$, i.e.,

$$
E_s(W_i - \zeta) \approx a_q \zeta^{q + \beta - cq} c^q (\lambda c)^{-q} \exp(-\lambda \zeta^c). \tag{B.16}
$$

For $E_s^2(W - \zeta)$, we get similarly four terms, the first of which $\approx a_q^2 \frac{(2q)!}{q!} \zeta^{2q(1-c)+\beta} c \, (\lambda c)^{-2q} \exp(-\lambda \zeta^c)$; the second is the corresponding integral from $\zeta + \epsilon$ to infinity and requires $\epsilon \zeta^{c-1} \lambda \to \infty$ to go to zero exponentially faster. The third one can analogously be bounded by

$$
\int_\zeta^{\zeta + \epsilon} \left| s^2(W_i - \zeta) - \left( \frac{a_q}{q!} \right)^2 (W_i - \zeta)^{2q} \right| p.d.f._W(W) dW
$$

$$
\cdot \quad O(1)\epsilon^{2q+1} \zeta^\beta \exp(-\lambda \zeta^c)(1 + o(\zeta^{-c}))
$$

and thus needs $\epsilon^{2q+1} \lambda (\zeta^{c-1})^{2q} \to 0$. And finally, the fourth integral can be bounded again by replacing $s^2(\cdot)$ by $1(\cdot)$, which using Lemma 3(a) goes to zero exponentially

23

faster than the ...rst one if $2^{\circ c_i}{}^1 !\ 1$. For $2 = \circ i\ \frac{(c_i\ 1)(2q+\{)}{2q+1}$ with $0 < \{ < 1$ the ...rst integral dominates $Es^2(W_i\ i\ ^\circ)$, i.e.,

$$Es^2(W_i\ i\ ^\circ)\ \tfrac{1}{4}\ a_q^2{}^{i\,2q}{}_q^{\Cent}{}_{\circ}{}^{2q+\circledR i\ 2cq}(_{\circ}c)^{i\ 2q}\exp(_{i\ \circ}{}^{\circ c}): \qquad (B.17)$$

This $^2$ also satis...es the requirement for (B.16) to hold. Combining (B.16) and (B.17) give avar$(^1{}_s)$:

Concerning the results presented in (b), if $1\ i\ F(W)\ \tfrac{1}{4}\ W^\circledR e^{i\ _\circ W^c}$ with $c\cdot\ 1$ and $!(W)\ \tfrac{1}{4}\ W^\&\exp(_i\ {}^1W^d)$ with $d\cdot\ 1$ all assumptions in Proposition 2 are satis...ed and abias$(^1{}_i)\ \tfrac{1}{4}$ abias$(^1{}_s)$. The derivations of the asymptotic bias in that case requires us to compute $_i\ E(!(W_i)1(W_i > ^\circ))$ (see (B.1)). Substituting $1\ i\ F(W)$ and $!(W)$

$$E(!(W_i)1(W_i > ^\circ)) = {}_\circ c\int_\circ^1 W^{\&+\circledR+c_i\ 1}\exp(_i\ {}^1W^d\ i\ {}_\circ W^c)(1 + o(^{\circ i\ c})dW: \quad (B.18)$$

When $d = c$, this expectation can be obtained straightforwardly using the analysis of the asymptotic variance given above. When $d\ \&\ c$, Lemma 3(b) gives

$$E(!(W_i)1(W_i > ^\circ))\ \tfrac{1}{4}\ \begin{cases} \frac{_\circ c}{{}^1 d}{}_\circ^{\&+\circledR+c_i\ d}\exp(_i\ {}^1{}_\circ{}^d\ i\ {}_\circ{}^\circ c)(1 + o(^{\circ i\ c})) & d > c \\ {}_\circ^{\&+\circledR}\exp(_i\ {}^1{}_\circ{}^d\ i\ {}_\circ{}^\circ c)(1 + o(^{\circ i\ c}) & d < c: \end{cases} \quad (B.19)$$

When $c$ or $d$ (or both) exceed 1, the asymptotic bias of the Heckman and AS estimator are not equal any more. In that case we need to extend the analysis above to derive $_i\ E(!(W_i)s(W_i\ i\ ^\circ))$. As in (B.15), we write $E(!(W_i)s(W_i\ i\ ^\circ))$ as a sum of four integrals, where $^2 > 0$

$$E(!(W_i)s(W_i\ i\ ^\circ))$$
$$= \tfrac{a_q}{q!}\int_\circ^1 (W\ i\ ^\circ)^q{}_\circ cW^{\&+\circledR+c_i\ 1}\exp(_i\ {}^1W^d\ i\ {}_\circ W^c)dW \qquad (B.20)$$
$$i\ \tfrac{a_q}{q!}\int_{\circ+2}^1 (W\ i\ ^\circ)^q{}_\circ cW^{\&+\circledR+c_i\ 1}\exp(_i\ {}^1W^d\ i\ {}_\circ W^c)dW$$
$$+ \int_{\circ+2}^{\circ+2}\Big[s(W\ i\ ^\circ)\ i\ \tfrac{a_q}{q!}(W\ i\ ^\circ)^q\Big]{}_\circ cW^{\&+\circledR+c_i\ 1}\exp(_i\ {}^1W^d\ i\ {}_\circ W^c)dW$$
$$+ \int_{\circ+2}^\circ s(W\ i\ ^\circ){}_\circ cW^{\&+\circledR+c_i\ 1}\exp(_i\ {}^1W^d\ i\ {}_\circ W^c)dW:$$

In the case where $c = d$ the result follows directly from Lemma 3(a). Using a similar discussion when $c\ \&\ d$, Lemma 3(b) gives us

$$E\left(\hat{\mathfrak{f}}(W_i)s(W_i - \delta^\circ)\right)$$

$$\frac{1}{4} \quad \begin{cases} a_q \frac{c}{(1-d)^{q+1}} \exp(-\sigma^c - \lambda^{1-d})\sigma^{q+\&+\circledR+c-d} \, d_q & d > c \\ a_q(-c)^{1-q} \exp(-\sigma^c - \lambda^{1-d})\sigma^{q+\&+\circledR-cq} & d < c: \end{cases} \qquad (B.21)$$

This completes the derivations required for the proof. $\quad 2$

**Proof of Proposition 4:** According to Proposition 3, the asymptotic MSE for all cases to be considered has the form

$$MSE = an^{-1}\sigma^{-\circledR}\exp(\lambda\sigma^c) + b\sigma^{2\gamma}\exp(-\lambda^{1-d}); \qquad (B.22)$$

where

$$\gamma = \begin{cases} \& + (c - d)(1 + q) & \text{if } d > c; d > 1 \\ \& + (c - d) & \text{if } d > c; d \cdot 1 \\ \& & \text{otherwise.} \end{cases} \qquad (B.23)$$

Thus

$$\frac{@MSE}{@\sigma} = an^{-1}\sigma^{-\circledR-1}\exp(\lambda\sigma^c)(-\circledR + \lambda c\sigma^c) + b\sigma^{2\gamma-1}\exp(-\lambda^{1-d})(2\gamma - \lambda^{1-d}\sigma^d): \quad (B.24)$$

The optimal bandwidth $\sigma^\ast$ solves the ...rst order condition in which we ignore terms that go to zero faster than the ones we keep, i.e., $\sigma^\ast$ solves:

(i) $\quad -\circledR a\exp(\lambda)n^{-1}\sigma^{-\circledR-1} + 2\gamma b\exp(-\lambda)\sigma^{2\gamma-1} = 0 \qquad$ if $c = 0; d = 0 \ (\circledR; \gamma < 0)$

(ii) $\quad \lambda can^{-1}\sigma^{-\circledR+c-1}\exp(\lambda\sigma^c) + 2\gamma b\exp(-\lambda)\sigma^{2\gamma-1} = 0 \qquad$ if $c > 0; d = 0 \ (\gamma < 0))$

(iii) $\quad -\circledR a\exp(\lambda)n^{-1}\sigma^{-\circledR-1} - \lambda^1 db\sigma^{2\gamma+d-1}\exp(-\lambda^{1-d}) = 0 \qquad$ if $c = 0; d > 0 \ (\circledR < 0)$

(iv) $\quad \lambda can^{-1}\sigma^{-\circledR+c-1}\exp(\lambda\sigma^c) - \lambda^1 db\sigma^{2\gamma+d-1}\exp(-\lambda^{1-d}) = 0 \quad$ if $c > 0; d > 0:$

We discuss each case separately.

Case (i). In this case we get an analytic solution $\sigma^\ast = \left(\frac{2\gamma b\exp(-\lambda^1)}{\circledR a\exp(\lambda)}n\right)^{\frac{1}{-\circledR-2\&}}$ and $MSE(\sigma^\ast) = a\exp(\lambda)\left(\frac{2\gamma b\exp(-\lambda^1)}{\circledR a\exp(\lambda)}\right)^{\frac{\circledR}{\circledR+2\&}} n^{\frac{2\gamma}{\circledR+2\gamma}}\left(1 + \frac{b\exp(-\lambda^1)}{a\exp(\lambda)}\right)$, so $\iota = \frac{2\gamma}{\circledR+2\gamma}:$ Here for any $\sigma = \mu n^{\frac{1}{-\circledR-2\&}}$, the corresponding rate equals $\iota$.

Case (ii). Substituting from the ...rst order conditions we can express $MSE(\sigma^\ast)$ as $\frac14 b\sigma^{\ast 2\gamma}\exp(-\lambda^1)(1 + \frac{2\gamma}{\lambda c}\sigma^{-\circ c}) \frac14 b\sigma^{\ast 2\gamma}\exp(-\lambda^1).$ If $\sigma = (\mu\ln n)^{1=c}$ the ...rst term in

25

the derivative is proportional to $n^{-1+\beta\mu}(\ln n)^{\frac{-\omega+c_i-1}{c}}$ while the second is proportional to $(\ln n)^{\frac{2'-1}{c}}$; thus for $\theta_L = (\mu_L \ln n)^{1-c}$ with $\mu_L < \frac{1}{2}$, the second negative term dominates the growth in the derivative (abias$^2$ dominates $MSE(\theta_L)$) and $\iota(L) = 0$. For $\theta_H = (\mu_H \ln n)^{1-c}$ $\mu_H > \frac{1}{2}$, the ﬁrst term dominates the derivative (avar dominates $MSE(\theta_H)$) and $\iota(H) = 1 - 2\mu_H$ (note: only if $\mu = \frac{1}{2}$ does the variance not increase with n). Since $\theta_L < \theta^* < \theta_H$, we get (a), (b), (c), and (d).

Case (iii). Here similarly, $MSE(\theta^*) = an^{-1\circ\pi i\,\circledR}\exp(\varsigma)(1 - \frac{\circledR a}{2^1 d}\circ\pi i\,d) \frac{1}{4} an^{-1\circ\pi i\,\circledR}\exp(\varsigma)$. If $\theta_L = (\mu_L \ln n)^{1-d}$ and $\mu_L < \frac{1}{2^1}$ the derivative is negative, $MSE(\theta_L)$ is dominated by abias$^2$ and declines at a rate with polynomial component $n^{-\iota(L)}$ with $\iota(L) = 2^1\mu_L$. If $\theta_H = (\mu_H \ln n)^{1-d}$ and $\mu_H > \frac{1}{2^1}$ the derivative is positive, avar dominates and declines with $\iota(H) = 1$; (a), (b), (c), and (d) follow.

Case (iv). Note that substituting from the ﬁrst order condition here $MSE(\theta^*) \frac{1}{4}$ $an^{-1\circ\pi i\,\circledR}\exp(\varsigma\circ\pi c) \frac{1}{4} - 2\frac{-1 d}{\varsigma c}\circ\pi 2'+d_i{}^c\exp(-2^1\circ\pi d)$. If $c < d$, set $\theta_L = (\mu_L \ln n)^{1-d}$; $0 < \mu_L < \frac{1}{2^1}$ and $\theta_H = (\mu_H \ln n)^{1-c}$; $0 < \mu_H < \frac{1}{2}$, then for $MSE(\theta_L) \frac{1}{4}$ abias$^2$ we have $\iota_L = 2^1\mu_L$; $MSE(\theta_H) \frac{1}{4}$ avar we have $\iota_L = 1 - 2\mu_H$. As $\mu_L$ ($\mu_H$) is selected close to $\frac{1}{2^1}$ (0) we approach $\iota = 1$ and (a), (b), (c), and (d) follow. If $c > d$, $\theta_L = (\mu_L \ln n)^{1-c}$; $0 < \mu_L < \frac{1}{2}$ and $\theta_H = (\mu_H \ln n)^{1-d}$; $0 < \mu_H < \frac{1}{2^1}$. We get $\iota(L) = \iota(H) = 0$: Since here the expression for $MSE(\theta^*)$ for $\theta^* > \theta_L$ grows faster than $n^{-1}\exp(\varsigma\mu_L \ln n) = n^{-1+\varsigma\mu_L}$ for any $\mu_L$, $\iota = 0$. Thus (a), (b), (c) and (d) follow. For $c = d$, set $\theta_L = (\mu_L \ln n)^{1-c}$; $0 < \mu_L < \frac{1}{\varsigma+2^1}$ and $\theta_H = (\mu_H \ln n)^{1-c}$; $\frac{1}{\varsigma+2^1} < \mu_H < \frac{1}{2}$. In the expression for $MSE(\theta^*)$ with $\theta^* > \theta_L$ we have $\iota < 1 - \mu_L\varsigma$ for all $0 < \mu_L < \frac{1}{\varsigma+2^1}$, therefore considering $\mu_L$ arbitrarily close to $\frac{1}{\varsigma+2^1}$ we get $\iota = \frac{2^1}{\varsigma+2^1}$: Thus (a), (b), (c), and (d) follow. If $\varsigma\circledR = 2'$ an analytic solution $\theta^* = (\frac{1}{\varsigma+2^1})^{1-c}(\ln n + \ln(\frac{2^1 b}{\varsigma a}))^{1-c}$ obtains.

This completes the proof. $\blacksquare$

# References

Andrews, D.W.K. (1991): "Asymptotic Normality of Series Estimators for Nonparametric and Semiparametric Regression Models," Econometrica, 59, 307–345.

Andrews, D.W.K. and M.M.A. Schafgans (1998): "Semiparametric Estimation of the Intercept of a Sample Selection Model," Review of Economic Studies, 65, 497–518.

Buchinsky, M. (1998): "The Dynamics of Changes in the Female Wage Distribution in the USA: A Quantile Regression Approach," Journal of Applied Econometrics, 13, 1–30.

Danielsson, J. and C.G. de Vries (1997): "Beyond the Sample: Extreme Quantile and Probability Estimation," unpublished manuscript.

Gallant, R. and D. Nychka (1978): "Semi-Nonparametric Maximum Likelihood Estimation," Econometrica, 55, 363–390.

Gradshteyn, I.S. and I.M. Ryzhik (1994): Table of Integrals, Series, and Products, ed. by A. Je¤rey. Translated from the Russian by Scripta Technica, Inc. London: Academic Press. Fifth Edition.

Han, A. K. (1987): "Non-parametric Analysis of a Generalized Regression Model," Journal of Econometrics, 35, 303–316.

Heckman, J.J. (1990): "Varieties of Selection Bias," American Economic Review, 80, 2, 313–318.

Hosking, J., J. Wallis, and E. Wood (1985): "Estimation of the Generalized Extreme-Value Distribution by the Method of Probability-Weighted Moments," Technometrics, 27, 251–261.

Huisman, R., K. Koedijk, C. Kool, and F. Palm (1997): "Fat Tails in Small Samples," unpublished manuscript.

Ichimura, H. (1993): "Semiparametric Least Squares (SLS) and Weighted SLS Estimation of Single Index Models," Journal of Econometrics, 58, 71-120.

Ichimura, H. and L.F. Lee (1990): "Semiparametric Least Squares Estimation of Multiple Index Models: Single Equation Estimation," in Nonparametric and Semiparametric Estimation Methods in Econometrics and Statistics, ed. by W.A. Barnett, J. Powell, and G.E. Tauchen. Cambridge: Cambridge University Press.

Klein, R.W. and R.H. Spady (1993): "An E¢cient Semiparametric Estimator for Binary Response Models," Econometrica, 61, 387–421.

Lee, L.-F. (1982): "Some Approaches to the Correction of Selectivity Bias," Review of Economic Studies, XLIX, 355–372.

Newey, W. (1988): "Two Step Series Estimation of Sample Selection Models," unpublished manuscript, Department of Economics, Princeton University.

Powell, J. (1989): "Semiparametric Estimation of Censored Selection Models," unpublished manuscript, University of Wisconsin-Madison.

Powell, J., J. Stock, and T. Stoker (1989): "Semiparametric Estimation of Index Coe¢cients," Econometrica, 57, 1435–1460.

Robinson, P.M. (1988): "Root-N-Consistent Semiparametric Regression," Econometrica, 56, 931–954.

Schafgans, M.M.A. (1998a): "Ethnic Wage Di¤erences in Malaysia: Parametric and Semiparametric Estimation of the Chinese–Malay Wage-Gap", Journal of Applied Econometrics, 13, 481–504.

_____ (1998b):"A Monte Carlo Study of the Semiparametric Estimation of the Intercept of a Sample Selection Model," unpublished manuscript.