

Honoré, Bo E.; Hu, Luojia

**Working Paper**

## Selection without exclusion

Working Paper, No. 2018-10

**Provided in Cooperation with:**

Federal Reserve Bank of Chicago

*Suggested Citation:* Honoré, Bo E.; Hu, Luojia (2018) : Selection without exclusion, Working Paper, No. 2018-10, Federal Reserve Bank of Chicago, Chicago, IL, <https://doi.org/10.21033/wp-2018-10>

This Version is available at:

<https://hdl.handle.net/10419/200580>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



Federal Reserve Bank of Chicago

## **Selection Without Exclusion**

*Bo E. Honoré and Luojia Hu*

July 2, 2018

WP 2018-10

<https://doi.org/10.21033/wp-2018-10>

*\*Working papers are not edited, and all opinions and errors are the responsibility of the author(s). The views expressed do not necessarily reflect the views of the Federal Reserve Bank of Chicago or the Federal Reserve System.*

# Selection Without Exclusion\*

Bo E. Honoré<sup>†</sup>      LuoJia Hu<sup>‡</sup>

July 2, 2018

## Abstract

It is well understood that classical sample selection models are not semi-parametrically identified without exclusion restrictions. Lee (2009) developed bounds for the parameters in a model that nests the semiparametric sample selection model. These bounds can be wide. In this paper, we investigate bounds that impose the full structure of a sample selection model with errors that are independent of the explanatory variables but have unknown distribution. We find that the additional structure in the classical sample selection model can significantly reduce the identified set for the parameters of interest. Specifically, we construct the identified set for the parameter vector of interest. It is a one-dimensional line-segment in the parameter space, and we demonstrate that this line segment can be short in principle as well as in practice. We show that the identified set is sharp when the model is correct and empty when model is not correct. We also provide non-sharp bounds under the assumption that the model is correct. These are easier to compute and associated with lower statistical uncertainty than the sharp bounds. Throughout the paper, we illustrate our approach by estimating a standard sample selection model for wages.

Key Word: Sample Selection, Exclusion Restrictions, Bounds, Partial Identification.

JEL Code: C10, C14.

---

\*This research was supported by the Gregory C. Chow Econometric Research Program at Princeton University and by the National Science Foundation (Grant Number SES-1530741). The opinions expressed here are those of the authors and not necessarily those of the Federal Reserve Bank of Chicago or the Federal Reserve System. Leah Plachinski and Rachel Anderson provided excellent research assistance. We thank participants at numerous seminars for helpful comments. The most recent version of this paper can be found at <http://www.princeton.edu/~honore/papers/SelectionWithoutExclusion.pdf>.

<sup>†</sup>Mailing Address: Department of Economics, Julis Romo Rabinowitz Building, Princeton, NJ 08544. Email: [honore@Princeton.edu](mailto:honore@Princeton.edu).

<sup>‡</sup>Mailing Address: Economic Research Department, Federal Reserve Bank of Chicago, 230 S. La Salle Street, Chicago, IL 60604. Email: [luhu@frbchi.org](mailto:luhu@frbchi.org).

# 1 Introduction

This paper considers identification in the classical sample selection model (Heckman (1976))

$$y_i^* = x_i' \beta + \varepsilon_i, \tag{1}$$

where  $y_i = y_i^*$  is observed if  $w_i' \gamma + \nu_i > 0$ . Early applications of the model assumed  $(\varepsilon_i, \nu_i)$  is independent of  $(x_i, w_i)$ , and distributed according to a bivariate normal distribution where both means are 0 and the variance of  $\nu_i$  is 1. This allows one to estimate  $\beta$  (and  $\gamma$ ) by maximum likelihood or by a two-step procedure. See Heckman (1979). Powell (1987) and others later considered semiparametric estimation of  $\beta$  under the assumption that  $(\varepsilon_i, \nu_i)$  is independent of  $(x_i, w_i)$  but without the normality assumption. See, for example, Powell (1994). The key identifying assumption is that  $x_i$  must have full rank conditional on  $w_i' \gamma$ . This is essentially an exclusion restriction that requires that  $w_i$  include variables that do not enter in  $x_i$ . Das, Newey, and Vella (2003) make a similar exclusion restriction assumption in a much more nonparametric setting.<sup>1</sup>

In this paper, we address the question of how much can be learned without an exclusion restriction like the one assumed in the literature discussed above. We consider this important because it is often difficult to find variables that both matter for selection and can be credibly excluded from the main equation. For example, Krueger and Whitmore (2001) assumed normality and wrote “Identification in these models is based on the assumption of normal errors, as there is no exclusion restriction.” Lee (2009) and Krueger and Whitmore (2001) considered set identification in a sample selection model which contains (1) as a special case<sup>2</sup>. Unfortunately, these sets are often too large to be informative. For example, Barrow and Rouse (2017)

---

<sup>1</sup>Escanciano, Jacho-Chvez, and Lewbel (2016) considered an identified sample selection model in which identification is essentially driven by nonlinearity. We consider our paper a complement to theirs.

<sup>2</sup>Manski (1989) constructed bounds in a model that is neither more general nor more restrictive than our setting. Blundell, Gosling, Ichimura, and Meghir (2007) also constructed bounds in a sample selection model, but in a much more nonparametric setting than the one considered here.

wrote “Unfortunately, Lee Bounds estimates (Lee, 2009) are quite wide and largely uninformative.” This is the motivation for this paper.

We first gain insights by studying the simplest case where the only explanatory variable is binary (Section 3). We demonstrate that in that model, the identified region for the parameter of interest can be quite small, and we provide conditions under which the upper or lower limits of the bound for the parameter coincide with the true parameter value. These results are then generalized to a model with a single potentially non-binary explanatory variable.

We next study the sample selection model with a more general set of explanatory variables in Section 4. We show that in this case, the identified set is one-dimensional. This observation is also implicit in Chamberlain (1986). Combining this insight with the results from Section 3, we then construct the identified set for the parameter vector. We show that if the model is correctly specified, our constructed identified region is sharp, and that it is necessarily empty when the model is misspecified.

The population version of the identified set for  $\beta$  can be small enough to be empirically interesting. However, the characterization of the sharp identified set for  $\beta$  relies heavily on the tail behavior of the distribution of  $y_i$  (conditional on selection). This will make estimation of the set based on a sample analog unattractive. We will therefore propose estimators of slightly larger sets.

Throughout the paper, we illustrate our approach by estimating a classical sample selection model for wages. We introduce this application in Section 2 and expand the analysis throughout the paper.

NOTATIONAL NOTE: Throughout this paper, we use  $f$  with a subscript letter to denote the density of that variable. If a variable,  $y$ , is subject to sample selection, i.e., it is observed with probability less than 1,  $f_y$  will integrate to the probability that  $y_i$  is observed. For the unobserved error terms,  $\varepsilon_i$  and  $\nu_i$ ,  $f_\varepsilon$  and  $f_\nu$  denote the underlying densities and they each integrate to 1.

## 2 Empirical Illustration: Wages and Ethnicity

Throughout the paper, we use a simple sample selection model for log-wages to illustrate our approach. The emphasis will be on the effect of ethnicity on wages. Inspired by Mora (2008), we investigate the wage-differential between third-generation Mexican-Americans and other Americans after controlling for sample selection.

Like Mora (2008), we use CPS data on wages from Arizona, California, New Mexico and Texas. Our data spans the years 2003 to 2016, and contains 129,907 women, of whom 26,698 are third-generation Mexican-Americans and 103,209 are non-Hispanic whites. There are 118,418 men. Of them, 21,402 are third-generation Mexican-Americans and 97,016 are non-Hispanic whites. For women, the percentage working is 64% for third-generation Mexican-Americans and 61% for non-Hispanic whites. For men, the shares are 71% and 67%, respectively.

Summary statistics are provided in Table 1. Appendix 2 provides details about the data.

## 3 Simplest Case: Single Regressor

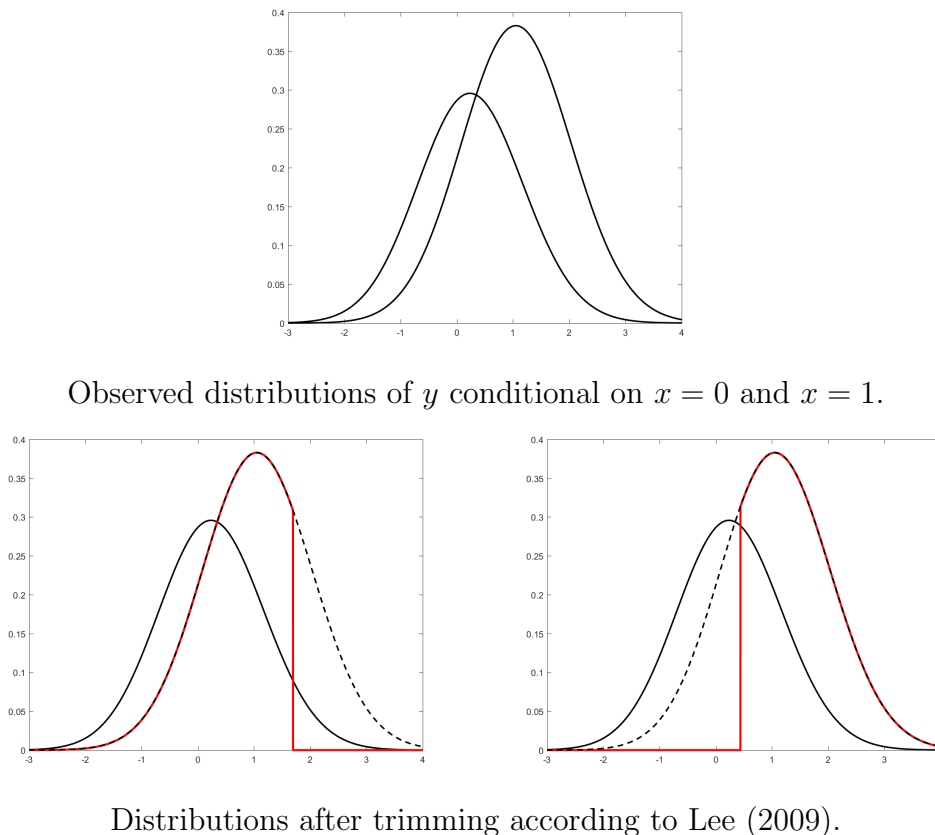
Consider first the simple case with a scalar binary explanatory variable,  $x_i$ :

$$y_i^* = x_i\beta + \varepsilon_i \tag{2}$$

where  $y_i = y_i^*$  is observed if  $x_i + v_i > 0$  and  $(\varepsilon_i, \nu_i)$  is independent of  $x_i$ . When  $x_i + v_i \leq 0$ ,  $y_i^*$  is not observed and  $y_i$  is undefined. The coefficient on the sample selection equation is only identified up to scale and its sign is identified. There is therefore no loss of generality by assuming that the coefficient on  $x_i$  is 1. Since point-mass or limited support of the distribution of  $\varepsilon_i$  generally help with identification, we assume that  $\varepsilon_i$  is continuously distributed with full support conditional on  $\nu_i$ . We will implicitly assume random sampling, and occasionally drop the subscript  $i$  to aid readability.

Lee (2009) considers a more general sample selection model in which both the distribution of  $y_i$  and the probability of selection depend on selection in a nonparametric manner and the object of interest is the average effect of the “treatment”  $x_i$ . His main assumption is a monotonicity assumption that requires that any individual who is selected into the sample when  $x_i = 0$  would also have been selection in a counterfactual scenario where  $x_i = 1$ . As such, he essentially considered the same model, but with (2) replaced by  $y_i = h(x_i, \varepsilon_i)$  for some unknown  $h$ . The average treatment effect,  $E[y_i^* | x_i = 1] - E[y_i^* | x_i = 0]$ , is not identified in this case, but Lee constructed the sharp identified set for the parameter  $E[y_i^* | x_i = 1, s_i] - E[y_i^* | x_i = 0, s_i]$  where  $s_i$  is the event that  $y_i$  would be observed whether  $x_i = 0$  or  $x_i = 1$ . Lee’s bounds are illustrated graphically in Figure 1 for a data-generating process with  $(\varepsilon_i, \nu_i)'$  distributed according to a bivariate normal distribution,  $\beta = 1$ ,  $E[\varepsilon_i] = 0$ ,  $E[\nu_i] = \frac{1}{2}$ ,  $V[\varepsilon_i] = 1$ ,  $V[\nu_i] = 1$  and  $\text{cov}(\varepsilon_i, \nu_i) = \frac{1}{2}$ . See also Example 1 below.

Figure 1: Construction of Lee Bounds for Data-Generating Process in Example 1



The model considered here implies the monotonicity assumption in Lee (2009). If  $y_i$  is observed when  $x_i = 0$ , then  $v_i$  must be greater than 0, as a result,  $y_i$  will also be observed for the same draw of  $\nu_i$  when  $x_i = 1$ . Hence the sample selection model (2) is the version of Lee’s setup in which the treatment effect is constant, and Lee’s bounds can be thought of as non-sharp bounds on  $\beta$ .

It is useful to define a binary variable for whether  $y_i$  is observed,  $d_i = 1 \{x_i + v_i > 0\}$ . For all  $c_1 < c_2$ , we then have

$$P(c_1 < \varepsilon_i \leq c_2, d_i = 1 | x_i = 0) \leq P(c_1 < \varepsilon_i \leq c_2, d_i = 1 | x_i = 1) \quad (3)$$

or

$$P(c_1 < y_i \leq c_2, d_i = 1 | x_i = 0) \leq P(c_1 < y_i - \beta \leq c_2, d_i = 1 | x_i = 1) \quad (4)$$

Equation (3) is illustrated in Figure 2 using the same data-generating process as above. The contour plot in the left panel shows the joint distribution of  $(\varepsilon_i, v_i)$  before selection and the solid line in the right panel depicts the corresponding marginal distribution of  $\varepsilon_i$ . The selection implies that  $y_i^*$  is not observed when  $v_i \leq -x_i$ . For  $x_i = 0$  and  $x_i = 1$ , this means that we “lose” the errors below the solid lines in the left panel of Figure 2. The dashed lines in the right panel of Figure 2 show the “density” of the remaining  $\varepsilon$ ’s. These densities integrate to the probability that  $y_i^*$  is observed conditional on  $x_i$ .

The restriction (4) can be expressed in terms of the density of the observed  $y$  conditional on  $x_i$ . Define

$$f_y(c | x_i) = f_{y^*}(c | x_i) P(d_i = 1 | y_i = c, x_i).$$

This is the “density” of the observed  $y$ , except that it does not integrate to 1 because  $y$  is not observed when  $d = 0$ .

With this notation, (4) can be expressed as<sup>3</sup>  $f_y(c | x_i = 0) \leq f_y(c + \beta | x_i = 1)$  for

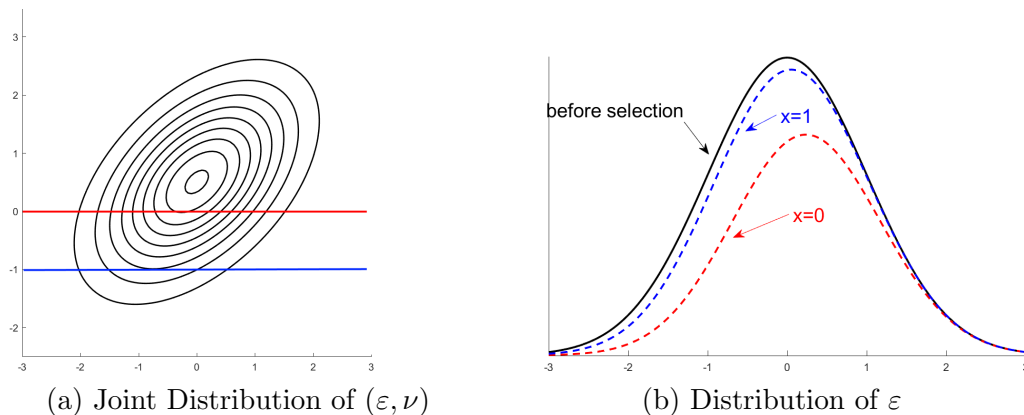
---

<sup>3</sup>This idea is reminiscent of the insight in Kitagawa (2015) who creates a test for instrument validity based on whether one product of a density and a probability lies above a second product of a density and a probability at all points. To map our insight into his, we would have to think of



all values of  $c$ . The following theorem establishes that this inequality contains all the available information. As a result it can be used to construct the identified region for  $\beta$ .

Figure 2: Distribution of  $\varepsilon$  Before and After Selection



**Theorem 1** Let  $x_i$  be a scalar binary random variable and let  $(\varepsilon_i, \nu_i)$  be independent of  $x_i$ . If  $y_i^* = x_i\beta + \varepsilon_i$  and if  $y_i = y_i^*$  is observed when  $d_i = 1 \{x_i + \nu_i > 0\}$  equals one, then the identified region for  $\beta$  is

$$\mathbf{B} = \{b \in \mathbb{R} : f_y(c | x_i = 0) \leq f_y(c + b | x_i = 1) \text{ for all values of } c\}$$

provided that  $P(d_i = 1 | x_i = 1) > 0$ .

**Proof.** This is a special case of Theorem 3 below. The proof here is more readable. The discussion above established that the true  $\beta$  belongs to  $\mathbf{B}$ . We will now argue that for any  $b$  in  $\mathbf{B}$ , there exists a joint distribution<sup>4</sup>  $\tilde{f}$  of  $(\varepsilon, \nu)$  such that  $(\tilde{f}, b)$  will be consistent with the observed  $f_y(\cdot | x = 1)$  and  $f_y(\cdot | x = 0)$ . First, define the marginal distribution of  $\varepsilon$  by

$$\tilde{f}_\varepsilon(c) = \frac{f_y(c + b | x = 1)}{P(d = 1 | x = 1)}.$$

---

(a) his outcome as our  $y - x\beta$ , (b) his instrument as our  $x$ , and (c) his treatment as our selection dummy.

<sup>4</sup>For ease of exposition, we have dropped the subscript  $i$  in the proof.

Next, define the conditional cumulative distribution function of  $\nu$  given  $\varepsilon$  at  $-1$  and  $0$  by

$$\tilde{F}_\nu(-1|\varepsilon) = 1 - P(d = 1|x = 1)$$

and

$$\tilde{F}_\nu(0|\varepsilon) = 1 - \frac{f_y(\varepsilon|x=0)}{f_y(\varepsilon+b|x=1)}P(d = 1|x = 1),$$

respectively. Since  $\frac{f_y(\varepsilon|x=0)}{f_y(\varepsilon-b|x=1)} \leq 1$ ,  $\tilde{F}_\nu(-1|\varepsilon) \leq \tilde{F}_\nu(0|\varepsilon)$ , the resulting  $\tilde{F}_\nu(\cdot|\varepsilon)$  is a CDF. It does not matter what the values of  $\tilde{F}(\cdot|\varepsilon)$  are at points other than  $0$  and  $-1$ .

With this  $(\tilde{f}, b)$

$$\begin{aligned} \tilde{f}_y(c|x=1) &= \tilde{f}_\varepsilon(c-b|x=1)\tilde{P}(d=1|y=c, x=1) \\ &= \tilde{f}_\varepsilon(c-b|x=1)\left(1 - \tilde{F}_\nu(-1|\varepsilon=c-b)\right) \\ &= \frac{f_y(c|x=1)}{P(d=1|x=1)}P(d=1|x=1) = f_y(c|x=1) \end{aligned}$$

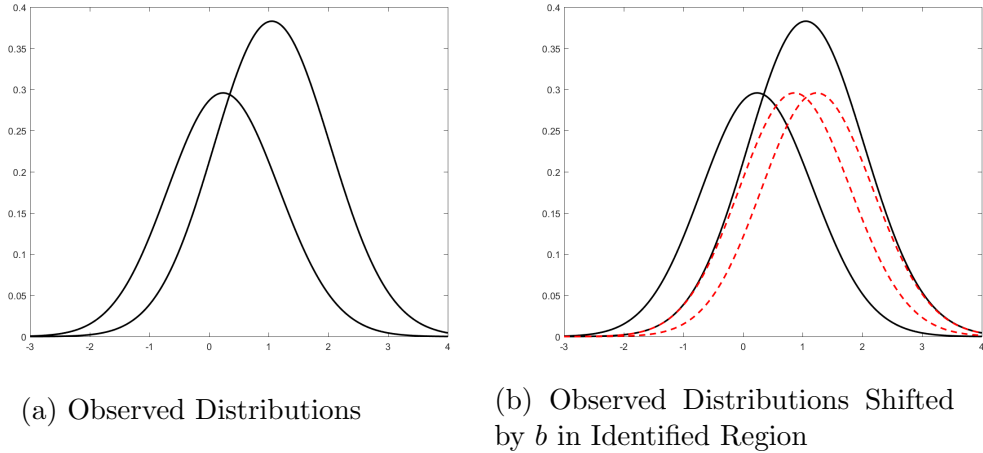
and

$$\begin{aligned} \tilde{f}_y(c|x=0) &= \tilde{f}_\varepsilon(c|x=0)\tilde{P}(d_i=1|y=c, x=0) \\ &= \tilde{f}_\varepsilon(c|x=0)\left(1 - \tilde{F}_\nu(0|\varepsilon=c)\right) \\ &= \frac{f_y(c+b|x=1)}{P(d=1|x=1)}\frac{f_y(c|x=0)}{f_y(c+b|x=1)}P(d=1|x=1) = f_y(c|x=0). \end{aligned}$$

In this construction, the distribution of  $\varepsilon$  conditional on  $\nu$  is continuous. This proves the theorem. ■

The construction of the identified region in Theorem 1 is illustrated graphically in Figure 3. The left side of Figure 3 shows the density of the observed  $y$  conditional on  $x$  multiplied by the conditional probability of selection for  $x_i$  equal to  $0$  and  $1$ . Theorem 1 characterized the identified region for  $\beta$  as the length of the horizontal shifts of one of the curves that will result in one of the curves being above the other. This is illustrated in the right hand side of Figure 3.

Figure 3: Our Bounds



**Example 1** Let  $(\varepsilon_i, \nu_i)'$  be distributed according to a bivariate normal distribution with  $\beta = 1$ ,  $E[\varepsilon_i] = 0$ ,  $E[\nu_i] = \frac{1}{2}$ ,  $V[\varepsilon_i] = 1$ ,  $V[\nu_i] = 1$  and  $\text{cov}(\varepsilon_i, \nu_i) = \frac{1}{2}$ . With these  $P(d_i = 1 | x_i = 0) = 0.691$  and  $P(d_i = 1 | x_i = 1) = 0.933$ . This is the situation depicted in Figure 3 and the identified region for  $\beta$  is  $[0.626, 1.00]$ . In contrast, the Lee bounds are<sup>5</sup>  $[0.389, 1.238]$ .

To estimate the identified region characterized by in Theorem 1, one needs to estimate the density of the observed  $y_i$  conditional on  $x_i$  for  $x_i = 0$  and  $x_i = 1$ , and then compare these densities for all values of the argument. This is troublesome because the densities will typically both be close to 0 in the tails, and small estimation errors will have a big effect on which one takes on the larger value. Moreover, the densities are estimated at a nonparametric rate in general. This suggests constructing an identified region by exploring (4) for a finite number of pairs of  $(c_1, c_2)$ . For example, one could calculate the deciles of the observed  $y$  conditional on  $x = 0$  and then use  $(c_1, c_2) = (d_{j-1}, d_j)$  for  $j = 1, \dots, 10$ , where  $d_j$  is the  $j$ -th decile,  $d_0 = -\infty$  and  $d_{10} = \infty$ .

**Example 2** (*Example 1 continued*) In this setup the crude bounds described above are  $[0.609, 1.025]$ .

---

<sup>5</sup>To calculate the bounds, we use equation (5) in Muthén (1990) after correcting a typo in the second line (the next to last subscript- $i$  should be subscript- $j$ ).

In Example 1, the upper bound of the identified set equals the true  $\beta$ . This is true in general when the true (unknown) distribution of the errors is a bivariate normal with positive correlation.

**Proposition 1** *When the distribution of the errors is bivariate normal with positive correlation, the upper limit of the identified region is the true parameter value. When the correlation is negative, the lower limit of the identified region is the true value. With no selection, i.e., independence of the errors, the identified region is the true value.*

**Proof.** See Appendix 1. ■

The proof of Proposition 1 is driven by the tail behavior of the normal distribution. As a result, the proposition can be generalized to

**Proposition 2** *Suppose that the density of  $\varepsilon$  has sufficiently thin tails that for  $a > 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow \infty$  and for  $a < 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow -\infty$ . Then*

1. *If the distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  whenever  $c_1 > c_2$ , then the upper limit of the identified region is the true value.*
2. *If the distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  whenever  $c_1 < c_2$ , then the lower limit of the identified region is the true value.*
3. *If  $\nu$  and  $\varepsilon$  are independent, then the identified region is the true value.*

**Proof.** See Appendix 1. ■

The assumption on the tail behavior of the marginal distribution of  $\varepsilon$  is slightly stronger than log-concavity and is implied by tail-behavior of the form  $\exp(-ax^\gamma)$  for  $a > 0$  and  $\gamma > 1$ . We interpret the stochastic dominance assumption in 1 as positive selection: larger values of  $\varepsilon$  are associated with higher probability of selection.

Likewise, we interpret the stochastic dominance assumption in 2 as negative selection. The approach in Proposition 2 is different from, but similar in spirit to, the approach in Heckman (1990) and Andrews and Schafgans (1998). Both approaches rely on “identification at infinity,” but while Heckman (1990) and Andrews and Schafgans (1998) need an exclusion restriction and rely on extreme values of the selection index, we do not need exclusion restrictions and rely on extreme values of the outcome variable.

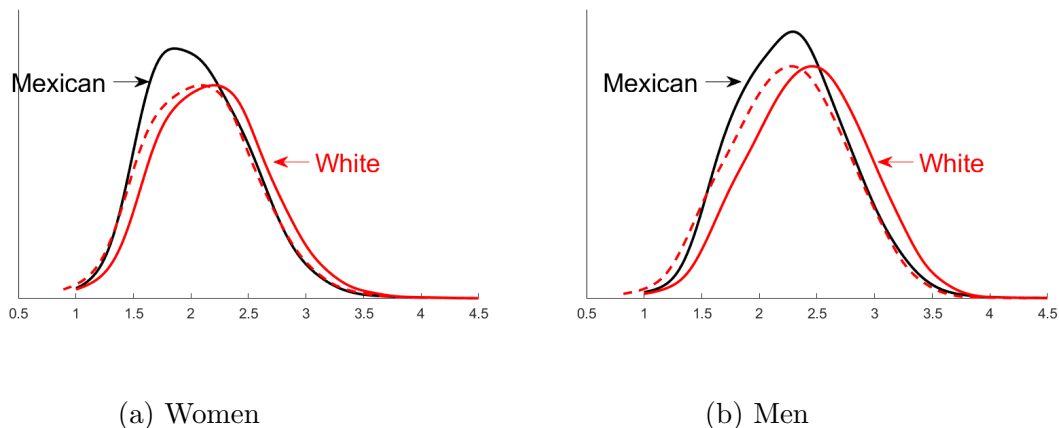
### 3.1 Empirical Illustration Part 1

Using the data described in Section 2, we plot the “densities” (the product of the density conditional on selection and the probability of selection) of log-wages for Mexican-Americans and Non-Hispanic white Americans by gender. We restrict the sample to individuals whose highest degree is high school. These are depicted by the solid lines in Figure 4. The areas under the Mexican-American curves are larger than for Non-Hispanic white Americans because the former are more likely to work for pay.

The dashed lines are the curves for the Non-Hispanic whites shifted by  $-0.11$  for women and by  $-0.18$  for men. The shifted curves for whites almost fit under the curves for the Mexican-Americans. In this case, the Lee bounds for the log-wage differentials between Mexican-Americans and Non-Hispanic white Americans case are  $(-0.210, -0.041)$  for women and  $(-0.249, -0.074)$  for men while the difference in means are  $-0.123$  for women and  $-0.162$  for men.

Heckman’s two-step estimator exploits variation in the conditional mean of the dependent variable. When the only explanatory variable is binary, there will be perfect collinearity between it and the sample selection correction term. The procedure therefore cannot be applied. In contrast, the maximum likelihood estimator for the log-wage differentials between Mexican-Americans and Non-Hispanic white Americans exploits information from the entire distribution of the dependent variable, and this estimator is therefore in principle applicable, although it is likely to be fragile.

Figure 4: Shifted and Unshifted Log-Wage Distributions



For example, when we tried to estimate the model using Stata<sup>6</sup>, the routine failed to converge for the men and located a false minimum for the women. Our own Matlab program estimated the coefficient on Mexican-American to be  $-0.174$  for women and  $-0.208$  for men.

### 3.2 Single Non-Binary Regressor

Theorem 1 applies to the case where  $x_i$  is binary. When  $x_i$  is not binary and unbounded from above, identification at infinity arguments like those in Andrews and Schafgans (1998) and Heckman (1990) yield point-identification of  $\beta$ . We therefore focus on the case where  $x_i$  is bounded from above.

When  $x_i$  is not binary and bounded from above, applying (4) to all pairs of values in the support of  $x_i$  yields bounds on the identified region of  $\beta$ . The following theorem establishes that the intersection of these bounds is sharp.

**Theorem 2** *Let  $(x_i, \varepsilon_i, \nu_i)$  be a random vector such that  $(\varepsilon_i, \nu_i)$  is independent of  $x_i$ . If  $y_i = x_i\beta + \varepsilon_i$  and if  $y_i = y_i^*$  is observed when  $d_i = 1 \{x_i + \nu_i > 0\}$  equals one and*

---

<sup>6</sup>More precisely the routine `heckman` in Stata Version 14 with all the default options.

the upper bound on the support of  $x_i$  is  $x_{\max}$ , then the identified region for  $\beta$  is <sup>7</sup>

$$\mathbf{B} = \{b \in \mathbb{R} : f_{y_i|x_i}(c + x_i b | \xi_1) \leq f_{y_i|x_i}(c + x_i b | \xi_2) \\ \text{for all values of } c \text{ and } \xi_1 < \xi_2 \text{ in the support of } x_i\}.$$

**Proof.** Follows from Theorem 3 below. ■

As discussed above, the identified set can also be expressed in terms of probabilities,

$$\mathbf{B} = \{b \in \mathbb{R} : P(c_1 < y_i - x_i b \leq c_2, d_i = 1 | x_i = \xi_1) \leq P(c_1 < y_i - x_i b \leq c_2, d_i = 1 | x_i = \xi_2) \\ \text{for all values of } c_1 < c_2 \text{ and } \xi_1 < \xi_2 \text{ in the support of } x_i\}.$$

We finally note that the conclusion of Proposition 2 carries over to general distribution of  $x_i$ . Specifically,

**Proposition 3** *Let  $x_i$  be a random variable and let  $(\varepsilon_i, \nu_i)$  be independent of  $x_i$ . If  $y_i = x_i \beta + \varepsilon_i$  and if  $y_i = y_i^*$  is observed when  $d_i = 1 \{x_i + v_i > 0\}$  equals one and the upper bound on the support of  $x_i$  is  $x_{\max}$ . Suppose that the density of  $\varepsilon$  has sufficiently thin tails that for  $a > 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow \infty$  and for  $a < 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow -\infty$ . Then*

1. *If the distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  whenever  $c_1 > c_2$ , then the upper limit of the identified region for  $\beta$  is the true value.*
2. *If the distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  whenever  $c_1 < c_2$ , then the lower limit of the identified region for  $\beta$  is the true value.*
3. *If  $\nu$  and  $\varepsilon$  are independent, then the identified region for  $\beta$  is the true value.*

**Proof.** See Appendix 1. ■

---

<sup>7</sup>Recall that  $f_{y_i|x_i}$  does not integrate to 1, since  $y_i$  is not always observed.

## 4 More General Sample Selection Model

We now return to the sample selection model with a  $k$ -dimensional vector of explanatory variables,  $x_i$ ,

$$y_i^* = x_i' \beta + \varepsilon = x_{i1} \beta_1 + x_{i2}' \beta_2 + \varepsilon_i \quad (5)$$

where  $y_i = y_i^*$  is observed if  $x_i' \gamma + v_i > 0$ . When the support of  $x_i' \gamma$  is unbounded from above, identification at infinity arguments like those in Andrews and Schafgans (1998) and Heckman (1990) yield point-identification of  $\beta$ . We therefore focus on the case where  $x_i' \gamma$  is bounded from above.

To fix ideas, suppose that  $\beta_1$  is the parameter of interest.

Conditions under which  $\gamma$  is identified up to scale are well-understood; see for example Powell (1994) and the references therein. In the following, we assume that these conditions hold and that the necessary scale normalization has been imposed by normalizing the first element of  $\gamma$  to be 1. We will then write  $\gamma = (1, \gamma_2)'$  to distinguish between the variable of interest,  $x_{i1}$ , and the other explanatory variables. As in the previous section, we assume independence between  $(\varepsilon_i, v_i)$  and  $(x_{i1}, x_{i2})$  and define  $g(z_i) = E[\varepsilon_i | v_i > -z_i, x_{i1}, x_{i2}, z_i]$ . We can then write

$$y_i = x_{i1} \beta_1 + x_{i2}' \beta_2 + g(z_i) + u_i \quad (6)$$

with  $z_i = x_i' \gamma$  and  $E[u_i | x_{i1}, x_{i2}] = 0$ .

In this section, we will argue that the vector  $\beta$  is identified except for a single scale parameter. In the following subsections we will then show that bounds can be obtained for this parameter. The intuition is very simple. Suppose we knew  $\beta_1$ . We could then define  $w_i^* = y_i^* - x_{i1} \beta_1 = x_{i2}' \beta_2 + \varepsilon_i$ . The variable  $x_{i1}$  would then be excluded from the model for  $w_i^*$ . On the other hand, by normalizing the first coefficient in the selection equation to be 1, we have already assumed that  $x_{i1}$  matters for selection. Hence we have the necessary exclusion restriction, and the parameter vector,  $\beta$ , is identified except for the one-dimensional component  $\beta_1$ . Here, we give a slightly different argument because it makes the empirical implementation easier.



Following, e.g. Robinson (1988), we start by noting that (6) implies that

$$y_i - E[y_i | z_i] = (x_{i1} - E[x_{i1} | z_i])\beta_1 + (x_{i2} - E[x_{i2} | z_i])'\beta_2 + u_i \quad (7)$$

Next note that

$$\begin{aligned} (x_{i1} - E[x_{i1} | z_i]) + (x_{i2} - E[x_{i2} | z_i])'\gamma_2 &= (x_{i1} - E[x_{i1} | z_i]) + (x_{i2}\gamma_2 - E[x_{i2} | z_i]'\gamma_2) \\ &= x_i'\gamma - E[x_i'\gamma | z_i] = x_i'\gamma - E[x_i'\gamma | x_i'\gamma] = 0. \end{aligned}$$

In other words,  $(x_{i1} - E[x_{i1} | z_i]) = -(x_{i2} - E[x_{i2} | z_i])'\gamma_2$ . Equation (7) can then be written as

$$\begin{aligned} y - E[y | z_i] &= (-(x_{i2} - E[x_{i2} | z_i])'\gamma_2)\beta_1 + (x_{i2} - E[x_{i2} | z_i])'\beta_2 + u_i \quad (8) \\ &= (x_{i2} - E[x_{i2} | z_i])'(\beta_2 - \gamma_2\beta_1) + u_i \end{aligned}$$

As a result, we can identify  $\alpha_2 = (\beta_2 - \gamma_2\beta_1)$  as the coefficient in a regression of  $y_i - E[y_i | z_i]$  on  $(x_{i2} - E[x_{i2} | z_i])$ . Since  $\gamma_2$  is identified, this implies that for a given value of  $\beta_1$ ,  $\beta_2$  is identified. In other words, the identification problem is essentially one-dimensional, and bounds on  $\beta_1$  will imply bounds of the whole  $\beta$  vector.

## 4.1 Sharp Bounds

With the result of the previous section, we can write

$$\begin{aligned} y_i^* &= x_{i1}\beta_1 + x_{i2}'\beta_2 + \varepsilon_i \\ &= x_{i1}\beta_1 + x_{i2}'(\alpha_2 + \gamma_2\beta_1) + \varepsilon_i \end{aligned}$$

or

$$y_i^* - x_{i2}'\alpha_2 = (x_{i1} + x_{i2}'\gamma_2)\beta_1 + \varepsilon_i \quad (9)$$

where  $\gamma_2$  and  $\alpha_2$  are identified as above and  $y_i = y_i^*$  (and hence  $y_i - x_{i2}'\alpha_2 = y_i^* - x_{i2}'\alpha_2$ ) is observed when  $d_i = 1 \{x_{i1} + x_{i2}'\gamma_2 + \nu_i > 0\}$ . We can then apply Theorem 2 to

bound  $\beta_1$  to the region

$$\begin{aligned} \mathbf{B} &= \{b_1 \in \mathbb{R} : P(c_1 < y_i - x'_{i2}\alpha_2 - (x_{i1} + x'_{i2}\gamma_2)b_1 \leq c_2, d_i = 1 | (x_{i1} + x'_{i2}\gamma_2) = \xi_1) \\ &\leq P(c_1 < y_i - x'_{i2}\alpha_2 - (x_{i1} + x'_{i2}\gamma_2)b_1 \leq c_2, d_i = 1 | (x_{i1} + x'_{i2}\gamma_2) = \xi_2) \\ &\text{for all values of } c_1 < c_2 \text{ and } \xi_1 \leq \xi_2 \text{ in the support of } x_i \}. \end{aligned} \quad (10)$$

The identified region for the whole vector,  $\beta$ , is then the one-dimensional line segment

$$\left\{ \begin{pmatrix} b_1 \\ \alpha_2 + \gamma_2 b_1 \end{pmatrix} : b_1 \in \mathbf{B} \right\}.$$

The identified region can also be written in terms of the density of the observed data:

$$\begin{aligned} \mathbf{B} &= \{b_1 \in \mathbb{R} : f_{y|x}(c + x'_{i2}\alpha_2 + (x_{i1} + x'_{i2}\gamma_2)b_1 | x_i = \xi_1) \\ &\leq f_{y|x}(c + x'_{i2}\alpha_2 + (x_{i1} + x'_{i2}\gamma_2)b_1 | x_i = \xi_2) \\ &\text{for all values of } c \text{ and all } \xi_1 \text{ and } \xi_2 \text{ in the support of } x_i \text{ that satisfy } \xi'_1\gamma \leq \xi'_2\gamma \}. \end{aligned} \quad (11)$$

The bounds implied by  $\mathbf{B}$  are sharp by Theorem 3.

**Theorem 3** *Suppose that (i)  $(\varepsilon_i, \nu_i)$  is continuously distributed with everywhere positive density, (ii)  $(\varepsilon_i, \nu_i)$  is independent of  $x_i$ , (iii)  $E[\varepsilon_i | \nu_i > a]$  is finite for all  $a$ , (iv) there is no proper linear subspace of  $R^k$  that contains  $x_i$  with probability 1, and (v)*

$$y_i^* = x'_i\beta + \varepsilon_i$$

*is observed if  $d_i = 1 \{x'_i\gamma + \nu_i \geq 0\}$  equals one for some  $\beta$  and some  $\gamma$  with  $\gamma_1 = 1$ . If  $\gamma$  is identified and the support of  $x'_i\gamma$  is bounded from above, then  $\mathbf{B}$  is the (sharp) identified region for  $\beta_1$ .*

**Proof.** The discussion in the text established that the true  $\beta_1$  belongs to  $\mathbf{B}$  and that

$\beta_2 = \alpha_2 + \gamma_2\beta_1$ . We next argue that for any  $b = (b_1, \alpha'_2 + \gamma'_2 b_1)'$  with  $b_1$  in  $\mathbf{B}$ , there exists a joint distribution<sup>8</sup>  $\tilde{f}$  of  $(\varepsilon, \nu)$  such that the model combined with  $(\tilde{f}, b)$  will result in the same distribution as the observed  $f_{y|x}(\cdot|\xi)$  for all  $\xi$  in the support of  $x$ .

Let  $\bar{x}$  be such that  $\bar{x}'\gamma$  is the upper bound of the support of  $x'\gamma$ .

First, define the marginal distribution of  $\varepsilon$  by

$$\tilde{f}_\varepsilon(c) = \frac{f_{y|x}(c + \bar{x}'b | x = \bar{x})}{P(d = 1 | x = \bar{x})} \quad (12)$$

The definition of  $\mathbf{B}$  guarantees that this gives the same  $\tilde{f}_\varepsilon(c)$  for all choices of  $\bar{x}$  such that  $\bar{x}'\gamma$  is the upper bound of the support of  $x'\gamma$ .

With this construction, the distribution of  $\varepsilon$  is continuous. We define the conditional cumulative distribution function of  $\nu$  given  $\varepsilon$  over the support of  $-x'\gamma$ . Let  $a$  be a point in that support and let  $x_a$  be such that  $x'_a\gamma = -a$ . We then define

$$\tilde{F}_\nu(a|\varepsilon = c) = 1 - \frac{f_{y|x}(c + x'_a b | x = x_a)}{f_{y|x}(c + \bar{x}'b | x = \bar{x})} P(d = 1 | x = \bar{x}). \quad (13)$$

The assumption that  $(\varepsilon, \nu)$  has positive density guarantees that the denominator is not zero, and the definition of  $\mathbf{B}$  guarantees that the particular choices of  $x_a$  and  $\bar{x}$  do not matter as long as they satisfy  $x'_a\gamma = -a$  and  $\bar{x}'\gamma$  is the upper bound of the support of  $x'\gamma$ .

Let  $a_1 < a_2$  be two point in the support of  $-x'\gamma$  and let  $x_{a_1}$  and  $x_{a_2}$  be such that  $x'_{a_1}\gamma = -a_1$  and  $x'_{a_2}\gamma = -a_2$ . Then  $f_{y|x}(c + x'_a b | x_{a_1}) \geq f_{y|x}(c + x'_a b | x_{a_2})$  by the definition of  $\mathbf{B}$ , and the constructed  $\tilde{F}(\cdot|\varepsilon)$  is non-decreasing and therefore satisfies the requirements for a cumulative distribution function. It does not matter what the values of  $\tilde{F}(\cdot|\varepsilon)$  are at points other than those in the support of  $-x'\gamma$ .

---

<sup>8</sup>For ease of exposition, we have dropped the subscript  $i$  in the proof.

With this  $(\tilde{f}, b)$ , the model yields

$$\begin{aligned}
\tilde{f}_{y|x}(\tilde{y}|\tilde{x}) &= \tilde{f}_\varepsilon(\tilde{y} - \tilde{x}'b) \tilde{P}(d = 1 | y = \tilde{y}, x = \tilde{x}) \\
&= \tilde{f}_\varepsilon(\tilde{y} - \tilde{x}'b) \left(1 - \tilde{F}_\nu(-\tilde{x}'\gamma | \varepsilon = \tilde{y} - \tilde{x}'b)\right) \\
&= \left(\frac{f_{y|x}(\tilde{y} - \tilde{x}'b + \bar{x}'b | x = \bar{x})}{P(d = 1 | x = \bar{x})}\right) \\
&\quad \left(\frac{f_{y|x}(\tilde{y} - \tilde{x}'b + \tilde{x}'b | x = \tilde{x})}{f_{y|x}(\tilde{y} - \tilde{x}'b + \bar{x}'b | x = \bar{x})} P(d = 1 | x = \bar{x})\right) \\
&= f_{y|x}(\tilde{y}|\tilde{x}).
\end{aligned}$$

This proves the theorem. ■

Theorem 3 states that  $\mathbf{B}$  characterizes the identified region for  $\beta_1$  when the model is correct. Theorem 4 below establishes that when  $\mathbf{B}$  is not empty, the linear sample selection model cannot be rejected by the data.

**Theorem 4** *Suppose that the data-generating process for the observed distribution of  $(d_i, y_i, x_i)$  is such that (i)  $y_i$  is only observed when  $d_i = 1$ , (ii)  $P(d_i = 1 | x_i)$  can be written as a non-decreasing function of  $x_i'\gamma$  for some  $\gamma = (1, \gamma_2)'$ , (iii) the support of  $x_i'\gamma$  is bounded from above, and (iv) the density of  $y_i$  given  $x_i$  is positive everywhere for all  $x_i$  in the support of  $x_i$ .*

*If, for some vector  $\alpha_2$ ,  $\mathbf{B}$  is not empty, then for every  $\beta$  in  $\mathbf{B}$ , there exists a distribution of  $(\varepsilon_i, \nu_i)$  such that the observed distribution is the same as the one generated from a model in which  $(\varepsilon_i, \nu_i)$  is independent of  $x_i$ ,*

$$y_i^* = x_i'\beta + \varepsilon_i$$

*and  $y_i = y_i^*$  is observed if  $d_i = 1 \{x_i'\gamma + \nu_i \geq 0\}$  equals 1.*

**Proof.** Let  $b_1$  be an element of  $\mathbf{B}$ ,  $b_2 = \alpha_2 + \gamma_2 b_1$  and  $b = (b_1, b_2)'$ . We will show that we can construct a joint distribution,  $\tilde{f}$ , of  $(\varepsilon, \nu)$  such that the observed distribution of  $(y_i, x_i)$  is the same as the distribution generated from a model where  $y^* = x'b + \varepsilon$  and  $y = y^*$  is observed if  $d = 1 \{x'\gamma + \nu \geq 0\}$  equals 1.

Let  $\bar{x}$  be a vector such that  $\bar{x}'\gamma$  is the upper bound of the support of  $x'\gamma$  and define the marginal distribution of  $\varepsilon$  by

$$\tilde{f}_\varepsilon(c) = \frac{f_{y|x}(c + \bar{x}'b | x = \bar{x})}{P(d = 1 | x = \bar{x})} \quad (14)$$

The definition of  $\mathbf{B}$  guarantees that this results in the same  $\tilde{f}$  for all  $\bar{x}$ .

We next define the conditional CDF of  $\nu$  given  $\varepsilon$  at all points in the support of  $-x'\gamma$ . For a point  $a$  in the support of  $-x'\gamma$ , let  $x_a$  be such that  $-x_a'\gamma = a$ . We then define

$$\tilde{F}_{\nu|\varepsilon}(a|c) = 1 - \frac{f_{y|x}(c + x_a'b | x = x_a)}{f_{y|x}(c + \bar{x}'b | x = \bar{x})} P(d = 1 | x = \bar{x}).$$

By the definition of  $\mathbf{B}$ ,  $f_{y|x}(c + x'b | x = x_a)$  is the same for all  $x_a$  such that  $x_a'\gamma = -a$ , so we do not need to worry about the particular choice of  $x_a$ . Also, note that by assumption, the probability that  $P(d = 1 | x)$  is the same for all  $x$  such that  $x'\gamma$  is the upper bound of the support of  $x'\gamma$ .

With these distributions, the density of the observed  $y$  from the model would be

$$\begin{aligned} \tilde{f}_{y|x}(\tilde{y}|\tilde{x}) &= \tilde{f}_\varepsilon(\tilde{y} - \tilde{x}'b) \tilde{P}(d = 1 | y = \tilde{y}, x = \tilde{x}) \\ &= \tilde{f}_\varepsilon(\tilde{y} - \tilde{x}'b) \left(1 - \tilde{F}(-\tilde{x}'\gamma | \varepsilon = \tilde{y} - \tilde{x}'b)\right) \\ &= \left(\frac{f_{y|x}(\tilde{y} - \tilde{x}'b + \bar{x}'b | x = \bar{x})}{P(d = 1 | x = \bar{x})}\right) \left(\frac{f_{y|x}(\tilde{y} - \tilde{x}'b + \tilde{x}'b | x'\gamma = \tilde{x}'\gamma)}{f_{y|x}(\tilde{y} - \tilde{x}'b + \bar{x}'b | x'\gamma = \bar{x}'\gamma)} P(d = 1 | x = \bar{x})\right) \\ &= f_{y|x}(\tilde{y}|\tilde{x}). \end{aligned}$$

In words, the distribution of  $y$  given  $x$  constructed from the model coincides with the observed distribution. This completes the proof. ■

## 4.2 Non-sharp Bounds

One could in principle estimate bounds on  $\beta_1$  based on the set,  $\mathbf{B}$ , above. We do not pursue this approach because the resulting estimates would depend on the tails of nonparametrically estimated densities. In this section, we instead present non-sharp bounds based on moments that can be estimated using sample averages or

U-statistics.

The discussion above implies that for  $\xi_1 \leq \xi_2$  in the support of  $x_{i1} + x'_{i2}\gamma_2$

$$\begin{aligned} E(1 \{c_1 < y_i - x'_{i2}\alpha_2 \leq c_2, d_i = 1\} | (x_{i1} + x'_{i2}\gamma_2) = \xi_1) \\ \leq E(1 \{c_1 < y_i - x'_{i2}\alpha_2 + (\xi_1 - \xi_2)b_1 \leq c_2, d_i = 1\} | (x_{i1} + x'_{i2}\gamma_2) = \xi_2) \end{aligned} \quad (15)$$

for any  $b_1$  in the identified set. Note that  $c_1$  and  $c_2$  in (15) can depend on  $b_1$ ,  $\xi_1$ , and  $\xi_2$ .

This implies the moment inequalities

$$\begin{aligned} E[1 \{c_1 < y_i - x'_{i2}\alpha_2 - (x'_i\gamma) \beta_1 \leq c_2, d_i = 1\} | (x_{i1} + x'_{i2}\gamma_2) \in A_1] \\ \leq E[1 \{c_1 < y_i - x'_{i2}\alpha_2 - (x'_i\gamma) \beta_1 \leq c_2, d_i = 1\} | (x_{i1} + x'_{i2}\gamma_2) \in A_2] \end{aligned} \quad (16)$$

for all sets  $A_1, A_2$  where all elements in  $A_1$  are strictly below the elements in  $A_2$ . For example, the  $A$ -sets could be defined by percentiles of  $x_{i1} + x'_{i2}\widehat{\gamma}_2$ . This statement can also be stated in terms of the conditional densities.

Equation (15) also implies that for  $b_1$  in the identified set, and a pair of observations,  $i$  and  $j$ ,

$$\begin{aligned} E[(1 \{c_1 < y_j - x'_{j2}\alpha_2 + (x_{i1} - x_{j1} + (x_{i2} - x_{j2})'\gamma_2) b_1 \leq c_2, d_j = 1\} \\ - 1 \{c_1 < y_i - x'_{i2}\alpha_2 \leq c_2, d_i = 1\}) 1 \{x_{i1} + x'_{i2}\gamma_2 < x_{j1} + x'_{j2}\gamma_2\}] \geq 0 \end{aligned} \quad (17)$$

for all values of  $c_1 < c_2$ , where  $c_1$  and  $c_2$  can depend on  $b_1$ ,  $x_i$  and  $x_j$ .

The moment inequalities (17) and (16) can be used to estimate non-sharp bounds for  $\beta$  in (5).

The following example suggests that the size of the identified set is likely to be small enough to be useful.

**Example 3** *Consider the data generating process*

- $(\nu_i, \varepsilon_i)$  bivariate normal with  $V[\nu_i] = 1$ ,  $V[\varepsilon_i] = 2$ ,  $\text{cov}(\nu_i, \varepsilon_i) = 1$ ,  $E[\nu_i] = \frac{1}{2}$  and  $E[\varepsilon_i] = 0$ .

- $x_{ik} = U_{ik} + Z_i$  for  $k = 1, 2, 3$ , where  $U_{ik} \sim U\left(-\frac{1}{2}, \frac{1}{2}\right)$  and  $Z_i \sim N\left(0, \frac{1}{25}\right)$  (all independent)
- $\beta = (1, 1, 1)'$  and  $\gamma = (0.45, 0.55, 0, 55)$  (before normalization)

We calculate the (non-sharp) identified region for  $\beta_1$  based on (16), the  $A$ 's based on quintiles of  $x_{j1} + x'_{j2}\gamma_2$  and  $c_1$  and  $c_2$  adjacent deciles of  $y_i - x'_{i2}\alpha_2 - (x_{i1} + x'_{i2}\gamma_2) b_1$  to be (0.658, 1.003). When we decreased the number of inequalities by only considering  $A_1 = (-\infty, \text{median}(x_{j1} + x'_{j2}\gamma_2))$  and  $A_2 = (\text{median}(x_{j1} + x'_{j2}\gamma_2), \infty)$  and  $c_1$  and  $c_2$  adjacent quintiles of  $y_i - x'_{i2}\alpha_2 - (x_{i1} + x'_{i2}\gamma_2) b_1$ , the (non-sharp) identified region for  $\beta_1$  increased to (0.529, 1.031).

We calculate the (non-sharp) identified region for  $\beta_1$  based on (17) and  $c_1$  and  $c_2$  adjacent deciles of  $y_i - x'_{i2}\alpha_2$  (conditional on  $d_i = 1$ ) to be (0.491, 1.009).

By comparison, the 5<sup>th</sup> and 95<sup>th</sup> percentiles of Heckman's two-step estimator for  $\beta_1$  based on 1,000 observations from this design are 0.332 and 1.714.<sup>9</sup>

### 4.3 Empirical Illustration Part 2

To investigate the usefulness of the approach from Section 4.2 in empirical settings, we return to the question in Section 2. In this application, the parameter of interest is the coefficient on being third-generation Mexican-American as opposed to non-Hispanic white. The other explanatory variables are age, age-squared, experience, experience-squared, education dummies (less than high school, some college, college, and advanced degree, with high school as the omitted category), dummies for being a veteran and being married, state dummies and year dummies.

We first estimate the model under the assumption of joint normality of the errors, using both maximum likelihood estimator and Heckman's two-step estimator. The estimation results are presented in Table 2.<sup>10</sup> To implement the idea in 4.2, we define

---

<sup>9</sup>The bounds based on (16) are calculated using a sample with 100,000,000 observations, while the bounds based on (17) are calculated using 100,000 observations. We use a smaller number for (17) because it is based on all pairs of observations. The percentiles of Heckman's two-step estimator are calculated in Matlab by Monte Carlo using 100,000 replications.

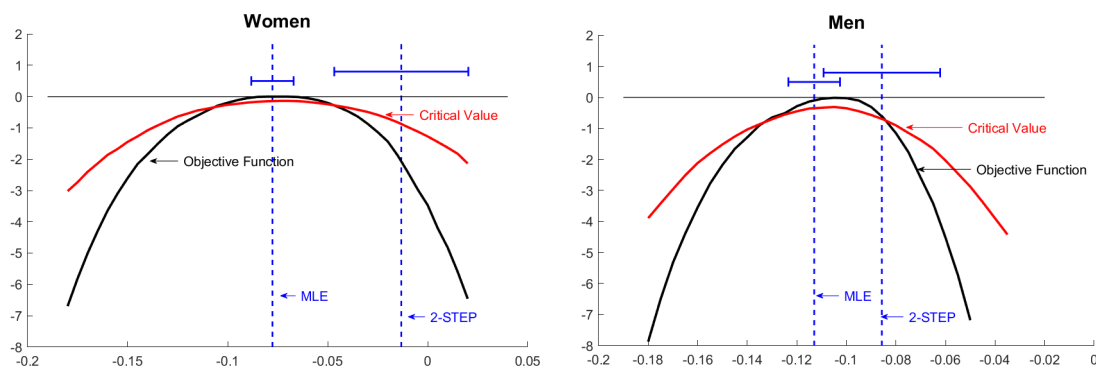
<sup>10</sup>Table 2 was produced using Stata.

a sample analog of the solutions to the population inequalities in (16) by the solutions to  $Q_{1n}(b_1) = 0$ , where

$$Q_{1n}(b_1) = - \sum_{\ell, k} \max \left\{ \widehat{E} [1 \{c_\ell < y_i - x'_{i2} \widehat{\alpha}_2 - (x'_i \widehat{\gamma}) b_1 \leq c_{\ell+1}, d_i = 1\} | x'_i \widehat{\gamma} \in A_k] - \widehat{E} [1 \{c_\ell < y_i - x'_{i2} \widehat{\alpha}_2 - (x'_i \widehat{\gamma}) b_1 \leq c_{\ell+1}, d_i = 1\} | x'_i \widehat{\gamma} \in A_{k+1}], 0 \right\}^2$$

Figure 5 displays the objective function and the 5%-critical value function calculated using sub-sampling (see Canay and Shaikh (2017)) with sub-sample size equal to 15,000 and 1,000 sub-samples. The parameter  $\gamma$  is estimated by logit maximum likelihood<sup>11</sup> and  $\alpha_2 = (\beta_2 - \gamma_2 \beta_1)$  is estimated from (8), where the conditional expectations are estimated by kernel regressions with standard normal kernel and bandwidth equal to 0.2 times the standard deviation of  $x'_i \widehat{\gamma}$  (in the sample where  $y_i$  is observed). We choose  $c_1 = -\infty$ ,  $c_2$  to  $c_9$  are the deciles of  $\{y_i - x'_{i2} \widehat{\alpha}_2 - (x'_i \widehat{\gamma}) b_1\}$  in the sample where  $y_i$  is observed and  $c_{10} = \infty$ . The sets  $A_k$  corresponds to the intervals between quintiles of  $x'_i \widehat{\gamma}$ . The figure also displays the maximum likelihood estimator and Heckman's two-step estimator for  $\beta_1$  along with their 95% confidence intervals.

Figure 5: Obj. Function of the Coefficient on 3rd Generation Mexican-American



Similarly, we define a sample analog of the solutions to the population inequalities

<sup>11</sup>Alternatively, one could use a semiparametric estimator such as Han (1987)'s maximum rank correlation estimator in the first step.



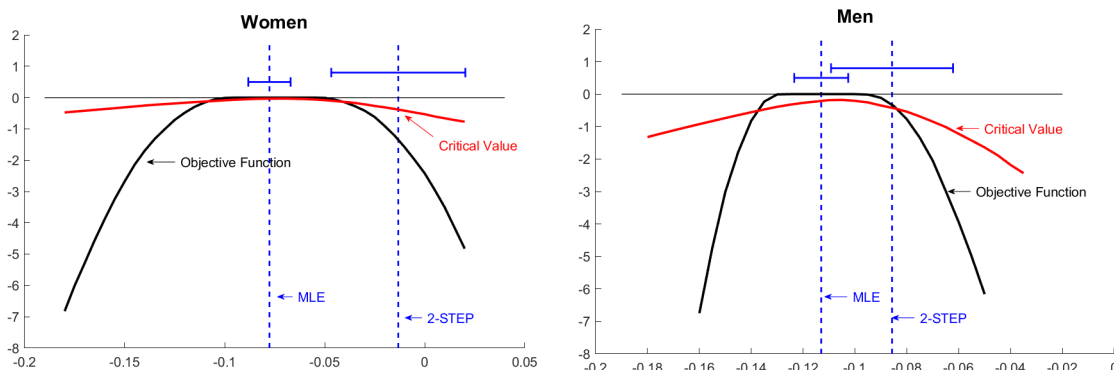
in (17) by the solutions to  $Q_{2n}(b_1) = 0$ , where

$$Q_{2n}(b_1) = - \sum_{\ell} \min \left\{ \widehat{E} \left[ (1 \{c_{\ell} < y_j - x'_{j2}\alpha_2 + (x_{i1} - x_{j1} + (x_{i2} - x_{j2})' \gamma_2) b_1 \leq c_{\ell+1}, d_j = 1\} \right. \right. \\ \left. \left. - 1 \{c_{\ell} < y_i - x'_{i2}\alpha_2 \leq c_{\ell+1}, d_i = 1\} \right) 1 \{x_{i1} + x'_{i2}\gamma_2 < x_{j1} + x'_{j2}\gamma_2\} \right], 0 \right\}^2 = 0$$

where  $c_1 = -\infty$ ,  $c_2$  to  $c_9$  are the deciles of  $\{y_i - x'_{i2}\widehat{\alpha}_2 - (x'_{i2}\widehat{\gamma}) b_1\}$  and  $c_{10} = \infty$ . Figure 6 displays the objective function and the 5%-critical value function. Because of the double-sum involved in calculating the sample analog of  $\widehat{E}$  in this case, we use a sub-sample size of 5,000.

The figure also displays the maximum likelihood estimator and Heckman's two-step estimator for  $\beta_1$  along with their 95% confidence intervals.

Figure 6: Obj. Function of the Coefficient on 3rd Generation Mexican-American



The implications of Figures 5 and 6 are similar. As expected, the implied confidence intervals for the identified sets are larger than the confidence intervals based on the parametric estimators. Our set estimate contains the maximum likelihood estimate for both samples. For men, it is also close to the two-step estimate. For women, the two-step estimate is, however, quite different from our estimated set as well as from the maximum likelihood estimate. This casts doubt on the validity of the normality assumption for women. On the other hand, the moment inequalities implied by the independence assumption (equation (15)) are not rejected by the data in either sample.

## 4.4 Lee-Style Bounds Using Powell-style Regressions

In this subsection, we discuss a different idea for constructing bounds for  $\beta_1$  in (5). This idea applies Lee (2009)'s trimming logic to the pairwise difference estimator in Powell (1987).

Recall that  $\alpha_2 = (\beta_2 - \gamma_2\beta_1)$  and  $\gamma_2$  are both identified, and that  $\gamma_1$  is normalized to 1. We therefore write  $\beta = (0, \alpha_2')' + \gamma\beta_1$ . Consider two observations  $i$  and  $j$  such that  $y_i$  and  $y_j$  are observed and  $x_{i1} > x_{j1}$ . Recall that  $z_i = x_i'\gamma$  and suppose that  $z_i > z_j$ . Similar to equation (9), we have

$$y_\ell - x'_{2\ell}\alpha_2 = (x_{1\ell} + x'_{2\ell}\gamma_2)\beta_1 + \varepsilon_\ell$$

for  $\ell = i, j$ . This implies

$$(y_i - x'_{i2}\alpha_2) - (y_j - x'_{j2}\alpha_2) = ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2)\beta_1 + (\varepsilon_i - \varepsilon_j). \quad (18)$$

Now suppose that we regress  $(y_i - x'_{i2}\alpha_2) - (y_j - x'_{j2}\alpha_2)$  on  $(x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2$ . This would lead to a biased estimator of  $\beta_1$  because  $(\varepsilon_i - \varepsilon_j)$  will not have mean 0. The problem is that the selection has trimmed  $\varepsilon_i$  and  $\varepsilon_j$  in different ways. The distribution of  $\varepsilon_i$  before selection is composed of a fraction  $(1 - p(z_i))$  for which  $y_i$  is not observed, a fraction  $p(z_i) - p(z_j)$  for which  $y_i$  is observed, but would not have been for observation  $j$ , and a fraction  $p(z_j)$  for which  $y_i$  is observed, and also would have been for observation  $j$ . This means that in the sample for which  $y_i$  is observed, the distribution of  $\varepsilon_i$  is composed of a fraction  $\frac{p(z_i) - p(z_j)}{p(z_i)}$  that has not been selected away, but would have been for observation  $j$ . So we need to trim a fraction,  $\frac{p(z_i) - p(z_j)}{p(z_i)}$ , of the distribution of  $\varepsilon_i$  in order for  $\varepsilon_i$  and  $\varepsilon_j$  to have the same distribution after selection and trimming. This is exactly Lee (2009)'s argument in the case where there is a single explanatory variable, which is binary.

Applying this insight to construct bounds on the parameter of interest in our regression context,  $\beta_1$ , we need to know exactly which observations to trim in order to get upper and lower bounds for  $\beta_1$ . Since the OLS estimator of  $\beta_1$  in (18) has the

form

$$\left( \sum_{i,j} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2)^2 \right)^{-1} \sum_{i,j} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2) ((y_i - x'_{i2} \alpha_2) - (y_j - x'_{j2} \alpha_2)),$$

it is clear that to get the upper bound for  $\beta_1$ , we need to trim the lower tail of  $(y_i - x'_{i2} \alpha_2)$  and the upper tail  $(y_j - x'_{j2} \alpha_2)$  when  $(x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2$  is positive, and we need to trim the upper tail of  $(y_i - x'_{i2} \alpha_2)$  and the lower tail  $(y_j - x'_{j2} \alpha_2)$  when  $(x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2$  is negative.

This leads to the estimator of the upper bound for  $\beta_1$

$$\widehat{\beta}_1^U = \left\{ \sum_{i,j} 1 \{nottrimmed_{ij}\} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2)^2 \right\}^{-1} \left\{ \sum_{i,j} 1 \{nottrimmed_{ij}\} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2) ((y_i - x'_{i2} \alpha_2) - (y_j - x'_{j2} \alpha_2)) \right\} \quad (19)$$

where  $1 \{nottrimmed_{ij}\}$  is defined in Appendix 3. The corresponding estimator of the lower upper bound for  $\beta_1$

$$\widehat{\beta}_1^L = \left\{ \sum_{i,j} 1 \{nottrimmed_{ij}\} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2)^2 \right\}^{-1} \left\{ \sum_{i,j} 1 \{nottrimmed_{ij}\} ((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2}) \gamma_2) ((y_i - x'_{i2} \alpha_2) - (y_j - x'_{j2} \alpha_2)) \right\} \quad (20)$$

where  $1 \{nottrimmed_{ij}\}$  is also defined in Appendix 3.

Of course, in practice,  $\gamma$  (and hence  $z_i$ ),  $F$  (defined in Appendix 3) and  $p$  are unknown and have to be estimated. The asymptotic distribution of the upper and lower bound for  $\beta_1$  can then found using Chen, Linton, and Van Keilegom (2003).

One could also estimate (non-sharp) bounds for  $\beta_1$  by using median regression rather than OLS after applying the same trimming as above. To see why this is the case, consider median regression of a variable  $y_i$  on  $x_i$  (without a constant). In the

population, this estimator would solve  $E[x_i \text{sign}(y_i - x_i b)] = 0$  for  $b$ . When  $x_i$  is positive, trimming the low values of  $y_i$  would lead to a larger value of the solution for  $b$ , and when  $x_i$  is negative, trimming the high values of  $y_i$  would lead to a larger value of the solution for  $b$ . This leads to exactly the same trimming as for the OLS estimator.

**Example 4** (*Example 3 continued*) *The bounds obtained by applying OLS above are*<sup>12</sup> (0.29, 1.05).

## 4.5 Empirical Illustration Part 3

Using the approach in Section 4.4, we estimate the lower and upper bounds on the coefficient on being third-generation Mexican-American for men to be  $-0.1546$  and  $-0.0568$ , respectively. For the women, the corresponding bounds are<sup>13</sup>  $-0.1342$  and  $-0.0144$ . These intervals contain the bounds estimated in Section 4.3. We leave it for future research to investigate whether the bounds here can be interpreted as bounds under milder assumptions than those made in Section 4.2.

## 4.6 Discrete Regressors

Identification of  $\gamma$  (up to scale) rules out data generating processes where  $x_i$  is a set of dummy variables, but  $\gamma$  is still set-identified in that case. Let  $\Gamma$  be the identified

---

<sup>12</sup>We calculate these bounds using 400,000 simulation draws. We use the true  $\gamma$  and  $p(z)$  (see Appendix 3). The parameter  $a_2$  is estimated as in Section 4.3 using a normal kernel and bandwidth equal to 0.02 times the standard deviation of  $x'_i \hat{\gamma}$  conditional on selection. The cumulative distribution function,  $F$ , is estimated using nonparametric regression with a normal kernel and the same bandwidth.

<sup>13</sup>To calculate these bounds, we first estimate  $\gamma$  and  $\alpha_2$  as in Section 4.3. The function  $p(z)$  is then estimated by nonparametric kernel regression of  $d_i$  on  $x'_i \hat{\gamma}$  using a normal kernel and bandwidth equal to 0.2 times the standard deviation of  $x'_i \hat{\gamma}$ . The cumulative distribution function  $F$  in Appendix 3 is also calculated by nonparametric kernel regression now with bandwidth equal to 0.2 times the standard deviation of  $x'_i \hat{\gamma}$  conditional on selection. We use a larger bandwidth than in Example 3 because  $\gamma$  is estimated and the sample is smaller.

set for  $\gamma$ . The identified region for  $\beta_1$  from Section 4.1 is then

$$\begin{aligned} \mathbf{B} &= \bigcup_{g \in \Gamma} \{b_1 \in \mathbb{R} : P(c_1 < y_i - x_{i2}\alpha_2 - (x_{i1} + x_{i2}g_2) b_1 \leq c_2, d_i = 1 \mid (x_{i1} + x_{i2}g_2) = \xi_1) \\ &\leq P(c_1 < y_i - x_{i2}\alpha_2 - (x_{i1} + x_{i2}g_2) b_1 \leq c_2, d_i = 1 \mid (x_{i1} + x_{i2}g_2) = \xi_2) \\ &\quad \text{for all values of } c_1 < c_2 \text{ and } \xi_1 < \xi_2\}. \end{aligned}$$

We next consider the regression style approach in Section 4.4 in the case where the vector  $x_i$  can take a finite number of different values. Specifically, assume that  $x_i$  can take the values  $(w_1, \dots, w_m)$  with probabilities  $(\pi_1, \dots, \pi_m)$  and sample frequencies  $(n_1, \dots, n_m)$ . For each pair,  $(w_k, w_\ell)$  with selection probabilities  $p(w_k) \geq p(w_\ell)$ , we need to trim a fraction,  $\frac{p(w_k) - p(w_\ell)}{p(w_k)}$ , of the observations corresponding to  $x_i = w_k$ . If we knew which observations to trim, we could express a consistent estimator of  $\beta$  as

$$\begin{aligned} &\left[ \sum_{k,l} n_k n_\ell \mathbf{1}\{p(w_k) \geq p(w_\ell)\} \frac{p(w_\ell)}{p(w_k)} (w_k - w_\ell) (w_k - w_\ell)' \right]^{-1} \\ &\sum_{i,j} \mathbf{1}\{p(x_i) \geq p(x_j)\} \mathbf{1}\{not\_trim_{ij}\} (x_i - x_j) (y_i - y_j) \\ &= A \sum_{i,j} \mathbf{1}\{p(x_i) \geq p(x_j)\} \mathbf{1}\{not\_trim_{ij}\} (x_i - x_j) (y_i - y_j) \\ &= \sum_{i,j} \mathbf{1}\{p(x_i) \geq p(x_j)\} \mathbf{1}\{not\_trim_{ij}\} A (x_i - x_j) (y_i - y_j) \\ &= \sum_{k,\ell} \mathbf{1}\{p(w_k) \geq p(w_\ell)\} A (w_k - w_\ell) \sum_{x_i=w_k, x_j=w_\ell} \mathbf{1}\{not\_trim_{ij}\} (y_i - y_j). \end{aligned}$$

Now suppose that we are interested in the first element of  $\beta$ . To get an upper bound, we would trim the lower tail of  $y_i$  if the first element of  $A(w_k - w_\ell)$  is positive, and the upper tail of  $y_i$  if the first element of  $A(w_k - w_\ell)$  is negative.

## 5 Concluding Remarks

This paper has studied identification in a classical semiparametric sample selection model in which both the selection mechanism and outcome of interest depend linearly on the same explanatory variables, and the errors are independent of the explanatory variables. This model is not semiparametrically point-identified, but the sharp identified set is one-dimensional. Toy-examples as well as an empirical application suggest that the identified set can be quite small in practice. In this respect, the practical take-away of this paper is similar to papers in different areas of economics which have demonstrated that the identified regions of non-identified parameters can be small enough to be useful in empirical applications. The papers by Haile and Tamer (2003), Honoré and Lleras-Muney (2006), and Blundell, Gosling, Ichimura, and Meghir (2007) are early examples of this.

The numerical calculations presented in this paper illustrate that the bounds obtained under the semiparametric model considered here are much tighter than those obtained in Lee (2009)'s nonparametric setting. We leave it for future research to investigate intermediate assumptions that are weaker than those imposed here, but strong enough to generate identified sets that are small enough to be empirically informative.

## References

- ANDREWS, D. W. K., AND M. M. A. SCHAFGANS (1998): "Semiparametric Estimation of the Intercept of a Sample Selection Model," *The Review of Economic Studies*, 65(3), 497–517.
- BARROW, L., AND C. E. ROUSE (2017): "Financial Incentives and Educational Investment: The Impact of Performance-Based Scholarships on Student Time Use," *Education Finance and Policy*.
- BLUNDELL, R., A. GOSLING, H. ICHIMURA, AND C. MEGHIR (2007): "Changes in

- the Distribution of Male and Female Wages Accounting for Employment Composition Using Bounds,” *Econometrica*, 75(2), 323–363.
- CANAY, I. A., AND A. M. SHAIKH (2017): “Practical and Theoretical Advances in Inference for Partially Identified Models,” in *Advances in Economics and Econometrics: Eleventh World Congress*, ed. by B. Honoré, A. Pakes, M. Piazzesi, and L. Samuelson, vol. 2, pp. 271–306. Cambridge University Press.
- CHAMBERLAIN, G. (1986): “Asymptotic Efficiency in Semi-Parametric Models with Censoring,” *Journal of Econometrics*, 32(2), 189–218.
- CHEN, X., O. LINTON, AND I. VAN KEILEGOM (2003): “Estimation of Semiparametric Models when the Criterion Function Is Not Smooth,” *Econometrica*, 71(5), 1591–1608.
- DAS, M., W. K. NEWEY, AND F. VELLA (2003): “Nonparametric Estimation of Sample Selection Models,” *The Review of Economic Studies*, 70(1), 33–58.
- ESCANCIANO, J. C., D. JACHO-CHVEZ, AND A. LEWBEL (2016): “Identification and estimation of semiparametric two-step models,” *Quantitative Economics*, 7(2), 561–589.
- HAILE, P., AND E. TAMER (2003): “Inference with an Incomplete Model of English Auctions,” *Journal of Political Economy*, 111(1), 1–51.
- HAN, A. (1987): “Nonparametric Analysis of a Generalized Regression Model,” *Journal of Econometrics*, 35, 303–316.
- HECKMAN, J. J. (1976): “The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables and a Simple Estimator for Such Models,” *Annals of Economic and Social Measurement*, 5(4), 475–92.
- (1979): “Sample Selection Bias as a Specification Error,” *Econometrica*, 47(1), 153–61.

- (1990): “Varieties of Selection Bias,” *The American Economic Review*, 80(2), 313–318.
- HONORÉ, B. E., AND A. LLERAS-MUNEY (2006): “Bounds in Competing Risks Models and the War on Cancer,” *Econometrica*, 74(6), 1675–1698.
- KITAGAWA, T. (2015): “A Test for Instrument Validity,” *Econometrica*, 83(5), 2043–2063.
- KRUEGER, A. B., AND D. M. WHITMORE (2001): “The Effect of Attending a Small Class in the Early Grades on College-test Taking and Middle School Test Results: Evidence from Project Star,” *The Economic Journal*, 111(468), 1–28.
- LEE, D. S. (2009): “Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects,” *The Review of Economic Studies*, 76(3), 1071–1102.
- MANSKI, C. F. (1989): “Anatomy of the Selection Problem,” *The Journal of Human Resources*, 24(3), 343–360.
- MORA, R. (2008): “A nonparametric decomposition of the Mexican American average wage gap,” *Journal of Applied Econometrics*, 23(4), 463–485.
- MUTHÉN, B. (1990): “Moments of the censored and truncated bivariate normal distribution,” *British Journal of Mathematical and Statistical Psychology*, 43(1), 131–143.
- POWELL, J. L. (1987): “Semiparametric Estimation of Bivariate Latent Models,” Working Paper no. 8704, Social Systems Research Institute, University of Wisconsin–Madison.
- (1994): “Estimation of Semiparametric Models,” in *Handbook of Econometrics*, ed. by R. F. Engle, and D. L. McFadden, no. 4 in Handbooks in Economics, pp. 2443–2521. Elsevier, North-Holland, Amsterdam, London and New York.
- ROBINSON, P. M. (1988): “Root-N-Consistent Semiparametric Regression,” *Econometrica*, 56(4), 931–954.



# Appendix 1: Proofs of Propositions

**Proof of Proposition 1.** This proposition is a special case of Proposition 2<sup>14</sup>. Here we provide a more readable proof that explicitly uses properties of the normal distribution. Recall that if

$$\begin{pmatrix} \nu \\ y \end{pmatrix} \Big| x \sim N \left( \begin{pmatrix} \mu \\ x'\beta \end{pmatrix}, \begin{pmatrix} 1 & \rho\sigma \\ \rho\sigma & \sigma^2 \end{pmatrix} \right),$$

then

$$\nu | y, x \sim N \left( \mu + \frac{\rho}{\sigma} (y - x'\beta), 1 - \rho^2 \right).$$

Hence

$$\begin{aligned} f(y | \nu > -x) &= \frac{f(y) P(\nu > -x | y)}{P(\nu > -x)} \\ &= \frac{1}{\sigma} \varphi \left( \frac{y - x'\beta}{\sigma} \right) \Phi \left( \frac{x + \mu + \frac{\rho}{\sigma} (y - x'\beta)}{\sqrt{1 - \rho^2}} \right) \Big/ \Phi(x + \mu) \end{aligned}$$

and therefore

$$f_y(\cdot | x) = \frac{1}{\sigma} \varphi \left( \frac{y - x'\beta}{\sigma} \right) \Phi \left( \frac{x + \mu + \frac{\rho}{\sigma} (y - x'\beta)}{\sqrt{1 - \rho^2}} \right).$$

Now consider at  $b$  in the identified region, **B**. For that  $b$ , the inequality  $f_y(c | x = 0) \leq f_y(c + b | x = 1)$  holds for all values of  $c$ . This can be written as

$$\frac{f_y(c | x = 0)}{f_y(c + b | x = 1)} \leq 1.$$

---

<sup>14</sup>Note to referees: We are happy to take this proof out if you think it is unnecessary given the proof of Proposition 2. We have kept it because the concreteness of the calculation helped us understand the results.

Under normality, the inequality becomes

$$\begin{aligned} \frac{f_y(c|x=0)}{f_y(c+b|x=1)} &= \frac{\frac{1}{\sigma}\varphi\left(\frac{c}{\sigma}\right)\Phi\left(\frac{\mu+\frac{\rho}{\sigma}c}{\sqrt{1-\rho^2}}\right)}{\frac{1}{\sigma}\varphi\left(\frac{c+b-\beta}{\sigma}\right)\Phi\left(\frac{1+\mu+\frac{\rho}{\sigma}(c+b-\beta)}{\sqrt{1-\rho^2}}\right)} \\ &= \exp\left((b-\beta)c + (b-\beta)^2/2\sigma^2\right) \frac{\Phi\left(\frac{\mu+\frac{\rho}{\sigma}c}{\sqrt{1-\rho^2}}\right)}{\Phi\left(\frac{1+\mu+\frac{\rho}{\sigma}(c+b-\beta)}{\sqrt{1-\rho^2}}\right)} \leq 1. \end{aligned}$$

Now assume that  $\rho > 0$  and consider the limit as  $c \rightarrow \infty$ . If  $b > \beta$ , the first term in the product increases to  $\infty$ , while the second term converges to 1. This contradicts the inequality, and we conclude that  $b \leq \beta$ . Hence  $\beta$  is the upper endpoint of  $\mathbf{B}$ .

When  $\rho < 0$ , we consider the limit as  $c \rightarrow -\infty$  and conclude that  $b \geq \beta$ . Hence  $\beta$  is the lower endpoint of  $\mathbf{B}$ .

Finally, when  $\rho = 0$ , the inequality becomes

$$(b-\beta)c \leq -\log\left(\frac{\Phi(\mu)}{\Phi(1+\mu)}\right) - (b-\beta)^2/2$$

for all values of  $c$ . This can only be true if  $b = \beta$ , and  $\beta$  is point-identified.

This completes the proof. ■

## Proof of Proposition 2.

Assumptions:

1. The distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  if  $c_1 > c_2$ .

2. The density of  $\varepsilon$  has sufficiently thin tails that for  $a > 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow \infty$  and for  $a < 0$ ,  $f(c)/f(c+a) \rightarrow \infty$  as  $c \rightarrow -\infty$ .

Recall that

$$\begin{aligned} f_y(c|x) &= f_{y^*}(c)P(\nu > -x|y^* = c) \\ &= f_\varepsilon(c-x\beta)P(\nu > -x|\varepsilon = c-x\beta). \end{aligned}$$

Now consider at  $b$  in the identified region,  $\mathbf{B}$ . For that  $b$ , the inequality  $f_y(c|x=0) \leq f_y(c+b|x=1)$  holds for all values of  $c$ . In other words

$$\frac{f_y(c|x=0)}{f_y(c+b|x=1)} = \frac{f_\varepsilon(c) P(\nu > 0|\varepsilon = c)}{f_\varepsilon(c+b-\beta) P(\nu > -1|\varepsilon = c+b-\beta)} \leq 1.$$

Suppose that  $b > \beta$ . Then

$$\frac{P(\nu > 0|\varepsilon = c)}{P(\nu > -1|\varepsilon = c+b-\beta)} > \frac{P(\nu > 0|\varepsilon = c)}{1}.$$

The right hand side is increasing in  $c$  by the stochastic dominance assumption. Hence it is bounded from below by some positive constant,  $k$ . Therefore

$$\frac{f_y(c|x=0)}{f_y(c+b|x=1)} > k \frac{f_\varepsilon(c)}{f_\varepsilon(c+b-\beta)},$$

where the ratio on the right hand side increases to  $\infty$  as  $c$  goes to  $\infty$ . This contradicts the inequality, and we conclude that no  $b$  in  $\mathbf{B}$  can be greater than the true  $\beta$ . Hence  $\beta$  is the upper endpoint of  $\mathbf{B}$ .

Now consider the case where the distribution of  $\nu$  given  $\varepsilon = c_1$  stochastically dominates the distribution of  $\nu$  given  $\varepsilon = c_2$  if  $c_1 < c_2$ . Suppose that  $b < \beta$ . Then again

$$\frac{P(\nu > 0|\varepsilon = c)}{P(\nu > -1|\varepsilon = c+b-\beta)} > P(\nu > 0|\varepsilon = c) > P(\nu > 0|\varepsilon = 0)$$

for all  $c < 0$ . Therefore

$$\frac{f_y(c|x=0)}{f_y(c+b|x=1)} > k \frac{f_\varepsilon(c)}{f_\varepsilon(c+b-\beta)}$$

for  $c < 0$ . Taking the limit as  $c \rightarrow -\infty$  brings the right hand side above 1, and we conclude that a  $b$  for which  $b < \beta$  cannot belong to the set  $\mathbf{B}$ . Hence  $\beta$  is the upper endpoint of  $\mathbf{B}$ .

Finally, when  $\varepsilon$  and  $\nu$  are independent, the inequality defining  $\mathbf{B}$  is

$$\frac{f_\varepsilon(c) P(\nu > 0 | \varepsilon = c)}{f_\varepsilon(c + b - \beta) P(\nu > -1 | \varepsilon = c + b - \beta)} = \frac{f_\varepsilon(c) P(\nu > 0)}{f_\varepsilon(c + b - \beta) P(\nu > -1)} \leq 1.$$

Taking the limit as  $c \rightarrow -\infty$  generates a contradiction when  $b < \beta$  and taking the limit as  $c \rightarrow \infty$  generates a contradiction when  $b > \beta$ . Therefore  $\mathbf{B} = \{\beta\}$ .

This completes the proof. ■

**Proof of Proposition 3.** First consider case 1 (with positive selection). The identified set can be written as

$$\begin{aligned} \mathbf{B} &= \{b \in \mathbb{R} : f_{y|x}(c + xb | \xi_1) \leq f_{y|x}(c + xb | \xi_2) \\ &\quad \text{for all values of } c \text{ and } \xi_1 < \xi_2 \text{ in the support of } x\} \\ &= \bigcap_{\xi_1 < \xi_2} \{b \in \mathbb{R} : f_{y|x}(c + \xi_1 b | \xi_1) \leq f_{y|x}(c + \xi_2 b | \xi_2) \text{ for all values of } c\} \\ &= \bigcap_{\xi_1 < \xi_2} \{b \in \mathbb{R} : f_{y-\xi_1 b|x}(c | \xi_1) \leq f_{y-\xi_1 b|x}(c + (\xi_2 - \xi_1)b | \xi_2) \text{ for all values of } c\}. \end{aligned}$$

Now consider one of the sets on the right hand side above. By Proposition 2, the upper limit of  $(\xi_2 - \xi_1)b$  will be  $(\xi_2 - \xi_1)\beta$ . This implies that the upper limit of all the sets in the intersection above is the true  $\beta$ . Hence, the upper limit on the intersection is  $\beta$ .

The proofs of cases 2 and 3 are similar. ■

## Appendix 2: Data Details

This analysis utilizes the Merged Outgoing Rotation Groups (MORG) files of the Current Population Survey (CPS), which were prepared by the National Bureau of Economic Research (NBER). Following Mora, we restrict our sample to non-Hispanic whites and Mexican-Americans between the ages of 25 and 62 (inclusive) who live in Arizona, California, New Mexico, or Texas. We further limit our analysis to those who have at least one parent born in the United States (i.e., third-generation Americans).

We also drop the top 1.67% of earners in each year's income distribution from our analysis, and we multiply top-coded earnings by 1.33. Finally, in our wage samples, we exclude self-employed workers, as well as individuals who report that they are working but do not report either hours worked or earnings.

The variables are:

- Log hourly wage: Calculated by taking the natural log of an individual's weekly earnings divided by his usual hours works, adjusted for inflation.
- Veteran status: Indicator variable that equals one if an individual ever reported serving in the U.S. military, and zero otherwise.
- Married: Indicator variable that equals one if an individual reports that she or he is either (1) a married civilian with spouse present, (2) a married Armed Forces member with spouse present, or (3) married with spouse absent or separated, and zero otherwise.
- Experience: For individuals who have completed at least seventh grade, their labor market experience is defined as their age (in single years) minus their education-years minus 6. Individuals whose educational attainment is less than seventh grade are assigned an experience level equal to their age minus thirteen.
- Education-years: Following Mora, we assign education-years based on the level of education attainment reported in the data as follows:
  - Less than 1st grade = 0 years of education
  - 1st – 4th grade = 2.5 years of education
  - 5th or 6th grade = 5.5 years of education
  - 7th or 8th grade = 7.5 years of education
  - 9th = 9 years of education
  - 10th = 10 years of education
  - 11th = 11 years of education

- 12th grade (no diploma) = 12 years of education
- High school graduate, diploma, or GED = 12 years of education
- Some college but no degree = 13.5 years of education
- Associate degree – occupational/vocational = 14 years of education
- Associate’s degree – academic program = 14 years of education
- Bachelor’s degree (i.e. BA, AB, BS) = 16.5 years of education
- Master’s degree (i.e. MA, MS, MEng, MSW, MBA) = 18 years of education
- Professional school degree (i.e. MD, DDS, DVM, LLB, JD) = 18 years of education
- Doctorate degree (i.e. PhD, EdD) = 20 years of education

## Appendix 3: Expressions for Trimming Indicators in Section 4.4

Define  $F(\eta; z_i)$  to be the cdf of  $y_i - x'_{i2}\alpha_2$  conditional on selection and on  $z_i (= x_{i1} + x'_{i2}\gamma_2)$ . The expression for  $1\{\text{nottrimmed}_{ij}\}$  in (19) is then

$$\begin{aligned}
1\{\text{nottrimmed}_{ij}\} &= 1\{((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2) > 0\} 1\{z_i > z_j\} \\
&\quad \cdot 1\left\{F(y_i - x'_{i2}\alpha_2; z_i, x_{i1} + x'_{i2}\gamma_2) > \frac{p(z_i) - p(z_j)}{p(z_i)}\right\} \\
&+ 1\{((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2) < 0\} 1\{z_i > z_j\} \\
&\quad \cdot 1\left\{F(y_i - x'_{i2}\alpha_2; z_i, x_{i1} + x'_{i2}\gamma_2) < 1 - \frac{p(z_i) - p(z_j)}{p(z_i)}\right\} \\
&+ 1\{((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2) > 0\} 1\{z_i < z_j\} \\
&\quad \cdot 1\left\{F(y_j - x'_{j2}\alpha_2; z_j, x_{j1} + x'_{j2}\gamma_2) < 1 - \frac{p(z_j) - p(z_i)}{p(z_j)}\right\} \\
&+ 1\{((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2) < 0\} 1\{z_i < z_j\} \\
&\quad \cdot 1\left\{F(y_j - x'_{j2}\alpha_2; z_j, x_{j1} + x'_{j2}\gamma_2) > \frac{p(z_j) - p(z_i)}{p(z_j)}\right\}
\end{aligned}$$

and the expression for  $1\{\text{nottrimmed}_{ij}\}$  in (20) is

$$\begin{aligned}
1\{\text{nottrimmed}_{ij}\} &= 1\left\{\left((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2\right) > 0\right\} 1\{z_i > z_j\} \\
&\quad \cdot 1\left\{F(y_i - x'_{i2}\alpha_2; z_i, x_{i1} + x'_{i2}\gamma_2) < 1 - \frac{p(z_i) - p(z_j)}{p(z_i)}\right\} \\
&+ 1\left\{\left((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2\right) < 0\right\} 1\{z_i > z_j\} \\
&\quad \cdot 1\left\{F(y_i - x'_{i2}\alpha_2; z_i, x_{i1} + x'_{i2}\gamma_2) > \frac{p(z_i) - p(z_j)}{p(z_i)}\right\} \\
&+ 1\left\{\left((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2\right) > 0\right\} 1\{z_i < z_j\} \\
&\quad \cdot 1\left\{F(y_j - x'_{j2}\alpha_2; z_j, x_{j1} + x'_{j2}\gamma_2) > \frac{p(z_j) - p(z_i)}{p(z_j)}\right\} \\
&+ 1\left\{\left((x_{i1} - x_{j1}) + (x'_{i2} - x'_{j2})\gamma_2\right) < 0\right\} 1\{z_i < z_j\} \\
&\quad \cdot 1\left\{F(y_j - x'_{j2}\alpha_2; z_j, x_{j1} + x'_{j2}\gamma_2) < 1 - \frac{p(z_j) - p(z_i)}{p(z_j)}\right\}.
\end{aligned}$$

Table 1: Summary Statistics

	Mexican Women		White Women		Mexican Men		White Men	
	mean	sd	mean	sd	mean	sd	mean	sd
California	0.38	0.49	0.46	0.50	0.37	0.48	0.47	0.50
Arizona	0.08	0.26	0.11	0.31	0.08	0.28	0.11	0.31
Texas	0.47	0.50	0.34	0.47	0.46	0.50	0.34	0.47
Real Wage	2.16	0.57	2.42	0.63	2.34	0.58	2.61	0.61
Working	0.64	0.48	0.61	0.49	0.71	0.45	0.67	0.47
Age	40.66	10.76	44.34	10.81	40.56	10.74	43.99	10.92
Experience	21.63	11.19	23.88	11.15	21.65	10.97	23.63	11.08
Less than HS	0.16	0.37	0.04	0.20	0.16	0.37	0.05	0.21
Some College	0.33	0.47	0.34	0.48	0.31	0.46	0.33	0.47
College	0.12	0.32	0.27	0.44	0.11	0.31	0.26	0.44
Advanced Degree	0.05	0.21	0.12	0.33	0.04	0.18	0.12	0.32
Married	0.53	0.50	0.62	0.49	0.55	0.50	0.61	0.49
Veteran	0.01	0.10	0.02	0.13	0.12	0.32	0.16	0.37
No. Observations	26,698		103,209		21,402		97,016	



Table 2: Parametric Sample Selection Model

	Women		Men	
	MLE	2-Step	MLE	2-Step
Mexican-American	-0.078 (0.005)	-0.013 (0.017)	-0.113 (0.005)	-0.084 (0.012)
Age	0.113 (0.007)	0.213 (0.026)	0.079 (0.006)	0.112 (0.014)
Age-squared	-0.047 (0.005)	-0.118 (0.018)	-0.047 (0.004)	-0.072 (0.010)
Experience	-0.070 (0.006)	-0.127 (0.016)	-0.025 (0.005)	-0.045 (0.009)
Experience-squared	0.006 (0.004)	0.036 (0.009)	-0.014 (0.004)	-0.007 (0.005)
Less than HS	-0.177 (0.015)	-0.372 (0.050)	-0.170 (0.012)	-0.222 (0.023)
Some College	0.033 (0.011)	0.017 (0.014)	0.051 (0.009)	0.043 (0.011)
College	0.155 (0.025)	0.084 (0.036)	0.235 (0.023)	0.205 (0.026)
Advanced Degree	0.199 (0.034)	0.113 (0.047)	0.257 (0.031)	0.194 (0.041)
Veteran	0.030 (0.016)	0.037 (0.020)	-0.001 (0.006)	0.015 (0.008)
Married	0.033 (0.005)	-0.079 (0.028)	0.136 (0.005)	0.185 (0.019)
California	0.204 (0.007)	0.178 (0.011)	0.151 (0.007)	0.140 (0.009)
Arizona	0.098 (0.009)	0.103 (0.012)	0.042 (0.009)	0.052 (0.010)
Texas	0.031 (0.008)	0.064 (0.013)	0.015 (0.008)	0.045 (0.014)
Year Dummies	yes	yes	yes	yes
No. Observations	127, 738	127, 738	118, 250	118, 250

Standard errors in parentheses.