

# Simultaneous Sample Selection Models\*

William Matcham<sup>†</sup>

September 24, 2022

## Abstract

I extend sample selection models by allowing the outcome to affect selection directly. I microfound the model then provide identification and estimation results for semiparametric and parametric models. The simultaneity between the outcome and selection generates additional endogeneity, and, unlike traditional sample selection models, my identification results require an excluded regressor in the *outcome* equation. Simulations confirm the finite sample performance of the new estimator and show sizeable differences in parameters compared to models that do not account for the direct effect of the outcome on the selection decision. I finish with an application relating to the examination process for patents and patent's potential quality. I show that traditional sample selection methods understate the positive effect of the inventing firm's size on patent quality.

**Keywords:** Endogenous sample selection, Heckman Selection, Structural Parameters, Patent Prosecution Process

**JEL Classification:** C34, C35, J31, O34

---

\*I especially thank Mark Schankerman for many discussions relating to this paper. I also thank Thomas Glinnan, Tim Cejka, Roberto Maura-Rivero, Kamila Nowakowicz, Taisuke Otsu, and Marcia Schafgans for comments and suggestions.

<sup>†</sup>Department of Economics, London School of Economics. Email: w.o.matcham@lse.ac.uk

# 1 Introduction

Economic datasets are rarely a random sample from the population. Resultantly, a rich literature in economics and other social sciences explicitly models the sample selection process. The typical approach projects the outcome variable and sample selection utility onto a set of exogenous regressors and, in doing so, does not estimate a direct effect of the outcome on the selection decision. In this paper, I address this gap in the literature, analyzing sample selection models that include the dependent variable in the selection equation. My central and novel contribution is to provide sufficient conditions for identification of the complete set of model parameters alongside consistent and asymptotically normal estimation methods. Through an empirical example in the context of the patent application process I illustrate the empirical relevance and importance of my contribution.

As the natural starting point, I consider a model where the two structural equations are linear in parameters and the outcome variable is continuous. The model contains the traditional sample selection issue: the researcher only observes the outcome for an endogenously selected subsample. The novelty is that the value of the outcome itself *directly* affects the selection decision. These two features combined create a simultaneity issue, rendering traditional identification results invalid. I provide sufficient conditions for identification in two cases. The first is the semiparametric case, where the joint distribution of the errors is unspecified. The second is the parametric case, where the errors are joint Gaussian with unknown parameters. Both identification results require an excluded regressor in the outcome equation, a condition stronger than what is typical in sample selection models.

I follow identification by discussing estimation. I discuss possible estimation methods in the semiparametric case and detail more formally how to estimate all parameters in the parametric case via likelihood methods. The likelihood is similar to the one from traditional sample selection estimation, except the contribution to the likelihood by those not selected consists entirely of reduced form parameters. Additionally, in the parametric case I show how to estimate the outcome equation parameters using a two-step regression method, analogous to Heckman (1976). I follow with simulations of the parametric model, which (1) confirm the finite-sample performance of the estimator and (2) illustrate substantial differences between the new estimator and existing methods when a direct effect exists.

I finish with an application relating to the patent system. In that context, researchers study

patent quality (scope), *conditional* on obtaining a patent (Lanjouw and Schankerman, 2004; Kuhn and Thompson, 2019; Feng and Jaravel, 2020). But inventors only enjoy patent protection if the expected utility from a given level of scope exceeds the utility from alternative intellectual property protection such as trade secrecy. When I allow for a direct effect of potential patent quality on the selection decision, I find a more substantial positive impact of the inventing firm’s size on patent quality. Existing methods, which correct only for *unobservables* driving both patent quality and selection, find a modest difference in patent quality between small and large firms because they miss the fact that, all else equal, small firms must have especially high-quality patents to be selected. My application emphasizes the importance of including a direct effect of the outcome on the selection equation where appropriate.

## 2 Related Literature

This paper relates to an influential literature on sample selection models, which started by analyzing women’s labour force participation (Heckman, 1974; Gronau, 1974; Lewis, 1974).<sup>1</sup> Seminal work by Heckman (Heckman, 1976, 1979, 1990) formalized the idea in these empirical examples, motivating several extensions as surveyed in Vella (1998). My paper generalizes traditional selection models in a way made precise in Section 4.1.

In my setup, since the outcome directly affects the selection decision, and observing the outcome depends on the selection decision, the model relates to simultaneous equation models. Amemiya (1974) considers simultaneous equation models with truncated dependent variables rather than endogenously selected ones. Amemiya (1984) provides a general survey of Tobit models. The closest to my setup is Amemiya’s “Type 2”.<sup>2</sup> In another related paper, Heckman (1978) presents a broad range of simultaneous equation models, including one with endogenous dummy variables in a simultaneous equation system. My setup resembles this paper’s, except in Heckman (1978) the outcome is always observed.

Recent work on endogenous sample selection models splits into two broad categories. One part of the literature relaxes the linearity of the outcome and selection equations and the

---

<sup>1</sup>For a detailed textbook discussion of self-selectivity, see Maddala (1986) and Gourieroux (2000), or shorter summaries in Amemiya (1984, 1985).

<sup>2</sup>This is a model given by  $y_{1i}^* = x'_{1i}\beta_1 + u_{1i}$  and  $y_{2i}^* = x'_{2i}\beta_2 + u_{2i}$  with  $y_{2i} = y_{2i}^*y_{1i}$  and  $y_{1i} = 1(y_{1i}^* > 0)$ .

parametric distributional assumptions made on the error terms, studying semiparametric and nonparametric sample selection models (Ahn and Powell, 1993; Andrews and Schafgans, 1998; Das, Newey, and Vella, 2003). I focus on linear equations and (semi)parametric distributional assumptions for the error term, but nonlinear and nonparametric extensions are possible.

The second, and even more recent, strand of the literature on sample selection models focuses on robustness. For example, Bastos, Barreto-Souza, and Genton (2022) and Carlson (2022) account for heteroskedasticity, Marchenko and Genton (2012) addresses heavy tails through the use of student-t errors, Ogundimu and Hutton (2016) studies skewness of outcomes, and Bastos and Barreto-Souza (2021) uses the bivariate Birnbaum-Sanders distribution to address nonnegativity of outcome variables. Roelsgaard and Taylor (2022) proposes a semiparametric machine learning estimator for sample selection models. All of these potential issues or novelties persist in my setup and future work can check if the proposed solutions carry through.

There are numerous *applications* of traditional sample selection models.<sup>3</sup> In ongoing work, Matcham and Schankerman (2022) adds an equation for post-patenting outcomes to the two-equation model of patent scope and selection I estimate in my application.

## 3 Motivating Examples

Below, I describe four economic examples that motivate my model.

### 3.1 Wages and labor Force Participation

A common motivation for the sample selection model is that wages are only observed for those who participate in the labor force, and factors that determine labor force participation correlate with factors that determine wages. In particular, the original structural

---

<sup>3</sup>Among others between 1970 and 1985, see: Heckman (1974, 1979), Nelson (1977), Cogan (1981), Hanoch (2014, 1976) Gordon and Blinder (1980), Griliches, Hall, and Hausman (1978), Kenny, Lee, Maddala, and Trost (1979), Willis and Rosen (1979), Lee (1978), Abowd and Farber (1982), Katz (1977), Nakosteen and Zimmer (1980), Poirier and Ruud (1981), Weisbrod (1980), Roberts, Maddala, and Enholm (1978), Trost (1977), Lee (1978), Rosen (1979), and King (1980).

model motivating the sample selection literature *does* include the outcome directly in the selection equation. The three equations are

$$W_M = X'\beta_X + X'_M\beta_M + \varepsilon_M \quad (1)$$

$$W_R = X'\alpha_X + X'_R\alpha_R + \varepsilon_R \quad (2)$$

$$D = 1 [W_M > W_R], \quad (3)$$

where  $W_M$  denotes market wage,  $W_R$  the reservation wage,  $D$  is an indicator of labor force participation, and  $(\varepsilon_M, \varepsilon_R)$  are jointly normal. In particular, researchers only observe market wages when they exceed individuals' rarely-observed reservation wage. Substituting  $W_R$  from equation (2) into (3) yields

$$D = 1 \left[ W_M\kappa_M + X'\kappa_X + X'_R\kappa_R - \varepsilon_R > 0 \right], \quad (4)$$

with  $(\kappa_M, \kappa_X, \kappa_R) = (1, -\alpha_X, -\alpha_R)$ . This is close to the econometric model I study in this paper. Further, substituting  $W_M$  from (1) into (4) gives

$$W_M = X'\beta_X + X'_M\beta_M + \varepsilon_M$$

$$D = 1 [Z'\gamma + \varepsilon_D > 0],$$

where  $Z = (X, X'_M, X'_R)'$ ,  $\gamma = (\beta'_X - \alpha'_X, \beta'_M, -\alpha'_R)'$ , and  $\varepsilon_D = \varepsilon_M - \varepsilon_R$ . This is the typical endogenous sample selection model, with  $X_R$  and all excluded regressors  $X_M$  added to the selection equation. The index  $Z'\gamma + \varepsilon_D$  is a reduced form representation of the difference between market wage  $W_M$  and reservation wage  $W_R$ . On common regressors, the parameters of the selection equation are the difference of the structural parameters  $\beta_X$  and  $\alpha_X$ . Further, because  $\varepsilon_D = \varepsilon_M - \varepsilon_R$ , the estimated covariance is not  $\text{cov}(\varepsilon_M, \varepsilon_R)$  but instead  $\text{cov}(\varepsilon_M, \varepsilon_M - \varepsilon_R)$ . For the case of jointly normal errors, my model will estimate all parameters in the model given by (1) - (3) along with  $\text{cov}(\varepsilon_M, \varepsilon_R)$ .

### 3.2 Entry Games

Consider an incumbent firm deciding whether to block ( $D = 0$ ) or allow ( $D = 1$ ) a potential firm's entry. If the new firm enters, they choose the optimal production quantity  $Y_N$ , and the incumbent chooses the quantity  $Y_I$  that maximizes their duopoly profit  $\pi_d(Y_I, Y_N)$ . If the potential entrant gets blocked, the incumbent obtains monopoly profit  $\pi_m$ , which includes a blocking cost. The incumbent blocks if  $\pi_m > \pi_d(Y_I, Y_N)$  so

$$D = 1 [\pi_m - \pi_d(Y_I, Y_N) > 0].$$

Hence the decision to block depends on  $Y_N$ . A researcher either analyzing multiple markets of this nature, or analyzing the same market over many periods will only observe the entrant's output  $Y_N$  if the entrant is not blocked, and  $Y_N$  directly affects the incumbent's decision to block. Therefore, researchers can estimate the structural parameters of the incumbent's blocking choice and the potential entrant's output using the econometric model I analyze in this paper.

### 3.3 Patent Scope

To apply for a patent in the United States, inventors must submit an application to the United States Patent Trademark Office. First, the Office assigns an examiner to the application. Then, the examiner debates with the applicant (and their attorney) on the allowable level of patent protection (scope) to avoid infringing on existing patented inventions.<sup>4</sup> An often-missed detail about the patent application process is that the examiner cannot outright end the application process. Instead, the process can only end when the applicant meets the examiner's demands or the applicant abandons. Relating the process to sample selection models, we only observe patent scope ( $Y$ ) when the applicant does not abandon ( $D$ ), and the applicant's decision to abandon is a function of the allowed patent quality, should it realize.

### 3.4 Survey Data Response

In settings where researchers collect sensitive survey data, the outcome variable may directly affect individuals' willingness and ability to be surveyed. Suppose a hospital is analyzing alcohol consumption in the community. The value of alcohol consumption directly affects willingness to participate, yet the hospital can only study the alcohol consumption of those who are willing and able to participate. As another example, research analyzing the correlates of domestic abuse know that domestic abuse itself has an effect on the willingness to report exposure.

---

<sup>4</sup>See [Graham, Marco, and Miller \(2018\)](#) for more detail about the application process.

## 4 Economic and Econometric Models

An individual decides whether to obtain an outcome  $Y$ , with their choice recorded in the dummy variable  $D$ . The potential amount of the outcome is  $Y = g_Y(Z_Y, \varepsilon_Y)$ , where  $Z_Y$  are exogenous features of the individual that affect the outcome quantity and  $\varepsilon_Y$  is a random shock affecting the outcome level.

The difference between the benefits and costs of realizing the outcome is  $f(Y, Z_D, \varepsilon_B)$ , where  $Z_D$  are exogenous features of the individual affecting benefits and costs, and  $\varepsilon_B$  is a random shock. The outside option delivers utility  $\varepsilon_D$ .<sup>5</sup> When the individual makes their decision,  $\varepsilon_D$  and  $\varepsilon_Y$  have realized but, since the individual is yet to know the full benefits,  $\varepsilon_B$  has not.<sup>6</sup> The individual chooses  $D = 1$  if

$$g_D(Y, Z_D, \varepsilon_D) := \mathbb{E}_{\varepsilon_B | (Y, \varepsilon_D, Z_D, Z_Y)} [f(Y, Z_D, \varepsilon_B) - \varepsilon_D] > 0$$

The full model is thus

$$\begin{aligned} Y &= g_Y(X_Y, X, \varepsilon_Y) \\ D &= 1 \left[ g_D(Y, X_D, X, \varepsilon_D) > 0 \right], \end{aligned}$$

where I split  $Z_Y$  and  $Z_D$  into the common regressors  $X$  and excluded regressors  $X_Y$  and  $X_D$ .<sup>7</sup> The value of  $Y$  is only observed if  $D = 1$ , though I assume  $X_Y$  is always observed. This is the general model of the paper. Taking  $g_Y$  and  $g_D$  as linear (as I will in all that follows) the model is

$$Y = \gamma'_Y X_Y + \gamma'_X X \varepsilon_Y = \gamma' w_Y + \varepsilon_Y \quad (5)$$

$$D = 1 \left[ \beta'_D X_D + \beta'_X X + \beta_Y Y + \varepsilon_D > 0 \right] = 1 \left[ \beta'_{-Y} w_D + \beta_Y Y + \varepsilon_D > 0 \right], \quad (6)$$

where  $\gamma = (\gamma'_Y, \gamma'_X)'$ ,  $w_Y = (X'_Y, X)'$ ,  $\beta_{-Y} = (\beta'_D, \beta'_X)'$ , and  $w_D = (X'_D, X)'$ .

---

<sup>5</sup>The outside option could be a function of exogenous variables as well as a random shock, but nothing is gained from this addition.

<sup>6</sup>For example, in the patent context, the inventor may know their level of allowable scope, but will not know the market returns that their invention will bring. Further, they will not know if their invention will become obsolete, so they will not know how long they will renew the patent rights for.

<sup>7</sup>I could combine exogenous regressors  $X$  and  $X_Y$  (respectively  $X$  and  $X_D$ ) into a vector, say  $Z_Y$  (respectively  $Z_D$ ), but for exposition, the identification result and proof that follows splits common and separate regressors.

## 4.1 Relation to Standard Sample Selection Model

The remainder of the paper focuses on identification, estimation, and application of the model as given in (5) - (6). Before that it is befitting to compare this model to traditional endogenous sample selection models. To see this, substituting (5) into (6) yields

$$\begin{aligned} Y &= \gamma'_Y X_Y + \gamma'_X X + \varepsilon_Y \\ D &= 1 \left[ \kappa'_D X_D + \kappa'_X X + \kappa'_Y X_Y + \tilde{\varepsilon} > 0 \right], \end{aligned} \quad (7)$$

where

$$\kappa_D = \beta_D, \quad \kappa_Y = \beta_Y \gamma_Y, \quad \kappa_X = \beta_X + \beta_Y \gamma_X, \quad \tilde{\varepsilon} = \beta_Y \varepsilon_Y + \varepsilon_D. \quad (8)$$

This is a standard sample selection model, where the  $\gamma$  parameters are structural. But the parameters of the selection equation and the correlation between errors are reduced form are not structural. The traditional sample selection model will not estimate the  $\beta$  parameters. When  $\beta_Y = 0$ , the model collapses to the standard sample selection model, but the case of  $\beta_Y \neq 0$  requires additional analysis. Further, if the researcher fails to include  $X_Y$  in the selection equation, traditional sample selection estimation methods will fail to estimate any of the structural parameters, even in the outcome equation.

## 5 Identification

Now I turn to identification. Throughout,  $Y$  is continuous – I discuss identification with a discrete outcome briefly in Appendix B. I study identification of equations (5) and (6), where  $Y$  is observed if  $D = 1$ . The distribution of the errors  $(\varepsilon_Y, \varepsilon_D)$  is either unspecified (semiparametric) or conditionally jointly normal distributed (parametric). I start with the parametric case, since much of the intuition in this specific, more straightforward, case carries over to the semiparametric case.

### 5.1 Parametric Case

In this case, the errors  $(\varepsilon_Y, \varepsilon_D)$  are conditionally jointly normal distributed. Letting  $Z = (X'_D, X'_Y, X)'$ , the Gaussian case requires three standard normalizations:  $\text{Var}(\varepsilon_D|Z) = 1$  and  $\mathbb{E}(\varepsilon_Y|Z) = \mathbb{E}(\varepsilon_D|Z) = 0$ . Hence,

$$\begin{pmatrix} \varepsilon_Y \\ \varepsilon_D \end{pmatrix} \Big| Z \sim \mathcal{N}_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \sigma_Y^2 & \\ \rho\sigma_Y & 1 \end{pmatrix} \right). \quad (9)$$



To understand the identification result, it helps to analyze the reduced form. Rearranging (7) for  $\tilde{\varepsilon}$  and standardizing yields

$$D = 1 \left[ \underbrace{\frac{\tilde{\varepsilon}}{\sqrt{v}}}_{\mathcal{N}(0,1)} > -\pi'Z \right] = 1 \left[ \frac{\beta_Y \varepsilon_Y + \varepsilon_D}{\sqrt{v}} > \tilde{Z} \right], \quad (10)$$

where  $v = \text{Var}(\beta_Y \varepsilon_Y + \varepsilon_D) = \beta_Y^2 \sigma_Y^2 + 1 + 2\beta_Y \sigma_Y \rho$ ,  $\pi = (\pi'_D, \pi'_Y, \pi'_X)'$ ,  $\tilde{Z} = -\pi'Z$ , and

$$\pi_D = \frac{\beta_D}{\sqrt{v}}, \quad \pi_Y = \frac{\beta_Y \gamma_Y}{\sqrt{v}}, \quad \pi_X = \frac{\beta_X + \beta_Y \gamma_X}{\sqrt{v}}. \quad (11)$$

In what follows, I refer to  $\pi$  as the *standardized reduced form* parameters and  $Z$  as the *exogenous* variables.

My parametric identification result rests on two substantive conditions. First, I rule out perfect multicollinearity in exogenous regressors. Second, I require at least one regressor  $X_{Y_j}$  appearing in the equation for  $Y$  but not in the equation for  $D$ . This is different to the traditional sample selection model, which does not strictly require any excluded regressors. However, similar to traditional sample selection estimation methods, excluded regressors in the selection equation are also desirable, else identification only comes from the nonlinearity in the inverse Mills Ratio. Before the statement of the theorem, I state the requirement of an excluded regressor in the outcome equation as a formal requirement.

**Assumption 1 (Excluded Outcome Regressor)** *Either  $\beta_Y = 0$ , or there exists  $j$  such that  $\gamma_{Y_j} \neq 0$  and  $X_{Y_j} \neq 0$ .*

The identification result for the parametric case follows. I prove the result and all other formal results in Appendix A.

**Theorem 1** *Consider the model defined by (5), (6), and (9).*

*Part 1 ( $\gamma$ ): Assume that:*

1. *(Moment Existence):  $Y$  is absolutely continuous with positive and finite variance.*
2. *(No Perfect Multicollinearity): There are  $K_1$  points denoted  $\{Z_{(k)}\}_{k=1}^{K_1}$  such that the*

matrix

$$\mathcal{X}_1 = \begin{pmatrix} Z'_{(1)} \\ Z'_{(2)} \\ \vdots \\ Z'_{(K_1)} \end{pmatrix}$$

has full rank, where  $Z_{(k)} = (X'_{Dk}, X'_{Yk}, X'_k)' \in \mathbb{R}^{K_1}$ . Further, the matrix

$$\mathbb{E} \left[ w_Y - \mathbb{E}(w_Y | \tilde{\lambda}(Z)) \right] \left[ w_Y - \mathbb{E}(w_Y | \tilde{\lambda}(Z)) \right]' \quad (12)$$

is positive definite, where  $\tilde{\lambda}(Z) = \frac{\phi(Z'\pi)}{\Phi(Z'\pi)}$  and  $\pi$  are the standardized reduced form parameters as defined in (11).

Then  $\gamma$  is identified.

*Part 2 ( $F_\varepsilon$  and  $\beta$ ):* In addition, suppose that Assumption 1 holds. Then  $\beta$  and the parameters of the joint distribution of errors, denoted  $F_\varepsilon$ , (i.e.  $\sigma_Y$  and  $\rho$ ) are identified.

The proof of the parametric case proceeds in four main steps. First, I identify the standardized reduced form parameters using standard parametric binary choice identification. With the standardized reduced form parameters, in the second step I identify the parameters of the outcome equation. The idea is that conditioning on the selected sample adds the inverse Mills ratio into the outcome equation, after which linear regression identification results apply. Third, I use the exclusion restriction to identify the variance of  $\varepsilon_Y$ ,  $\beta_Y$ , and  $\rho$ . Finally, with everything else identified, the remaining parameters of the selection equation,  $\beta_{-Y}$ , follow trivially.

## 5.2 Semiparametric Case

Now I make no distributional assumption on the errors and derive alternative conditions for identification of the model parameters. The goal is to identify  $\beta$  and  $\gamma$ , and any possible components of the joint distribution of the errors.

Before the formal presentation of the theorem for identification of  $\gamma$  and  $\beta$ , I introduce some notation. Let  $\mathcal{Z}$  denote the support of  $Z$  and  $X_{D1}$  denote the first element of  $X_D$  (which is also the first element of  $Z$ ). Let  $Z_{-1}$  denote the subvector of  $Z$  including all

components of  $Z$  except  $X_{D1}$ . Let

$$S = \left\{ Z \in \mathcal{Z} \mid \mathbb{P}(D = 0|Z) \in (0, 1) \right\}.$$

Finally, let  $S_{-1}$  denote the projection of  $S$  on  $Z_{-1}$ .

**Theorem 2** *Consider the model defined by (5) and (6).*

*Part 1 ( $\gamma$ ): Assume that:*

1. *(Moment Existence):  $Y$  is absolutely continuous with finite first moment.*
2. *(Independence):  $(\varepsilon_Y, \varepsilon_D)$  is independent of  $Z$ . The support of  $(\varepsilon_Y, \varepsilon_D)$  is a convex set in  $\mathbb{R}^2$  with a non-empty interior.*
3. *(Sufficient Continuous Variation in Excluded Regressor):  $S$  is not contained in any proper linear subspace of  $\mathbb{R}^{\dim(Z)}$  and  $\mathbb{P}(S) > 0$ . There is a covariate in  $X_D$  (say,  $X_{D1}$  without loss of generality) such that for  $Z_{-1} \in S_{-1}$ , the support of  $X_{D1}|Z_{-1}$  intersected with  $S$  contains a non-empty interval  $(\underline{X}_{D1}, \bar{X}_{D1})$  with  $\underline{X}_{D1} < \bar{X}_{D1}$ .*
4. *(Normalization)  $\beta_{D1} = 1$*
5. *(No Perfect Multicollinearity) For reduced form parameters  $\kappa = (\kappa'_D, \kappa'_Y, \kappa'_X)'$  and  $\check{Z} = -\kappa'Z$ ,*

$$\mathbb{E} [w_Y - \mathbb{E}(w_Y|\check{Z})] [w_Y - \mathbb{E}(w_Y|\check{Z})]'$$
(13)

*is positive definite*

*Then  $\gamma$  is identified.*

*Part 2 ( $\beta$ ): In addition, suppose that Assumption 1 holds. Then  $\beta$  is identified.*

My formal identification result includes similar assumptions to the parametric case, such as the requirement of an excluded regressor in the outcome equation. However, I also require independence of errors and regressors, a normalization restriction on one parameter, and sufficient continuous variation in a regressor. Normalizing one parameter is standard in the semiparametric identification literature. Sufficient continuous variation in a regressor is also standard in the discrete-response literature and is implied by more stringent conditions on covariates (Manski, 1985, 1988; Horowitz, 2010; Komarova and Matcham, 2022). Intuitively, my condition guarantees some minimum desirable variation in the covariates.

Further, I assume from the outset that there is a also excluded regressor in the *selection* equation. It is even more pertinent in the semiparametric case that, even though identification may be possible without an excluded regressor in the selection equation, it would only be due to functional form rather than anything more substantive.

The proof of the semiparametric case proceeds similarly to the parametric equivalent. First, I identify the reduced form parameters of the selection equation using binary choice identification conditions from the semiparametric literature of discrete choice (Manski, 1985, 1988; Horowitz, 2010; Komarova and Matcham, 2022). In the second step, I work with the outcome equation. In this case, conditioning on the selected sample adds an *unknown function* to the selection equation, thereby mimicking the partially linear model (Robinson, 1988; Hardle, Liang, and Gao, 2002). Identification follows straightforwardly from the standard linear regression identification after “partialling out” the equivalent of the inverse Mills ratio. The “no perfect multicollinearity” assumption ensures the identification of the residual-based regression. Having identified the reduced form parameters and the outcome equation parameters, the excluded regressor in the outcome equation immediately aids in identifying  $\beta_Y$  and therefore the remaining  $\beta_{-Y}$  parameters.

Identification of the joint distribution of the errors is more difficult. I cannot identify the full distribution of the error terms because  $Y$  is unobserved when  $D = 0$ . It is immediate (and hence not proved) that with all other parameters identified, I can identify the marginal distribution of the reduced form error  $\varepsilon_D + \beta_Y \varepsilon_Y$  along with the distribution of  $\varepsilon_Y$  conditional on  $D = 1$ .

## 6 Estimation

Like identification, I split my work on estimation into the semiparametric and parametric case. I provide formal results in the parametric case, and discuss approaches to estimation in the semiparametric case without formal results on the asymptotic distribution of the estimator. Throughout this section, I assume access to a random sample  $\mathcal{W} = \{Y_i D_i, D_i, X_{Y_i}, X_{D_i}, X_i\}_{i=1}^n$  where observing  $Y_i D_i$  is equivalent to observing  $Y_i$  only when  $D_i = 1$ .<sup>8</sup>

---

<sup>8</sup>It is typical in sample selection models for the assigned value of  $Y$  when  $D = 0$  to be irrelevant. That is the case here, though in a number of settings such as the patent scope example above,  $Y$  will take the

## 6.1 Parametric Case

In this subsection, I present a likelihood-based method to estimate the parameters of the econometric model as defined by (5) - (9), and I describe a regression-based approach to estimate the parameters of the outcome equation.

### 6.1.1 Likelihood Methods

Letting  $\theta = (\gamma', \beta', \sigma_Y, \rho)'$ , the conditional likelihood function for  $n$  independent observations is

$$\mathcal{L}(\theta; \mathcal{W}) = \prod_{i=1}^n \wp_{0i}^{1-D_i} \wp_{1i}^{D_i}, \quad (14)$$

with the two  $\wp$  terms defined below. The term  $\wp_{0i}$  is the contribution to the likelihood when  $D_i = 0$  and  $\wp_{1i}$  represents the contribution when  $D_i = 1$ . Given the assumptions on  $(\varepsilon_{Yi}, \varepsilon_{Di})$  and letting  $Z_i = (X'_{Yi}, X'_{Di}, X'_i)'$ , from (10)

$$\begin{aligned} \wp_{0i} &= \mathbb{P} \left( \frac{\beta_Y \varepsilon_{Yi} + \varepsilon_{Di}}{\sqrt{v}} \leq -\pi' Z_i | Z_i \right) \\ &= \Phi(-\pi' Z_i) \end{aligned} \quad (15)$$

where the equality in (15) follows from

$$\beta_Y \varepsilon_{Yi} + \varepsilon_{Di} \sim \mathcal{N}_1 \left( 0, \underbrace{\beta_Y^2 \sigma_Y^2 + 1 + 2\rho\beta_Y\sigma_Y}_v \right).$$

Equation (15) reveals that the contribution to the likelihood from individuals with  $D_i = 0$  is a function of the standardized reduced form parameters  $\pi$ , so offers no assistance in separating structural parameters from reduced form. Next, as derived in section A.3,

$$\wp_{1i} = \frac{1}{\sigma_Y} \phi \left( \frac{Y_i - w'_{Yi} \gamma}{\sigma_Y} \right) \Phi \left( \frac{\beta'_{-Y} w_{Di} + \beta_Y Y_i + \frac{\rho}{\sigma_Y} (Y_i - w'_{Yi} \gamma)}{\sqrt{1 - \rho^2}} \right).$$

The log-likelihood therefore has a closed form and yields, under standard MLE regularity conditions, consistent and asymptotically normal estimates of  $\theta$ :

$$\sqrt{n} (\hat{\theta} - \theta) \xrightarrow{D} \mathcal{N}(0, V^{-1}),$$

with  $V = -\mathbb{E} \left( \frac{\partial^2 \log(\mathcal{L})}{\partial \theta \partial \theta'} \right)$ .

---

value 0 when  $D = 0$ .

### 6.1.2 Regression Methods

Let  $\lambda_i = \frac{\phi(\tilde{Z}_i)}{\Phi(-\tilde{Z}_i)}$ , where  $\tilde{Z}_i = -\pi'Z_i$ . Following from the identification proof,

$$\mathbb{E}(Y_i|D_i = 1, Z_i) = \gamma'w_{Y_i} + \tau\lambda_i.$$

Resultantly, the outcome equation for the subset with  $D_i = 1$  is

$$Y_i = \gamma'w_{Y_i} + \tau\lambda_i + V_i, \quad (16)$$

where  $\mathbb{E}(V_i|D_i = 1, Z_i) = 0$ . Hence, the following two-step procedure estimates  $\gamma$  and  $\tau$ :

1. Run a probit regression of  $D$  on  $Z$ . Let these reduced form estimates of  $\pi$  be  $\hat{\pi}$ .
2. Run a regression using data for  $D = 1$  of  $Y$  on  $X_Y$ ,  $X$ , and  $\hat{\lambda}(Z) = \frac{\phi(\hat{\pi}'Z)}{\Phi(\hat{\pi}'Z)} = \frac{\phi(\hat{Z})}{\Phi(-\hat{Z})}$ .

In Appendix A, I provide further details on this procedure. In particular, under general conditions, discussed extensively in [Amemiya \(1973\)](#) and [Jennrich \(1969\)](#),

$$\sqrt{n_1} \begin{pmatrix} \hat{\gamma} - \gamma \\ \hat{\tau} - \tau \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, B\psi B'),$$

where asymptotics are  $n, n_1 \rightarrow \infty$  with  $\frac{n_1}{n} \rightarrow k \in (0, 1)$ . Here  $n_1 = \sum_i D_i$  is the number of observations with  $D_i = 1$  (the selected subsample),

$$B = \text{plim}_{n_1 \rightarrow \infty} \left( \frac{W'_{Y+} W_{Y+}}{n_1} \right)^{-1},$$

where  $W_{Y+}$  is the matrix stacking  $X_Y$ ,  $X$ , and  $\hat{\lambda}$ , and

$$\psi = \text{plim}_{\substack{n \rightarrow \infty \\ n_1 \rightarrow \infty}} \left[ \frac{\sigma_Y^2}{n_1} \begin{pmatrix} \sum_i w_{Y_i} w'_{Y_i} \eta_i & \sum_i w_{Y_i} \lambda_i \eta_i \\ \sum_i \lambda_i w'_{Y_i} \eta_i & \sum_i \lambda_i^2 \eta_i \end{pmatrix} + \tau^2 \left( \frac{n_1}{n} \right) \begin{pmatrix} \frac{1}{n_1^2} \sum_i \sum_j w_{Y_i} w'_{Y_j} \vartheta_{ij} & \frac{1}{n_1^2} \sum_i \sum_j w_{Y_i} \Upsilon_{ij} \\ \frac{1}{n_1^2} \sum_i \sum_j w'_{Y_i} \Upsilon_{ij} & \frac{1}{n_1^2} \sum_i \sum_j \Omega_{ij} \end{pmatrix} \right],$$

where

$$\begin{aligned} \eta_i &= 1 + \frac{\tau^2(\tilde{Z}_i \lambda_i - \lambda_i^2)}{\sigma_Y^2}, \\ \vartheta_{ij} &= \partial \lambda_i \partial \lambda_j Z'_i V_{1st} Z_j, \\ \Upsilon_{ij} &= \lambda_i \vartheta_{ij}, \\ \Omega_{ij} &= \lambda_i \lambda_j \vartheta_{ij}, \end{aligned}$$

$V_{1st}$  is the asymptotic variance covariance matrix of  $\hat{\pi}$  and  $\partial\lambda_i = \partial\lambda_i/\partial\tilde{Z}$  is the derivative of  $\lambda_i$  with respect to  $\tilde{Z}_i$ . Conventional regression standard errors are invalid and must be adjusted to match  $B\psi B'$  by replacing  $\sigma_Y$ ,  $\lambda_i$ , and  $\tau$  with appropriate estimators.

## 6.2 Semiparametric Case

In what follows I discuss one approach to estimate the model parameters when there are no assumptions on the distribution of the error terms. The estimation method is a semiparametric version of the Heckman (1976) (and the regression approach in section 6.1.2). It applies semiparametric binary choice estimation to the selection equation, then uses partially linear semiparametric estimation to estimate the parameters of the outcome equation.

More precisely, the first step uses semiparametric binary choice methods, such as Klein and Spady (1993), to estimate  $\kappa$  from equation

$$\mathbb{P}(D_i = 0|Z_i) = F_{\varepsilon_D + \beta_Y \varepsilon_Y}(-\kappa'Z_i).$$

Denote the estimates by  $\check{\kappa}$ . The second step applies any of the estimation methods for a partially linear model (Robinson, 1988; Hardle, Liang, and Gao, 2002) to estimate  $\gamma$  from the partially linear model

$$Y_i = \gamma'w_{Y_i} + g(\check{Z}_i) + U_i$$

where  $\check{Z}_i = -\check{\kappa}'Z_i$ . The asymptotic distribution of the first-step estimator follows from Klein and Spady (1993). Deriving the asymptotic distribution of the second-step estimator requires adapting the asymptotic distribution results as in Robinson (1988) (or other partially linear estimation procedures) for the estimation error from  $\check{\kappa}$ .

## 7 Simulation

Before an empirical application, I present Monte Carlo simulations. The simulations showcase the finite sample performance of the estimator compared to existing methods and confirm the intuitions I gave for the identification result. For three designs, I run 500 simulations of (5) - (9), each with a sample size of 2000. Throughout  $\bar{D}$  represents the sample average of  $D_i$  and therefore the proportion of observations that are selected.

## Simulation 1: No Direct Effect

In all three simulations, variables are continuous and scalar-valued:  $X$  contains a constant and a regressor uniformly distributed on  $[-5, 5]$ ,  $X_D$  is t-distributed with 10 degrees of freedom, and  $X_Y$  follows a logistic distribution with location and scale equal to two and one, respectively. I set  $\beta_{-Y} = (\beta_D, \beta_0, \beta_X) = (-8, 1, 3)'$ ,  $\gamma = (\gamma_Y, \gamma_0, \gamma_X) = (-2, 3, -1)$ ,  $\rho = 0.5$ , and  $\sigma = 2$ . To confirm that I nest traditional sample selection models, in Simulation 1 I set  $\beta_Y = 0$  to rule out a direct effect of  $Y$  on the probability of treatment. Table 1 Panel A presents the simulation results and the selection probability. Row (1) provides means and standard deviation of the ML estimator as described in Section 6. Row (2) represents traditional sample selection methods, using MLE on a model imposing  $\beta_Y = 0$ .<sup>9</sup> Both methods estimate the parameters precisely, with small bias. Existing methods are marginally more precise: they benefit from not estimating  $\beta_Y$ , instead forcing it to equal 0.

## Simulation 2: Direct Effect

The design of Simulation 2 mirrors Simulation 1, except  $\beta_Y$  is equal to two. Table 1 Panel B reports the results. Since the design meets the identification conditions, the new method estimates the parameters with minimal bias and standard deviation. Existing sample selection methods fail to estimate  $\beta$ . Appendix C Table 3 presents the results of a simulation similar to Simulation 2, except that  $X_Y$  is discrete. Since the identification result does not require continuous regressors, it is unsurprising that the new method still estimates parameters well, with slightly more noise owing to the reduced variation in the regressor.

## Simulation 3: No Excluded Regressor

Finally, Simulation 3 mirrors Simulation 2 except with  $\gamma_Y = 0$ , so there is no excluded regressor. The results in Table 1 Panel C confirm the intuition of the identification result: both methods cannot estimate  $\beta_{-Y}$ ,  $\beta_Y$ , or  $\rho$ . In particular, traditional sample selection methods estimate  $\gamma$  parameters correctly because  $\gamma_Y = 0$  implies that the selection equation does not contain the omitted variable  $X_Y$ .

---

<sup>9</sup>Results from using the two-step Heckman estimator are almost identical to MLE results throughout.



TABLE 1: Simulation Results

---

*Panel A: No Direct Effect:  $\bar{D} = 0.53$*

---

<b>Parameter:</b>	$\sigma$	$\rho$	$\gamma_0$	$\gamma_X$	$\gamma_Y$	$\beta_0$	$\beta_X$	$\beta_D$	$\beta_Y$
<b>True Value:</b>	2	0.5	3	-1	-2	1	3	-8	0
<b>(1) New:</b>	2.06 (0.06)	0.49 (0.12)	3.00 (0.11)	-1.00 (0.03)	-2.00 (0.04)	1.04 (0.15)	3.08 (0.30)	-8.22 (0.78)	0.00 (0.03)
<b>(2) Existing:</b>	2.00 (0.04)	0.50 (0.12)	3.00 (0.11)	-1.00 (0.03)	-2.00 (0.04)	1.03 (0.13)	3.07 (0.28)	-8.20 (0.75)	–

---

*Panel B: Direct Effect:  $\bar{D} = 0.47$*

---

<b>Parameter:</b>	$\sigma$	$\rho$	$\gamma_0$	$\gamma_X$	$\gamma_Y$	$\beta_0$	$\beta_X$	$\beta_D$	$\beta_Y$
<b>True Value:</b>	2	0.5	3	-1	-2	1	3	-8	2
<b>(1) New:</b>	2.07 (0.08)	0.49 (0.13)	3.00 (0.09)	-1.00 (0.02)	-2.00 (0.04)	1.04 (0.26)	3.10 (0.51)	-8.27 (1.31)	2.07 (0.35)
<b>(2) Existing:</b>	2.14 (0.07)	0.80 (0.04)	2.15 (0.15)	-1.00 (0.02)	-1.57 (0.05)	-0.12 (0.03)	0.12 (0.01)	-0.96 (0.04)	–

---

*Panel C: No Excluded Regressor:  $\bar{D} = 0.76$*

---

<b>Parameter:</b>	$\sigma$	$\rho$	$\gamma_0$	$\gamma_X$	$\gamma_Y$	$\beta_0$	$\beta_X$	$\beta_D$	$\beta_Y$
<b>True Value:</b>	2	0.5	3	-1	0	1	3	-8	2
<b>(1) New:</b>	2.27 (0.39)	0.84 (0.18)	2.99 (0.23)	-1.01 (0.25)	-0.02 (0.24)	1.54 (0.64)	0.59 (0.88)	-2.62 (1.64)	0.31 (0.74)
<b>(2) Existing:</b>	2.00 (0.04)	0.98 (0.01)	3.00 (0.06)	-1.00 (0.02)	0.00 (0.02)	1.53 (0.05)	0.22 (0.01)	-1.75 (0.06)	–

---

Notes: Table 1 reports the sample mean and sample standard deviations (in parentheses) of the estimates of the model parameters, over 500 repeated samples. The “New” row provide estimates from using the newly proposed model. The “Existing” row estimates a traditional sample selection model which forces  $\beta_Y = 0$ .

## 8 Application

Finally, I apply the methods described in this paper to estimate a model of patent quality (scope) and selection into patenting. Section 3.3 provides background on the patent application process, and explains why it fits into my model. I access data on approximately 300,000 patent applications filed to the United States Patent and Trademark Office between 2012-2014. For each application, I observe whether the applicant abandoned ( $D = 0$ ) or obtained a patent ( $D = 1$ ). As a measure of scope ( $Y$ ), I use the number of independent claims on the granted patent.<sup>10</sup> Scope is not observed for patents that are not granted. I also observe a dummy for whether the firm applying has fewer than 500 employees ( $X$ ), which is clearly determined before the realization of  $D$ . Firm size affects scope through many channels, for example, larger firms may have better quality inventions. Size also affects the decision to abandon because smaller firms may not be able to pay the legal fees associated with multiple rounds of negotiation with a patent examiner.

For excluded regressors I create two leniency measures for examiners assigned to applications. Leniency instruments are a common identification strategy in applied economic research.<sup>11</sup> In the case of patents, there is support in the literature that examiners are as good as randomly assigned to applications, that is, that applications of a certain type or quality are not assigned to specific examiners.<sup>12</sup> The first of the two leniency measures ( $X_D$ ) is the average grant rate of the examiner assigned to the application. This is a common instrument used to explain the decision to patent (see e.g. [Sampat and Williams \(2019\)](#), [Farre-Mensa, Hegde, and Ljungqvist \(2020\)](#), [Gaulé \(2018\)](#), amongst others). The second measure of leniency is average number of independent claims on applications granted by the examiner assigned to the application. The “exclusion restriction” is that the leniency measure only affects the decision to patent through its effect on scope.

Table 2 provides the estimation results. A couple of findings warrant attention. First, the value of  $\beta_Y$  is statistically significant and, in being one fifth of the constant in the selection

---

<sup>10</sup>This variable is one of the four indicators used in the index of patent quality in [Lanjouw and Schankerman \(2004\)](#).

<sup>11</sup>[Kling \(2006\)](#) is one of the original applications; see [Frandsen, Lefgren, and Leslie \(2019\)](#) for a non-exhaustive list of applications.

<sup>12</sup>I exclude examiners who have conducted fewer than 50 examinations, though results are robust to other threshold choices.

equation, implies that there is a direct effect of scope on the decision to abandon or obtain a patent. Second, the outcome equation coefficient on the small entity indicator in the general model is 9% larger than the respective coefficient in the traditional sample selection model. This implies that traditional sample selection models understate the negative effects of being a small firm on patent scope. Traditional sample selection models ignore that the small firms who select into patenting must have especially high patent scope relative to small firms who abandon, because patent scope itself affects the patenting decision. Upon controlling for the direct effect, the positive effects of firm size are more apparent.

TABLE 2: Parameters: patent application process

Variable	SSM	New
<i>Patent issuance (D)</i>		
CONSTANT ( $\beta_0$ )	-1.17 (0.01)	-1.63 (0.02)
SMALL ENTITY ( $\beta_X$ )	-0.38 (0.01)	-0.36 (0.01)
GRANT LENIENCY ( $\beta_{-D}$ )	1.96 (0.01)	2.21 (0.03)
PATENT SCOPE ( $\beta_Y$ )	NA	0.30 (0.01)
<i>Patent Scope (Y)</i>		
CONSTANT ( $\gamma_0$ )	-1.21 (0.04)	-1.85 (0.02)
SMALL ENTITY ( $\gamma_X$ )	-0.67 (0.01)	-0.73 (0.01)
CLAIMS LENIENCY ( $\gamma_Y$ )	0.94 (0.01)	1.13 (0.01)
$\rho$	0.93 (0.00)	0.91 (0.00)
$\sigma$	1.74 (0.00)	1.82 (0.01)
$n$	332,199	332,199
$n_1 = \sum D_i$	149,523	149,523

Notes: Table 2 reports estimates from the patent application process model. The “SSM” column provides estimates from using MLE on the standard sample selection model with no  $Y$  in the selection equation. The column labelled “New” presents estimates from the newly proposed estimation method as described in the paper. Standard errors are in parentheses next to coefficients.

## 9 Conclusion

In this paper, I study identification and estimation of the complete set of structural parameters in sample selection models where the outcome directly affects selection. I provide model identification conditions, along with likelihood and two-step regression estimation methods. My simulations, and the empirical example relating to the patent application process, show the importance of including the outcome in the selection equation where appropriate. I recommend if practitioners believe that the outcome affects selection, they should use the methods I describe - even including all excluded regressors  $X_Y$  in the selection equation does not preserve the correlation or selection equation parameters.

This paper generates a couple of avenues for future research. First, as discussed in Appendix B, when the outcome variable  $Y$  is discrete, identification and estimation are more involved. Even with a simplification to the model so that  $Y^*$  enters the selection index  $D^*$  rather than  $Y$ , conditions for identification in this paper do not generalize. Also, though similar to derive, the likelihood requires simulation methods.

Second, throughout I considered *parametric* and *semiparametric* identification and estimation in a *linear* system. Extending the model to nonlinearity in  $Y$  and  $D$ , together with tackling nonparametric identification, is an avenue for further research.<sup>13</sup>

---

<sup>13</sup>See Das, Newey, and Vella (2003) for nonparametric estimation of sample selection models.

## References

The numbers at the end of every reference link to the pages citing the reference.

ABOWD, J. M. AND H. S. FARBER (1982): “Job Queues and the Union Status of Workers,” *ILR Review*, 35, 354–367. [4](#)

AHN, H. AND J. L. POWELL (1993): “Semiparametric estimation of censored selection models with a nonparametric selection mechanism,” *Journal of Econometrics*, 58, 3–29. [4](#)

AMEMIYA, T. (1973): “Regression Analysis when the Dependent Variable is Truncated Normal,” *Econometrica*, 41, 997–1016. [14](#), [33](#)

——— (1974): “Multivariate Regression and Simultaneous Equation Models when the Dependent Variables Are Truncated Normal,” *Econometrica*, 42, 999–1012. [3](#)

——— (1984): “Tobit models: A survey,” *Journal of Econometrics*, 24, 3–61. [3](#)

——— (1985): *Advanced Econometrics*, Harvard University Press. [3](#)

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, 497–517. [4](#)

BASTOS, F. AND W. BARRETO-SOUZA (2021): “Birnbaum–Saunders sample selection model,” *Journal of Applied Statistics*, 48, 1896–1916. [4](#)

BASTOS, F. D. S., W. BARRETO-SOUZA, AND M. G. GENTON (2022): “A Generalized Heckman Model With Varying Sample Selection Bias and Dispersion Parameters,” *Statistica Sinica*. [4](#)

CARLSON, A. (2022): “gtsheckman: Generalized Two Step Heckman Estimator,” *Unpublished Working Paper*. [4](#)

COGAN, J. F. (1981): “Fixed Costs and Labor Supply,” *Econometrica*, 49, 945–963. [4](#)

DAS, M., W. K. NEWEY, AND F. VELLA (2003): “Nonparametric Estimation of Sample Selection Models,” *The Review of Economic Studies*, 70, 33–58. [4](#), [20](#)

FARRE-MENSA, J., D. HEGDE, AND A. LJUNGQVIST (2020): “What Is a Patent Worth? Evidence from the U.S. Patent “Lottery”,” *The Journal of Finance*, 75, 639–682. [18](#)

- FENG, J. AND X. JARAVEL (2020): “Crafting Intellectual Property Rights: Implications for Patent Assertion Entities, Litigation, and Innovation,” *American Economic Journal: Applied Economics*, 12, 140–81. [3](#)
- FRANDBSEN, B. R., L. J. LEFGREN, AND E. C. LESLIE (2019): “Judging Judge Fixed Effects,” *NBER Working Paper Series*. [18](#)
- GAULÉ, P. (2018): “Patents and the Success of Venture-Capital Backed Startups: Using Examiner Assignment to Estimate Causal Effects,” *The Journal of Industrial Economics*, 66, 350–376. [18](#)
- GORDON, R. H. AND A. S. BLINDER (1980): “Market wages, reservation wages, and retirement decisions,” *Journal of Public Economics*, 14, 277–308. [4](#)
- GOURIEROUX, C. (2000): *Econometrics of Qualitative Dependent Variables*, Cambridge University Press. [3](#)
- GRAHAM, S. J., A. C. MARCO, AND R. MILLER (2018): “The USPTO Patent Examination Research Dataset: A window on patent processing,” *Journal of Economics & Management Strategy*, 27, 554–578. [6](#)
- GREENE, W. (2003): *Econometric Analysis*, Pearson Education. [27](#), [28](#)
- GRILICHES, Z., B. H. HALL, AND J. A. HAUSMAN (1978): “Missing Data and Self-Selection in Large Panels,” *Annales de l’insée*, 137–176. [4](#)
- GRONAU, R. (1974): “Wage Comparisons-A Selectivity Bias,” *Journal of Political Economy*, 82, 1119–43. [3](#)
- HANOCH, G. (1976): *A Multivariate Model of Labor Supply: Methodology for Estimation*, Santa Monica, CA: RAND Corporation. [4](#)
- (2014): *Chapter 3. Hours And Weeks In The Theory Of Labor Supply*., Princeton University Press, 119–165. [4](#)
- HARDLE, W., H. LIANG, AND J. GAO (2002): *Partially Linear Models*, Physica Heidelberg. [12](#), [15](#)
- HECKMAN, J. (1974): “Shadow Prices, Market Wages, and Labor Supply,” *Econometrica*, 42, 679–94. [3](#), [4](#)

- (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, Volume 5, number 4*, 475–492. [2](#), [3](#), [15](#)
- (1978): “Dummy Endogenous Variables in a Simultaneous Equation System,” *Econometrica*, 46, 931–59. [3](#)
- (1990): “Varieties of Selection Bias,” *The American Economic Review: Papers and Proceedings*, 80, 313–318. [3](#)
- HECKMAN, J. J. (1979): “Sample Selection Bias as a Specification Error,” *Econometrica*, 47, 153–161. [3](#), [4](#)
- HOROWITZ, J. L. (2010): *Semiparametric and Nonparametric Methods in Econometrics*, Springer Series in Statistics, Springer-Verlag. [11](#), [12](#), [31](#)
- JENNRICH, R. I. (1969): “Asymptotic Properties of Non-Linear Least Squares Estimators,” *The Annals of Mathematical Statistics*, 40, 633 – 643. [14](#), [33](#)
- KATZ, A. (1977): “Evaluating Contributions of the Employment Service to Applicant Earnings,” *Labor Law Journal*, 28, 472–8. [4](#)
- KENNY, L. W., L.-F. LEE, G. S. MADDALA, AND R. P. TROST (1979): “Returns to College Education: An Investigation of Self-Selection Bias Based on the Project Talent Data,” *International Economic Review*, 20, 775–789. [4](#)
- KING, M. A. (1980): “An Econometric Model of Tenure Choice and Demand for Housing as a Joint Decision,” in *Econometric Studies in Public Finance*, National Bureau of Economic Research, Inc, NBER Chapters, 137–159. [4](#)
- KLEIN, R. W. AND R. H. SPADY (1993): “An Efficient Semiparametric Estimator for Binary Response Models,” *Econometrica*, 61, 387–421. [15](#)
- KLING, J. R. (2006): “Incarceration Length, Employment, and Earnings,” *American Economic Review*, 96, 863–876. [18](#)
- KOMAROVA, T. AND W. MATCHAM (2022): “Multivariate ordered discrete response models,” *Unpublished Working Paper*. [11](#), [12](#)

- KUHN, J. AND N. THOMPSON (2019): “How to Measure and Draw Causal Inferences with Patent Scope,” *International Journal of the Economics of Business*, 26, 5–38. [3](#)
- LANJOUW, J. O. AND M. SCHANKERMAN (2004): “Patent Quality and Research Productivity: Measuring Innovation with Multiple Indicators,” *The Economic Journal*, 114, 441–465. [3](#), [18](#)
- LEE, L.-F. (1978): “Unionism and Wage Rates: A Simultaneous Equations Model with Qualitative and Limited Dependent Variables,” *International Economic Review*, 19, 415–33. [4](#)
- LEWIS, H. G. (1974): “Comments on Selectivity Biases in Wage Comparisons,” *Journal of Political Economy*, 82, 1145–55. [3](#)
- MADDALA, G. S. (1986): *Limited-Dependent and Qualitative Variables in Econometrics*, Cambridge University Press. [3](#)
- MANSKI, C. (1988): “Identification of Binary Response Models,” *Journal of the American Statistical Association*, 83, 729–738. [11](#), [12](#), [30](#)
- MANSKI, C. F. (1985): “Semiparametric analysis of discrete response: Asymptotic properties of the maximum score estimator,” *Journal of Econometrics*, 27, 313–333. [11](#), [12](#)
- MARCHENKO, Y. V. AND M. G. GENTON (2012): “A Heckman Selection-t Model,” *Journal of the American Statistical Association*, 107, 304–317. [4](#)
- MATCHAM, W. AND M. SCHANKERMAN (2022): “On the Validity of the Leniency IV,” *Working Paper*. [4](#)
- NAKOSTEEN, R. A. AND M. ZIMMER (1980): “Migration and Income: The Question of Self-Selection,” *Southern Economic Journal*, 46, 840–851. [4](#)
- NELSON, F. D. (1977): “Censored regression models with unobserved, stochastic censoring thresholds,” *Journal of Econometrics*, 6, 309–327. [4](#)
- OGUNDIMU, E. O. AND J. L. HUTTON (2016): “A Sample Selection Model with Skew-normal Distribution,” *Scandinavian Journal of Statistics*, 43, 172–190. [4](#)
- POIRIER, D. J. AND P. A. RUUD (1981): “On the appropriateness of endogenous switching,” *Journal of Econometrics*, 16, 249–256. [4](#)



- ROBERTS, R. B., G. MADDALA, AND G. ENHOLM (1978): “Determinants of the Requested Rate of Return and the Rate of Return Granted in a Formal Regulatory Process,” *Bell Journal of Economics*, 9, 611–621. [4](#)
- ROBINSON, P. M. (1988): “Root-N-Consistent Semiparametric Regression,” *Econometrica*, 56, 931–954. [12](#), [15](#)
- ROELSGAARD, S. AND L. TAYLOR (2022): “A Semiparametric Machine Learning Estimator for Sample Selection Models,” *Unpublished Working Paper*. [4](#)
- ROSEN, H. S. (1979): “Housing decisions and the U.S. income tax: An econometric analysis,” *Journal of Public Economics*, 11, 1–23. [4](#)
- SAMPAT, B. AND H. L. WILLIAMS (2019): “How Do Patents Affect Follow-On Innovation? Evidence from the Human Genome,” *American Economic Review*, 109, 203–36. [18](#)
- TROST, R. (1977): “Demand for Housing: A Model Based on Inter-related Choices Between Owning and Renting,” *Unpublished Ph.D dissertation, University of Florida*. [4](#)
- VELLA, F. (1998): “Estimating Models with Sample Selection Bias: A Survey,” *Journal of Human Resources*, 33, 127–169. [3](#)
- WEISBROD, B. (1980): *Wage Differentials Between the Private For-profit and Non-profit Sectors: The Case of Lawyers*, Discussion papers, University of Wisconsin-Madison, Institute for Research on Poverty. [4](#)
- WILLIS, R. J. AND S. ROSEN (1979): “Education and Self-Selection,” *Journal of Political Economy*, 87, S7–S36. [4](#)

# A Proofs and Derivations

## A.1 Proof of Theorem 1

*Proof.* The proof contains four parts:

1. Identify the standardized reduced form parameters  $\pi$  as given in (11)
2. Identify  $\gamma$
3. Identify  $\sigma_Y$
4. Identify  $\beta_Y, \rho$ , and  $\beta_{-Y}$

Each step works sequentially: I identify the standardized reduced form parameters first not because they are of interest *per se*, but because they aid in the identification of  $\gamma$ ,  $\sigma_Y, \rho$ , and eventually  $\beta$ .

### Part 1: Reduced Form Parameters

From (10),

$$\mathbb{P}(D = 1|Z) = \mathbb{P}\left(\frac{\beta_Y \varepsilon_Y + \varepsilon_D}{\sqrt{v}} > \tilde{Z} \middle| Z\right) = \Phi(-\tilde{Z}) = \Phi(\pi'Z),$$

so

$$\Phi^{-1}[\mathbb{P}(D = 1|Z)] = \pi'Z.$$

Let  $Z$  have dimension  $K_1$ . Then, using  $K_1$  different vectors of  $Z$  each denoted  $Z_{(k)}$  for  $k = 1, \dots, K_1$ , there are  $K_1$  equations

$$\underbrace{\Phi^{-1}[\mathbb{P}(D = 1|Z_{(k)})]}_{q_{(k)}} = \pi'Z_{(k)}.$$

Let  $Q$  stack  $q_{(k)}$  and  $\mathcal{X}_1$  stack  $Z_{(k)}$ . Then

$$Q = \mathcal{X}_1 \pi.$$

The matrices  $Q$  and  $\mathcal{X}_1$  are known from the distribution of observables. By the assumptions of the theorem,  $\mathcal{X}_1$  has full rank and thus  $\pi$  is identified. The standardized reduced form parameters  $\pi$  are the ones recovered from a probit regression of  $D$  on  $Z$ .

## Part 2: $\gamma$ Parameters

Now I identify  $\gamma$ , given that I have identified  $\pi$ . Consider the conditional expectation of  $Y$ , conditional on  $D = 1$  and  $Z$ , i.e.  $\mathbb{E}(Y|D = 1, Z)$ . Because I condition on  $D = 1$ , this object is known from the distribution of observables. Substituting  $Y$  from (5),

$$\mathbb{E}(Y|D = 1, Z) = \gamma'w_Y + \mathbb{E}(\varepsilon_Y|D = 1, Z), \quad (17)$$

where, for reference,  $w_Y = (X'_Y, X')'$ . Recall that  $D = 1$  if and only if  $\frac{\beta_Y\varepsilon_Y + \varepsilon_D}{\sqrt{v}} > \tilde{Z}$  where  $\tilde{Z}$  is identified because  $Z$  is observed and  $\pi$  is identified. So, owing to the joint normality of  $\varepsilon_Y$  and  $\frac{\beta_Y\varepsilon_Y + \varepsilon_D}{\sqrt{v}}$ , there holds that (Greene, 2003)

$$\begin{aligned} \mathbb{E}(\varepsilon_Y|D = 1, Z) &= \mathbb{E}\left(\varepsilon_Y \left| \frac{\beta_Y\varepsilon_Y + \varepsilon_D}{\sqrt{v}} > \tilde{Z}, Z\right.\right) \\ &= \text{cov}\left(\varepsilon_Y, \frac{\beta_Y\varepsilon_Y + \varepsilon_D}{\sqrt{v}}\right) \frac{\phi(\tilde{Z})}{1 - \Phi(\tilde{Z})} \\ &= \underbrace{\frac{\beta_Y\sigma_Y^2 + \rho\sigma_Y}{\sqrt{v}}}_{\tau} \underbrace{\frac{\phi(\tilde{Z})}{\Phi(-\tilde{Z})}}_{\lambda(\tilde{Z})}, \end{aligned} \quad (18)$$

where  $\lambda(\tilde{Z})$  is the inverse Mills ratio. I also use the notation  $\tilde{\lambda}(Z)$  to denote the same quantity, i.e.

$$\tilde{\lambda}(Z) = \lambda(\tilde{Z}) = \frac{\phi(\tilde{Z})}{\Phi(-\tilde{Z})}$$

Putting (17) together with (18) yields

$$\mathbb{E}(Y|D = 1, Z) = \gamma'w_Y + \tau\tilde{\lambda}$$

This implies that a regression of  $Y$  on  $w_Y$  and  $\tilde{\lambda}$  for observations with  $D = 1$  delivers  $\gamma$  and  $\tau$ . More formally, define

$$U = Y - \gamma'w_Y - \tau\tilde{\lambda}.$$

Then

$$Y = \gamma'w_Y + \tau\tilde{\lambda} + U, \quad \mathbb{E}(U|D = 1, w_Y, \tilde{\lambda}) = 0. \quad (19)$$

Now, (19) implies by the Tower property of conditional expectations that

$$\mathbb{E}(Y|D = 1, \tilde{\lambda}) = \gamma'\mathbb{E}(w_Y|D = 1, \tilde{\lambda}) + \tau\tilde{\lambda}. \quad (20)$$

Subtracting (20) from (19) yields

$$\underbrace{Y - \mathbb{E}(Y|D = 1, \tilde{\lambda})}_{\tilde{Y}} = \gamma' \left[ \underbrace{w_Y - \mathbb{E}(w_Y|D = 1, \tilde{\lambda})}_{\tilde{w}_Y} \right] + U \quad (21)$$

As indicated in (21), defining  $\tilde{Y} = Y - \mathbb{E}(Y|D = 1, \tilde{Z})$  and  $\tilde{w}_Y = w_Y - \mathbb{E}(w_Y|D = 1, \tilde{Z})$  implies the standard linear regression model

$$\tilde{Y} = \gamma' \tilde{w}_Y + U \quad (22)$$

The well-known identification conditions apply, namely that

$$\gamma = \mathbb{E}(\tilde{w}_Y \tilde{w}_Y')^{-1} \mathbb{E}(\tilde{w}_Y \tilde{Y})$$

provided that  $\mathbb{E}(\tilde{w}_Y \tilde{w}_Y')$  is positive definite, which is guaranteed by the assumptions in the theorem, namely (12). Hence  $\gamma$  is identified.

### Part 3: $\sigma_Y$ parameter

It is natural next to look at  $\text{Var}(Y|D = 1, Z)$ . Using properties of conditional multivariate normals (Greene, 2003),

$$\text{Var}(Y|D = 1, Z) = \sigma_Y^2 + \tau^2 \tilde{Z} \lambda(\tilde{Z}) - \tau^2 \lambda(\tilde{Z})^2.$$

From this  $\sigma_Y^2$  is identified after rearrangement.

### Part 4: $\beta_Y, \rho$ and $\beta_{-Y}$

From (11) it follows that

$$\begin{aligned} \beta_D &= \pi_D \sqrt{v} \\ \beta_X &= \pi_X \sqrt{v} - \beta_Y \gamma_X. \end{aligned}$$

So, after showing that  $\beta_Y$  and  $\rho$  are identified (which implies that  $v = \beta_Y^2 \sigma_Y^2 + 1 + 2\beta_Y \sigma_Y \rho$  is identified), it is immediate that  $\beta_{-Y}$  is identified. Hence to finish the proof, I must show that  $\beta_Y$  and  $\rho$  are identified.

It follows from (11) and the assumptions of the theorem that for some  $j$  with  $\gamma_{Yj} \neq 0$ ,

$$\frac{\beta_Y}{\sqrt{v}} = \frac{\pi_{Yj}}{\gamma_{Yj}} = \xi. \quad (23)$$

Additionally, recall that  $\tau$  is identified, where

$$\tau = \frac{\beta_Y \sigma_Y^2 + \rho \sigma_Y}{\sqrt{v}}. \quad (24)$$

First, suppose  $\xi = 0$  so that  $\beta_Y = 0$ . This implies that  $v = 1$  and  $\tau = \rho \sigma_Y$  or  $\rho = \tau / \sigma_Y$ . So when  $\xi = 0$ ,  $\beta_Y$  and  $\rho$  are identified. From now on, suppose  $\xi \neq 0$  so that  $\beta_Y \neq 0$ . Solving (23) and (24) simultaneously yields<sup>14</sup>

$$\beta_Y^2 = \frac{\xi^2}{1 + \sigma_Y^2 \xi^2 - 2\tau \xi}$$

and

$$\rho = \beta_Y \left( \frac{\tau}{\sigma_Y \xi} - \sigma_Y \right).$$

From this, we know that

$$\beta_Y = \pm \sqrt{\frac{\xi^2}{1 + \sigma_Y^2 \xi^2 - 2\tau \xi}}.$$

But, since  $\beta_Y = \xi \sqrt{v}$ , the sign of  $\beta_Y$  is the same as the sign of  $\xi$ . This means that  $\beta_Y$  is identified and therefore so is  $\rho$ .

□

[Back to Theorem 1](#)

## A.2 Proof of Theorem 2

*Proof.* The proof proceeds in three steps:

1. Identify the reduced form parameters  $\kappa$  of the selection equation in (7)
2. Identify  $\gamma$
3. Identify  $\beta_Y, \rho$ , and  $\beta_{-Y}$

### Part 1: Reduced Form Parameters

Rearranging (7) yields

$$\mathbb{P}(D = 0|Z) = \mathbb{P}(\tilde{\varepsilon} \leq -\kappa'Z) = F_{\tilde{\varepsilon}}(-\kappa'Z) = F_{\tilde{\varepsilon}}(\check{Z}) \quad (25)$$

---

<sup>14</sup>Note that  $1 + \sigma_Y^2 \xi^2 - 2\tau \xi = v^{-1}$  so this term cannot be nonpositive.

where  $\kappa = (\kappa'_D, \kappa'_Y, \kappa'_X)'$  and  $\check{Z} = -\kappa'Z$ . Equation (8) defines the terms in  $\kappa$ .

The assumptions of the theorem guarantee that  $\mathbb{P}(D = 0|Z)$  will not be degenerate for all  $Z \in S$  (in the sense that it will not be values of 0 or 1 only). Equation (25) forms the basis of the identification strategy. By the assumptions on  $(\varepsilon_Y, \varepsilon_D)$ , the function  $F_{\varepsilon}$  is strictly monotone. Hence, for any  $Z, \hat{Z} \in S$ ,

$$\mathbb{P}(D = 0|\hat{Z}) > \mathbb{P}(D = 0|Z) \iff \kappa'\hat{Z} < \kappa'Z. \quad (26)$$

The identification follows [Manski \(1988\)](#), similar to single-index models with a monotone link function.

I prove the identification of  $\kappa$  by contradiction. Towards a contradiction, assume  $k \neq \kappa$  are both consistent with observables. Both “parameters” are such that  $\kappa_1 = \beta_{D1} = 1$  and  $k_1 = b_{D1} = 1$ . Note that, conveniently,  $Z_1 = X_{D1}$ . By the theorem, there exists a positive probability set of  $Z_{-1}^0 \in S_{-1}$  such that  $\kappa'_{-1}Z_{-1}^0 \neq k'_{-1}Z_{-1}^0$ . Without loss of generality, assume

$$\kappa'_{-1}Z_{-1}^0 > k'_{-1}Z_{-1}^0$$

Then for any  $Z_1^0 = X_{D1}^0$  that complements  $Z_{-1}^0$  to a point in  $S$ , there holds that

$$X_{D1}^0 + \kappa'_{-1}Z_{-1}^0 > X_{D1}^0 + k'_{-1}Z_{-1}^0$$

In particular, take an  $X_{D1}^0 \in (\underline{X}_{D1}, \bar{X}_{D1})$ . Then, due to the continuity of  $X_{D1}$  on  $(\underline{X}_{D1}, \bar{X}_{D1})$ , there exists  $\hat{X}_{D1}^0 \neq X_{D1}^0$  such that  $(\hat{X}_{D1}^0, Z_{-1}^0) \in S$  and the two inequalities

$$X_{D1}^0 + \kappa'_{-1}Z_{-1}^0 > \hat{X}_{D1}^0 + \kappa'_{-1}Z_{-1}^0 \quad (27)$$

and

$$\hat{X}_{D1}^0 + \kappa'_{-1}Z_{-1}^0 > X_{D1}^0 + k'_{-1}Z_{-1}^0 \quad (28)$$

hold. Yet, because of (26), since  $\kappa$  is consistent with observables, inequality (27) implies that

$$\mathbb{P}(D = 0|(X_{D1}^0, Z_{-1}^0)) < \mathbb{P}(D = 0|(\hat{X}_{D1}^0, Z_{-1}^0)). \quad (29)$$

Similarly, since  $k$  is consistent with observables, inequality (28) implies that

$$\mathbb{P}(D = 0|(\hat{X}_{D1}^0, Z_{-1}^0)) < \mathbb{P}(D = 0|(X_{D1}^0, Z_{-1}^0)). \quad (30)$$

Inequalities (29) and (30) together imply a contradiction, so  $\kappa$  is identified.

## Part 2: $\gamma$ parameters

Now I identify  $\gamma$ , given that I have identified  $\kappa$ . Consider the conditional expectation of  $Y$ , conditional on  $D = 1$ ,  $w_Y$ , and  $\check{Z}$ , i.e.  $\mathbb{E}(Y|D = 1, w_Y, \check{Z})$ . Because I condition on  $D = 1$ , this object is known from the distribution of observables. Substituting  $Y$  from (5),

$$\mathbb{E}(Y|D = 1, w_Y, \check{Z}) = \gamma'w_Y + \mathbb{E}(\varepsilon_Y|D = 1, w_Y, \check{Z}).$$

Recall that  $D = 1$  if and only if  $\tilde{\varepsilon} > \check{Z}$  where  $\check{Z}$  is identified because  $Z$  is observed and  $\kappa$  is identified. As such,  $\mathbb{E}(\varepsilon_Y|D = 1, w_Y, \check{Z}) = \mathcal{G}(\check{Z})$ . The argument now follows from [Horowitz \(2010\)](#).

Define

$$U = Y - \gamma'w_Y - \mathcal{G}(\check{Z}).$$

Then

$$Y = \gamma'w_Y + \mathcal{G}(\check{Z}) + U, \quad \mathbb{E}(U|D = 1, w_Y, \check{Z}) = 0. \quad (31)$$

Now, (31) implies by the Tower property of conditional expectations that

$$\mathbb{E}(Y|D = 1, \check{Z}) = \gamma'\mathbb{E}(w_Y|D = 1, \check{Z}) + \mathcal{G}(\check{Z}). \quad (32)$$

Subtracting (32) from (31) yields

$$\underbrace{Y - \mathbb{E}(Y|D = 1, \check{Z})}_{\check{Y}} = \gamma' \underbrace{[w_Y - \mathbb{E}(w_Y|D = 1, \check{Z})]}_{\check{w}_Y} + U \quad (33)$$

As indicated in (33), defining  $\check{Y} = Y - \mathbb{E}(Y|D = 1, \check{Z})$  and  $\check{w}_Y = w_Y - \mathbb{E}(w_Y|D = 1, \check{Z})$  implies the standard linear regression model

$$\check{Y} = \gamma'\check{w}_Y + U \quad (34)$$

The well-known identification conditions apply, namely that

$$\gamma = \mathbb{E}(\check{w}_Y\check{w}_Y')^{-1} \mathbb{E}(\check{w}_Y\check{Y})$$

provided that  $\mathbb{E}(\check{w}_Y\check{w}_Y')$  is positive definite, which is guaranteed by the assumptions in the theorem, namely (13). Hence  $\gamma$  is identified.

### Part 3: $\beta$ parameters

Since  $\beta_D = \kappa_D$ ,  $\beta_D$  is identified. From (8) it follows that

$$\beta_X = \kappa_X - \beta_Y \gamma_X$$

So, after showing that  $\beta_Y$  is identified, it is immediate that  $\beta_X$  is identified. Hence to finish the proof, I must show that  $\beta_Y$  is identified.

It follows from (8) and the assumptions of the theorem that for some  $j$  with  $\gamma_{Yj} \neq 0$ ,

$$\beta_Y = \frac{\kappa_{Yj}}{\gamma_{Yj}}.$$

Therefore,  $\beta_Y$  and immediately  $\beta_X$  are identified, concluding the proof. □

[Back to Theorem 2](#)

### A.3 Likelihood Term

For brevity I omit conditioning on  $Z_i$ . Let  $M_i = \beta'_{-Y} w_{Di} + \beta_Y \gamma'_Y w_{Yi}$ . Then

$$\begin{aligned} \wp_{1i} &= \int_{-M_i - \beta_Y(Y_i - w'_{Yi}\gamma)}^{\infty} f_{\varepsilon_D, \varepsilon_Y}(\varepsilon_D, Y_i - w'_{Yi}\gamma) d\varepsilon_D \\ &= \int_{-M_i - \beta_Y(Y_i - w'_{Yi}\gamma)}^{\infty} f_{\varepsilon_D|\varepsilon_Y}(\varepsilon_D|Y_i - w'_{Yi}\gamma) f_{\varepsilon_Y}(Y_i - w'_{Yi}\gamma) d\varepsilon_D \\ &= \frac{1}{\sigma_Y} \phi\left(\frac{Y_i - w'_{Yi}\gamma}{\sigma_Y}\right) \int_{-M_i - \beta_Y(Y_i - w'_{Yi}\gamma)}^{\infty} \phi\left(\frac{\varepsilon_D - \frac{\rho}{\sigma_Y}(Y_i - w'_{Yi}\gamma)}{\sqrt{1 - \rho^2}}\right) d\varepsilon_D \quad (35) \\ &= \frac{1}{\sigma_Y} \phi\left(\frac{Y_i - w'_{Yi}\gamma}{\sigma_Y}\right) \left[1 - \Phi\left(\frac{-\beta'_{-Y} w_{Di} - (\beta_Y + \frac{\rho}{\sigma_Y})Y_i + \frac{\rho}{\sigma_Y} w'_{Yi}\gamma}{\sqrt{1 - \rho^2}}\right)\right] \\ &= \frac{1}{\sigma_Y} \phi\left(\frac{Y_i - w'_{Yi}\gamma}{\sigma_Y}\right) \Phi\left(\frac{\beta'_{-Y} w_{Di} + (\beta_Y + \frac{\rho}{\sigma_Y})Y_i - \frac{\rho}{\sigma_Y} w'_{Yi}\gamma}{\sqrt{1 - \rho^2}}\right). \end{aligned}$$

Equality (35) follows from

$$\varepsilon_D|\varepsilon_Y = s \sim \mathcal{N}\left(-\frac{\rho}{\sigma_Y}s, 1 - \rho^2\right)$$

when  $(\varepsilon_D, \varepsilon_Y)$  are joint normal as in (9).

[Back to likelihood](#)



## A.4 Two-Step Parametric Estimation Details

Rewriting (16) with an estimated value of  $\lambda_i$  yields

$$Y_i = \gamma' w_{Y_i} + \tau \hat{\lambda}_i + \underbrace{\tau(\lambda_i - \hat{\lambda}_i)}_{\check{V}_i} + V_i,$$

so that the effective error term is  $\check{V}_i = \tau(\lambda_i - \hat{\lambda}_i) + V_i$ .

As mentioned,  $\pi$  is estimated using all  $n$  observations by a MLE probit, and so

$$\sqrt{n}(\hat{\pi} - \pi) \xrightarrow{d} \mathcal{N}(0, V_{1st}).$$

Since  $\lambda_i$  is a twice-continuously differentiable function of  $\pi$ , by the Delta-Method,

$$\sqrt{n}(\hat{\lambda}_i - \lambda_i) \xrightarrow{d} \mathcal{N}(0, \Sigma_i)$$

where  $\Sigma_i = (\partial \lambda_i)^2 Z_i' V_{1st} Z_i$ , for  $\partial \lambda_i = \partial \lambda_i / \partial \tilde{Z}_i$ .

The task is to derive the asymptotic distribution of

$$\sqrt{n_1} \begin{pmatrix} \hat{\gamma} - \gamma \\ \hat{\tau} - \tau \end{pmatrix} = \begin{pmatrix} \frac{1}{n_1} \sum_i w_{Y_i} w'_{Y_i} & \frac{1}{n_1} \sum_i w_{Y_i} \hat{\lambda}_i \\ \frac{1}{n_1} \sum_i w_{Y_i} \hat{\lambda}_i & \frac{1}{n_1} \sum_i \hat{\lambda}_i^2 \end{pmatrix}^{-1} \begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_i w_{Y_i} \check{V}_i \\ \frac{1}{\sqrt{n_1}} \sum_i \hat{\lambda}_i \check{V}_i \end{pmatrix}.$$

All  $n$  observations are used to estimate the probit, but the regression in the second step only uses  $n_1$  observations. As such, I take probability limits are taken with both  $n, n_1 \rightarrow \infty$  and  $n_1/n \rightarrow k \in (0, 1)$ . Under conditions on regressors discussed in [Amemiya \(1973\)](#) and [Jennrich \(1969\)](#),

$$\text{plim}_{n_1 \rightarrow \infty} \begin{pmatrix} \frac{1}{n_1} \sum_i w_{Y_i} w'_{Y_i} & \frac{1}{n_1} \sum_i w_{Y_i} \hat{\lambda}_i \\ \frac{1}{n_1} \sum_i w'_{Y_i} \hat{\lambda}_i & \frac{1}{n_1} \sum_i \hat{\lambda}_i^2 \end{pmatrix}^{-1} = \text{plim}_{n_1 \rightarrow \infty} \begin{pmatrix} \frac{1}{n_1} \sum_i w_{Y_i} w'_{Y_i} & \frac{1}{n_1} \sum_i w_{Y_i} \lambda_i \\ \frac{1}{n_1} \sum_i w'_{Y_i} \lambda_i & \frac{1}{n_1} \sum_i \lambda_i^2 \end{pmatrix}^{-1} = B$$

and

$$\begin{pmatrix} \frac{1}{\sqrt{n_1}} \sum_i w_{Y_i} \check{V}_i \\ \frac{1}{\sqrt{n_1}} \sum_i \hat{\lambda}_i \check{V}_i \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, \psi),$$

where

$$\psi = \text{plim}_{\substack{n \rightarrow \infty \\ n_1 \rightarrow \infty}} \left[ \frac{\sigma_Y^2}{n_1} \begin{pmatrix} \sum_i w_{Y_i} w'_{Y_i} \eta_i & \sum_i w_{Y_i} \lambda_i \eta_i \\ \sum_i w'_{Y_i} \lambda_i \eta_i & \sum_i \lambda_i^2 \eta_i \end{pmatrix} + \tau^2 \left( \frac{n_1}{n} \right) \begin{pmatrix} \frac{1}{n_1^2} \sum_i \sum_j w_{Y_i} w'_{Y_j} \vartheta_{ij} & \frac{1}{n_1^2} \sum_i \sum_j w_{Y_i} \Upsilon_{ij} \\ \frac{1}{n_1^2} \sum_i \sum_j w'_{Y_i} \Upsilon_{ij} & \frac{1}{n_1^2} \sum_i \sum_j \Omega_{ij} \end{pmatrix} \right]$$

where

$$\begin{aligned}\eta_i &= 1 + \frac{\tau^2(\tilde{Z}_i\lambda_i - \lambda_i^2)}{\sigma_Y^2} \\ \vartheta_{ij} &= \partial\lambda_i \partial\lambda_j Z_i' V_{1st} Z_j \\ \Upsilon_{ij} &= \lambda_i \vartheta_{ij} \\ \Omega_{ij} &= \lambda_i \lambda_j \vartheta_{ij}.\end{aligned}$$

Putting these two results together, and using the Cramer convergence theorem,

$$\sqrt{n_1} \begin{pmatrix} \hat{\gamma} - \gamma \\ \hat{\tau} - \tau \end{pmatrix} \xrightarrow{d} \mathcal{N}(0, B\psi B'),$$

as required.

## B Discrete Response

### B.1 Model and Identification

When  $Y$  is discrete, say binary for simplicity, the model is

$$\begin{aligned}Y &= 1 \left[ \underbrace{\gamma'_X X + \gamma'_Y X_Y}_{Y^*} + \varepsilon_Y > 0 \right] = 1 [\gamma' w_Y + \varepsilon_Y > 0] \\ D &= 1 \left[ \underbrace{\beta'_X X + \beta'_D X_D + \beta_Y Y}_{D^*} + \varepsilon_D > 0 \right] = 1 [\beta_{-Y} w_D + \beta_Y Y + \varepsilon_D > 0]\end{aligned}$$

Where all terms are defined as in the main text. Now, when  $D = 0$  we observe nothing on  $Y$ , and when  $D = 1$  we observe  $Y = 1$  or  $Y = 0$ .

Identification will revolve around two observable objects:

1.  $\mathbb{P}(D = 0|Z)$
2.  $\mathbb{P}(D = 1 \cap Y = 1|Z)$  (or alternatively  $\mathbb{P}(Y = 1|D = 1, Z)$ )

Approaching identification as in Theorem 1 will not work because  $Y$  is not linear in parameters. Changing the model so that  $Y^*$  enters  $D^*$ , i.e.

$$Y = 1 \left[ \underbrace{\gamma'_X X + \gamma'_Y X_Y + \varepsilon_Y}_{Y^*} > 0 \right] = 1 [\gamma' w_Y + \varepsilon_Y > 0] \quad (36)$$

$$D = 1 \left[ \underbrace{\beta'_X X + \beta'_D X_D + \beta_Y Y^* + \varepsilon_D}_{D^*} > 0 \right] = 1 [\beta_{-Y} w_D + \beta_Y Y^* + \varepsilon_D > 0] \quad (37)$$

makes progress possible, but identification is still not immediate. Though a probit regression for the  $D$  equation delivers standardized reduced form parameters, a probit regression for the  $Y$  equation for  $D = 1$  does not deliver the  $\gamma$  parameters, and there are no higher moments of observables to identify the error correlation.

## B.2 Estimation

The log likelihood for the model (36) - (37) with conditionally standard bivariate normal errors (with correlation  $\rho$ ) is

$$\log(\mathcal{L}_{\text{disc}}) = \sum_i (1 - D_i) \log(\wp_{0i}) + D_i Y_i \log(\wp_{11i}) + D_i (1 - Y_i) \log(\wp_{10i}),$$

where  $\wp_{0i} = 1 - \Phi(\pi' Z_i)$  is the same as in the main text and

$$\begin{aligned} \wp_{11i} &= \mathbb{P}(D_i = 1 \cap Y_i = 1 | Z_i) \\ &= \underbrace{\mathbb{P}(Y_i = 1 | D_i = 1, Z_i)}_{\wp_{1|1i}} \underbrace{\mathbb{P}(D_i = 1 | Z_i)}_{\wp_{1i}}. \end{aligned}$$

Here

$$\wp_{1i} = \Phi(\pi' Z_i),$$

Therefore, the log-likelihood is

$$\begin{aligned} \log(\mathcal{L}_{\text{disc}}) &= \sum_i (1 - D_i) \log(\wp_{0i}) + D_i \log(\wp_{1i}) \\ &\quad + \sum_i D_i Y_i \log(\wp_{1|1i}) + D_i (1 - Y_i) \log(\wp_{0|1i}) \end{aligned}$$

The term  $\wp_{0|1i}$  is given by

$$\begin{aligned}
\wp_{0|1i} &= \mathbb{P} \left( \gamma' w_{Yi} + \varepsilon_{Yi} < 0 \mid \frac{\beta_Y \varepsilon_{Yi} + \varepsilon_{Di}}{\sqrt{v}} > -\pi' Z_i, Z_i \right) \\
&= \int_{-\pi' Z_i}^{\infty} \mathbb{P} \left( \varepsilon_{Yi} < -\gamma' w_{Yi} \mid \frac{\beta_Y \varepsilon_{Yi} + \varepsilon_{Di}}{\sqrt{v}} = x, Z_i \right) \phi(x) dx \\
&= \int_{-\pi' Z_i}^{\infty} \Phi \left( \frac{-\gamma' w_{Yi} - \vartheta x}{\sqrt{1 - \rho^2}} \right) \frac{\phi(x)}{1 - \Phi(-\pi' Z_i)} dx,
\end{aligned}$$

where  $\vartheta = v^{-1/2}(\beta_Y + \rho)$  (the same as  $\tau$  previously except with  $\sigma_Y = 1$ ) and similarly

$$\wp_{1|1i} = \int_{-\pi' Z_i}^{\infty} \left[ 1 - \Phi \left( \frac{-\gamma' w_{Yi} - \vartheta x}{\sqrt{1 - \rho^2}} \right) \right] \frac{\phi(x)}{1 - \Phi(-\pi' Z_i)} dx.$$

These integrals have no closed form but can be approximated with simulation methods. Results from simulated maximum likelihood estimation apply: the number of simulations must grow in proportion to the sample size to ensure that the asymptotic distribution of the estimator matches that of exact MLE.<sup>15</sup>

---

<sup>15</sup>Simulations of this model using Halton draws, which perform moderately well, are available upon request

## C Extra Simulation Results

TABLE 3: Simulation Results Design 2 with Discrete Regressor

---

*Design 2 with Discrete Regressor:  $\bar{D} = 0.84$*

---

<b>Parameter:</b>	$\sigma$	$\rho$	$\gamma_0$	$\gamma_X$	$\gamma_Y$	$\beta_0$	$\beta_X$	$\beta_D$	$\beta_Y$
<b>True Value:</b>	2	0.5	3	-1	-2	1	3	-8	2
<b>(1) New:</b>	2.07 (0.06)	0.50 (0.16)	3.00 (0.07)	-1.00 (0.02)	-2.00 (0.04)	1.08 (0.36)	3.21 (0.68)	-8.58 (1.74)	2.14 (0.47)
<b>(2) Existing:</b>	2.01 (0.04)	0.77 (0.04)	3.25 (0.06)	-1.00 (0.02)	-1.76 (0.05)	1.73 (0.06)	0.16 (0.01)	-1.26 (0.06)	–

---

Notes: Table 3 reports the sample mean and sample standard deviations (in parentheses) of the estimates of the model parameters, over 500 repeated samples. The “New” row provide estimates from using the newly proposed model. The “Existing” row estimates a model that assumes  $\beta_Y = 0$ .

[Back to Simulation Section](#)