

Article

# A Concentrated, Nonlinear Information-Theoretic Estimator for the Sample Selection Model

Amos Golan<sup>1,\*</sup> and Henryk Gzyl<sup>2</sup>

- <sup>1</sup> Department of Economics and the Info-Metrics Institute, American University, Kreeger Hall 104, 4400 Massachusetts Ave., NW, Washington, DC 20016-8029, USA
- <sup>2</sup> Centro de Finanzas, IESA, Caracas, Venezuela; E-Mail: henryk.gzyl@iesa.edu.ve
- \* Author to whom correspondence should be addressed; E-Mail: agolan@american.edu.

Received: 16 April 2010; in revised form: 3 June 2001 / Accepted: 11 June 2010 / Published: 14 June 2010

**Abstract:** This paper develops a semi-parametric, Information-Theoretic method for estimating parameters for nonlinear data generated under a sample selection process. Considering the sample selection as a set of inequalities makes this model inherently nonlinear. This estimator (i) allows for a whole class of different priors, and (ii) is constructed as an unconstrained, concentrated model. This estimator is easy to apply and works well with small or complex data. We provide a number of explicit analytical examples for different priors' structures and an empirical example.

**Keywords:** concentrated model; inequalities; information; maximum entropy; priors; sample selection

# 1. Introduction and Basic Model

The sample selection problem appears often in empirical studies of labor supply, individuals' wages and other topics. For small sample sizes the existing parametric [1] and semi-parametric estimators [2–5] have difficulties. Recently, [6], henceforth GMP, developed a semi-parametric, Information-Theoretic (IT) estimator for the sample-selection problem that performs well when the sample is small. This estimator is based on the IT generalized maximum entropy (GME) approach of [7] and [8]. GMP used a large number of sampling experiments to investigate and compare the small-sample behavior of their estimator relative to other estimators. GMP concluded that their IT estimator is the most stable estimator while the likelihood estimators predicted better within the sample for large enough samples. Their IT estimator outperformed the AP estimator in most cases and in all small samples. Another set of experiments within the nonlinear framework appear in [9].

GMP specified their IT-GME model with bounds on the parameters and with finite and discrete support. Though, their IT-GME estimator performs relatively well, it still has some of the basic shortcomings of that estimator. It has finite and bounded supports for both signal and noise, it is not flexible enough to incorporate infinitely large bounds and continuous support spaces, and it is constructed as a constrained optimization estimator. The objective here is to extend the estimator discussed in GMP in three directions. First, we allow unbounded support spaces for all parameters. Second, we accommodate for a whole class of (discrete and continuous) priors. Third, we construct our estimator as an unconstrained concentrated model.

#### 1.1. The Basic Sample Selection Model

For simplicity, we follow a common labor model discussed in [10]. Suppose individual h (h=1,...,N) values staying (working) at home at wage  $y_{1h}^*$  and can earn  $y_{2h}^*$  in the marketplace. If  $y_{2h}^* > y_{1h}^*$ , the individual works in the marketplace,  $y_{1h} = 1$ , and we observe the market value,  $y_{2h} = y_{2h}^*$ . Otherwise,  $y_{1h} = 0$  and  $y_{2h} = 0$ .

The individual's value at home or in the marketplace depends (linearly) on demographic characteristics  $(\mathbf{x})$ :

$$y_{lh}^* = x_{lh}^t \beta_l + \varepsilon_{lh} \tag{1}$$

$$y_{2h}^* = x_{2h}^t \beta_2 + \varepsilon_{2h} \tag{2}$$

where  $\mathbf{x}_{1h}$  and  $\mathbf{x}_{2h}$  are K<sub>1</sub> and K<sub>2</sub>-dimensional vectors,  $\beta_1$  and  $\beta_2$  are K<sub>1</sub> and K<sub>2</sub>-dimensional vectors of unknowns and "*t*" stands for "transpose". This model can be expressed as

$$y_{1h} = \begin{cases} 1 & \text{if } y_{2h}^* > y_{1h}^* \\ 0 & \text{if } y_{2h}^* \le y_{1h}^* \end{cases}$$
(3)

$$y_{2h} = \begin{cases} x_{2h}^{t} \beta_{2} + \mathcal{E}_{2h} & \text{if } y_{2h}^{*} > y_{1h}^{*} \\ 0 & \text{if } y_{2h}^{*} \le y_{1h}^{*} \end{cases}$$
(4)

Our objective is to estimate  $\beta_1$  and  $\beta_2$ . Typically the researcher is interested primarily in  $\beta_2$ .

Unlike the more traditional models, GMP constructed their model as a solution to a constrained optimization problem such that the information represented by the set of censored equations (3)-(4) enters the estimation as inequalities:

$$x_{2h}^{t} \beta_{2} + \varepsilon_{2h} = y_{2h}, \text{ if } y_{1h} = 1$$
 (5)

$$x_{2h}^{t}\beta_{2} + \varepsilon_{2h} > x_{1h}^{t}\beta_{1} + \varepsilon_{1h}$$
, if  $y_{1h} = 1$  (6)

$$x_{2h}^t \beta_2 + \varepsilon_{2h} \le x_{1h}^t \beta_1 + \varepsilon_{1h}, \quad \text{if } y_{1h} = 0 \tag{7}$$

In our formulation, we use inequalities as well to represent all available information in the set of censored equations.

#### 2. The Information-Theoretic Estimator

Rewrite equations (1)-(2) as finding  $\gamma_1$  and  $\gamma_2$  in  $y_1^* = A_1 \gamma_1$  and  $y_2^* = A_2 \gamma_2$ , where the dependent variable is censored and where  $\gamma_1 = \begin{pmatrix} \beta_1 \\ \epsilon_1 \end{pmatrix}$ ,  $A_1 = [X_1 I]$ ,  $\gamma_2 = \begin{pmatrix} \beta_2 \\ \epsilon_2 \end{pmatrix}$  and  $A_2 = [X_2 I]$ . We formulate the

censored model (5)-(7) in the following way.

Let the constraint sets  $C_i = C_{i,s} \times C_{i,n}$  (i = 1, 2). For each *i*,  $C_{i,s}$  is an auxiliary closed, convex set used to model the a-priori constraints on the  $\beta$ 's. Similarly, the closed convex set  $C_{i,n}$  is part of the specification of the "physical" nature of the noise and contains all possible realizations of  $\varepsilon$ . We view the coordinates  $\zeta_i$  in  $C_{i,s}$  and  $v_i$  in  $C_{i,n}$  as values of random variables distributed according to some probability measure  $dP_i(\zeta_i) \equiv dP_i(\zeta_i, v_i)$  such that their expectations (E) are

$$\beta_i = E_{P_i}[\varsigma_i] \text{ and } \varepsilon_i = E_{P_i}[v_i]$$
 (8)

We note that the qualifier "prior" assigned to the Q probability measures is not the traditional Bayesian view. Rather, the  $Q_s$  is just a mathematical construct to transform the estimation problem into a variational problem. The  $Q_n$ , however, could be viewed as the probability measure describing the statistical nature of the noise. The process of estimation of the noise involves a tilting of the prior measure.

Given some (any) prior measures  $dQ_i(\xi_i) \equiv dQ_i(\zeta_i, v_i) = dQ_{i,s}(\zeta_i)dQ_{i,n}(v_i)$  we search for densities  $\rho_i(\xi_i)$  such that  $dP_i = \rho_i(\xi_i)dQ_i(\xi_i)$  satisfies the system (3)-(4). This yields the parameter estimator  $\boldsymbol{\beta}^*$  and the estimated residuals  $\boldsymbol{\varepsilon}^*$ .

Next, let  $S(P_i, Q_i)$  denotes the differential entropy divergence measure between the priors,  $Q_i$ , and the post-data (posteriors)  $P_i$ . This is just the continuous version of the Kullback-Liebler information divergence measure, also known as relative entropy (see [11–13]). Since the data are naturally divided into observed and unobserved parts, we divide the data into two subsets: *J* and  $J^c$  of  $\{1, 2, ..., N\}$ . Next, rewrite the data (3)-(4)  $y_1^* = A_1\gamma_1$  and  $y_2^* = A_2\gamma_2$  as

$$y_1^* = \begin{pmatrix} y_1 \\ \overline{y}_1 \end{pmatrix} = \begin{pmatrix} B_1 \\ \overline{B}_1 \end{pmatrix} E_{P_1}[\xi_1]; \text{ and } y_2^* = \begin{pmatrix} y_2 \\ \overline{y}_2 \end{pmatrix} = \begin{pmatrix} B_2 \\ \overline{B}_2 \end{pmatrix} E_{P_2}[\xi_2]$$
(9)

where the matrices  $B_1$  and  $B_2$  correspond to the rows of the matrices  $A_i$  (*i*=1, 2) labeled by the indices for which observations are available. For the indices in *J* the values  $y_2$  are observed and  $y_{2h}^* > y_{1h}^*$ , whereas for the values in  $J^c$  all we know is that  $\overline{y}_{1h}^* > \overline{y}_{2h}^*$ .

Our "Basic (Primal) Problem" is the solution to

$$\sup_{(P_1, P_2)} \{ S(P_1, Q_1) + S(P_2, Q_2) | y_2 = B_2 E_{P_2}[\xi_2], B_2 E_{P_2}[\xi_2] > B_1 E_{P_1}[\xi_1]; \overline{B}_2 E_{P_2}[\xi_2] \le \overline{B}_1 E_{P_1}[\xi_1] \}$$
(10)

where the inequalities between vectors are taken to be component wise.

Next, we formulate the problem as a concentrated (unconstrained) entropy problem. To do so, we view the basic primal problem as a two stage problem, call it an "*equivalent primal problem*." In the equivalent model the first stage consists of the standard Generalized Entropy problem (the equality portion of the model) for which a dual can be easily formulated.

The Equivalent primal problem is a solution to the two stage optimization problem:

$$Sup\{Sup\{S(P_{1},Q_{1})+S(P_{2},Q_{2}) | y_{2} = B_{2}E_{P_{2}}[\xi_{2}], \eta_{1} = B_{1}E_{P_{1}}[\xi_{1}]; \\ \overline{\eta}_{2} = \overline{B}_{2}E_{P_{2}}[\xi_{2}]; \overline{\eta}_{1} = \overline{B}_{1}E_{P_{1}}[\xi_{1}]\} | y_{2} > \eta_{1}, \overline{\eta}_{2} \leq \overline{\eta}_{1}\}$$

$$(11)$$

Theorem 2.1. The equivalent primal problem (11) is equivalent to the following (dual) problem

$$Inf\{lnZ(\lambda_1, -\overline{\lambda}, \lambda_2, \overline{\lambda}) + \langle \lambda_1 + \lambda_2, y_2 \rangle | \lambda_1 \in \mathbb{R}^{|J|}_+, \lambda_2 \in \mathbb{R}^{|J|} \text{ and } \overline{\lambda} \in \mathbb{R}^{N-|J|}_+\}$$
(12)

where  $\langle a, b \rangle$  denotes the Euclidean scalar (inner) product of the vectors **a** and **b**,

$$Z(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},\overline{\lambda}_{2}) = \iint_{C_{1}\times C_{2}} e^{-\langle \lambda_{1},B_{1}\xi_{1} \rangle} e^{-\langle \overline{\lambda}_{1},\overline{B}_{1}\xi_{1} \rangle} e^{-\langle \lambda_{2},B_{2}\xi_{2} \rangle} e^{-\langle \overline{\lambda}_{2},\overline{B}_{2}\xi_{2} \rangle} dQ_{1}(\xi_{1}) dQ_{2}(\xi_{2})$$
$$= \int_{C_{1}} e^{-\langle \lambda_{1},B_{1}\xi_{1} \rangle} e^{-\langle \overline{\lambda}_{1},\overline{B}_{1}\xi_{1} \rangle} dQ_{1}(\xi_{1}) \int_{C_{2}} e^{-\langle \lambda_{2},B_{2}\xi_{2} \rangle} e^{-\langle \overline{\lambda}_{2},\overline{B}_{2}\xi_{2} \rangle} dQ_{2}(\xi_{2}) = Z_{1}(\lambda_{1},\overline{\lambda}_{1}) Z_{2}(\lambda_{2},\overline{\lambda}_{2})$$

and  $(\lambda_1, \overline{\lambda}_1, \lambda_2, \overline{\lambda}_2)$  are the four sets of Lagrange multipliers associated with (11). To carry out the procedure specified in (12) first set  $\overline{\lambda}_1 = -\overline{\lambda}_2 = -\overline{\lambda}$ , and then carry out the minimization. **Proof:** See Appendix.

To confirm the uniqueness of the solution to problem (12), observe that the function

$$\ell(\lambda_1, \overline{\lambda}_1, \lambda_2, \overline{\lambda}_2) = \ln Z(\lambda_1, \overline{\lambda}_1, \lambda_2, \overline{\lambda}_2) + \langle \lambda_1, \eta_1 \rangle + \langle \overline{\lambda}_1, \overline{\eta}_1 \rangle + \langle \lambda_2, y_2 \rangle + \langle \overline{\lambda}_2, \overline{\eta}_2 \rangle$$

is strictly convex on its domain  $\Psi(Q) = \{\lambda = (\lambda_1, \overline{\lambda_1}, \lambda_2, \overline{\lambda_2}) | Z(\lambda_1, \overline{\lambda_1}, \lambda_2, \overline{\lambda_2}) < \infty\}$ , and if  $\ell(\lambda_1, \overline{\lambda_1}, \lambda_2, \overline{\lambda_2}) \to \infty$  as  $\lambda = (\lambda_1, \overline{\lambda_1}, \lambda_2, \overline{\lambda_2}) \to \partial \Psi$ , then problem (12) has a unique solution, where " $\partial$ " is the boundary of the set  $\Psi$ . This is always true in the cases we consider here. A simple example in which it does not hold is  $l(\lambda) = e^{\lambda} + \lambda y$  with  $\mathbb{R}$  as domain. This has no minimum for a positive y.

Solving (12) yields  $\lambda_1^*, \overline{\lambda}^*, \lambda_2^*$ , which in turn yields the optimal maximum entropy (posterior) density.

$$\rho^*(\xi_1,\xi_2) = \frac{e^{-\langle \lambda_1^*,B_1\xi_1 \rangle - \langle \overline{\lambda},\overline{B}_1\xi_1 \rangle} e^{-\langle \lambda_2^*,B_2\xi_2 \rangle + \langle \overline{\lambda},\overline{B}_2\xi_2 \rangle}}{Z(\lambda_1^*,\overline{\lambda})Z(\lambda_2^*,-\overline{\lambda})} = \rho_1^*(\xi_1)\rho_2^*(\xi_2)$$

This density is naturally factored into a product of the maximum entropy densities of the two sets of equations. Therefore,  $\xi_1$  and  $\xi_2$  are independent with respect to the reconstructed density  $dP^*(\xi_1,\xi_2) = \rho_1^*(\xi_1)\rho_2^*(\xi_2)dQ_1(\xi_1)dQ_2(\xi_2)$ , and with respect to the original priors. Once  $P^*$  is solved, we follow (8), or (9), to get

$$\begin{pmatrix} \boldsymbol{\beta}_i^* \\ \boldsymbol{\varepsilon}_i^* \end{pmatrix} = E_{P_i^*}[\boldsymbol{\xi}_i^*] = \int_{C_i} \boldsymbol{\xi}_i \boldsymbol{\rho}_i^*(\boldsymbol{\xi}_i) dQ_i(\boldsymbol{\xi}_i); \quad i = 1, 2$$
(13)

With that generic formulation, we show below three analytic examples that cover a wide range of possible priors and support spaces for  $\beta$  and  $\epsilon$ .

### 3. Large Sample Properties

Denote by  $\beta_{iN}^*$  the estimator of the true  $\beta_i$  when the sample size is *N*. Throughout this section we add a subscript *N* to all quantities introduced in Section 2 to remind us that the size of the data set is *N*. We show that  $\beta_{iN}^* \to \beta_i$  and  $\sqrt{N} \left( \beta_{iN}^* - \beta_i \right) \to N(0, V_i)$  as  $N \to \infty$  in some appropriate way. The proof is similar in logic to Proposition 3.2 in [14]. We assume:

Assumption 3.1. For every sample size N, the minimizers of (12) are all in the interior of their domains:  $\lambda_1^* \in int(\mathbb{R}^{|J|}_+)$  and  $\overline{\lambda}_1^* \in int(\mathbb{R}^{N-|J|}_+)$  where "int" stands for interior.

Assumption 3.2. Let  $\frac{1}{N}X_i^tX_i = \frac{1}{N}\left(B_i^tB_i + \overline{B}_i^t\overline{B}_i\right) = \frac{J}{N}\left(\frac{1}{J}B_i^tB_i\right) + \frac{N-J}{N}\left(\frac{1}{N-J}\overline{B}_i^t\overline{B}_i\right)$ . Then, assume there exists  $\alpha \in (0,1)$  such that (i)  $N \to \infty$  and  $J \to \infty$  such that  $\left(\frac{J}{N}\right) \to \alpha$  and (ii) assume that there exists two matrices  $W_i^o$  and  $W_i^u$  such that  $\frac{1}{J}B_i^tB_i \to W_i^o$  and  $\frac{1}{N-J}\overline{B}_i^t\overline{B}_i \to W_i^u$ . Note that  $\frac{1}{N}X_i^tX_i \to W_i = \alpha W_i^o + (1-\alpha)W_i^u$ .

Proposition 3.1. (Convergence in distribution.) Under Assumptions 3.1 and 3.2

a)  $\beta_{iN}^* \xrightarrow{D} \beta_i$  as  $N \to \infty$ , for i=1, 2. b)  $\sqrt{N} \left( \beta_{iN}^* - \beta_i \right) \xrightarrow{D} N(0, V_i)$  as  $N \to \infty$ 

where  $\xrightarrow{D}$  stands for convergence in distribution and  $V_i = \sum_i W_i^{-1} \sum_i$ , where  $\sum_i$  is the covariance matrix of  $\zeta_i$  with respect to  $dQ_i(\zeta_i, V_i)$ .

The approximate finite sample variance is

$$\sigma_i^{*^2} = \frac{1}{N - K_i} \, \boldsymbol{\varepsilon}_i^{*_i} \boldsymbol{\varepsilon}_i^* \text{ for } i = 1, 2 \text{ and } \boldsymbol{\varepsilon}_i^* = E_{P_i^*}[v_i] \text{ as is shown in (8) or similarly in (12)}$$

## 4. Analytic Examples

We discuss three examples, corresponding to assuming that the  $\beta$ 's are either unbounded (Normal), bounded below (Gamma) and bounded below and above (Bernoulli). Under the normal priors, the minimum described in (12) can be explicitly computed. In the other cases, a numerical computation is necessary.

#### 4.1. Normal Priors

Let the constraint space be  $C = C_s \times C_n = \mathbb{R}^K \times \mathbb{R}^N$ . Using the traditional view and centering the support spaces at zero, the prior—a product of two normal distributions—is  $dQ(\xi) = \frac{\exp\left(-\langle \xi, \Sigma^{-2}\xi \rangle/2\right)}{(2\pi)^{(N+K)/2} (\det \Sigma^2)^{1/2}} d\xi$ . The covariance  $\Sigma$  has two diagonal blocks:  $K \times K$  and  $N \times N$ . Without loss of generality, we assume that these two matrices are  $\sigma_s^2 I_K$  and  $\sigma_n^2 I_N$  respectively. Our basic model holds for the general covariance structure  $\Sigma = \begin{bmatrix} \Sigma_1 \\ \Sigma_2 \end{bmatrix}$ .

Formulating these priors within our model yields

$$\ln Z_{1}(\lambda_{1},\overline{\lambda}_{1}) = \frac{1}{2} \Big\{ \Big\langle \lambda_{1}, B_{1} \Sigma_{1}^{2} B_{1}^{t} \lambda_{1} \Big\rangle - 2 \Big\langle \overline{\lambda}, \overline{B}_{1} \Sigma_{1}^{2} B_{1}^{t} \lambda_{1} \Big\rangle + \Big\langle \overline{\lambda}, \overline{B}_{1} \Sigma_{1}^{2} \overline{B}_{1}^{t} \overline{\lambda} \Big\rangle \Big\}$$

$$\ln Z_2(\lambda_2, -\overline{\lambda}_1) = \frac{1}{2} \Big\{ \Big\langle \lambda_2, B_2 \Sigma_2^2 B_2^t \lambda_2 \Big\rangle + 2 \Big\langle \overline{\lambda}, \overline{B}_2 \Sigma_2^2 B_2^t \lambda_2 \Big\rangle + \Big\langle \overline{\lambda}, \overline{B}_2 \Sigma_2^2 \overline{B}_2^t \overline{\lambda} \Big\rangle \Big\}$$

where  $\Sigma_1^2$  is a diagonal (K + N) × (K + N) matrix, the first block being a  $K \times K$  matrix with entries equal to  $\sigma_{1,s}^2$  and the second block is a  $N \times N$  matrix with entries equal to  $\sigma_{1,n}^2$ . That is, the priors on signal and noise spaces are iid normal random variables. Thus, problem (12) consists of finding the minimum of

$$lnZ(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},-\overline{\lambda}_{1})+\langle\lambda_{1}+\lambda_{2},y_{2}\rangle=\langle\lambda_{1}+\lambda_{2},y_{2}\rangle$$
$$+\frac{1}{2}\{\langle\lambda_{1},B_{1}\Sigma_{1}^{2}B_{1}^{t}\lambda_{1}\rangle-2\langle\overline{\lambda},\overline{B}_{1}\Sigma_{1}^{2}B_{1}^{t}\lambda_{1}\rangle+\langle\overline{\lambda},\overline{B}_{1}\Sigma_{1}^{2}\overline{B}_{1}^{t}\overline{\lambda}\rangle\}$$
$$+\frac{1}{2}\{\langle\lambda_{2},B_{2}\Sigma_{2}^{2}B_{2}^{t}\lambda_{2}\rangle+2\langle\overline{\lambda},\overline{B}_{2}\Sigma_{2}^{2}B_{2}^{t}\lambda_{2}\rangle+\langle\overline{\lambda},\overline{B}_{2}\Sigma_{2}^{2}\overline{B}_{2}^{t}\overline{\lambda}\rangle\}$$

over the set described in (12).

To verify that the minimizer of (12) occurs in the interior of the constraint set, we look at the first order conditions

$$B_{1}\Sigma_{1}^{2}B_{1}^{t}\lambda_{1} - B_{1}\Sigma_{1}^{2}\overline{B}_{1}^{t}\lambda + y_{2} = 0$$

$$B_{2}\Sigma_{2}^{2}B_{2}^{t}\lambda_{2} + B_{2}\Sigma_{2}^{2}\overline{B}_{2}^{t}\overline{\lambda} + y_{2} = 0$$

$$-\overline{B}_{1}\Sigma_{1}^{2}\overline{B}_{1}^{t}\overline{\lambda} + \overline{B}_{1}\Sigma_{1}^{2}B_{1}^{t}\lambda_{1} - \overline{B}_{2}\Sigma_{2}^{2}\overline{B}_{2}^{t}\overline{\lambda} - \overline{B}_{2}\Sigma_{2}^{2}B_{2}^{t}\lambda_{2} = 0$$
(14)

A feasible solution to (12) may lie inside the domain of the constraints and provides a solution.

Once the system is solved for  $\lambda_1^*, \overline{\lambda}_1^*, \lambda_2^*$ , the estimated densities are

$$dP_i^*(\xi_i) = \frac{e^{-\|\sum_i^{-1}\xi_i + \sum_i h_i^*\|^2/2}}{(2\pi)^{(N+K_i)/2} (\det \Sigma_i^2)^{1/2}} d\xi_i$$

which, as expected, are normally distributed. Defining

$$h_i = B_i^t \lambda_i + \overline{B}_i^t \overline{\lambda}_i^t = A_i^t \mu_i \equiv \begin{pmatrix} X_i^t \\ I \end{pmatrix} \begin{pmatrix} \lambda_i^* \\ \overline{\lambda}_i^* \end{pmatrix}$$

(recall that  $\overline{\lambda}_2^* = -\overline{\lambda}_1^* = -\overline{\lambda}$  due to the constraints) we use (13) to get

$$\begin{pmatrix} \boldsymbol{\beta}_i^* \\ \boldsymbol{\varepsilon}_i^* \end{pmatrix} = E_{\boldsymbol{P}_i^*}[\boldsymbol{\xi}_i^*] = -\boldsymbol{\Sigma}_i^2 \boldsymbol{A}_i^t \boldsymbol{\mu}_i = -\begin{pmatrix} \boldsymbol{\sigma}_{i,s}^2 & \boldsymbol{0} \\ \boldsymbol{0} & \boldsymbol{\sigma}_{i,n}^2 \end{pmatrix} \boldsymbol{A}_i^t \boldsymbol{\mu}_i = -\begin{pmatrix} \boldsymbol{\sigma}_{i,s}^2 \boldsymbol{X}_i^t \boldsymbol{\mu}_i \\ \boldsymbol{\sigma}_{i,n}^2 \boldsymbol{\mu}_i \end{pmatrix}$$

for i = 1, 2 and where  $\sigma_{i,s}^2$  and  $\sigma_{i,n}^2$  are matrices.

### 4.2. Gamma Priors

Let  $\beta$ 's be bounded below by 0. This can be easily generalized by an appropriate shifting of the support of the distributions. To show the generality of our model, we let the prior on the noise be normal thereby showing that one can use different priors for the signal and the noise.

The signal and noise constraint spaces respectively are  $C_s = \mathbb{R}_+^K = [0,\infty)^K$  and  $C_n = \mathbb{R}^N$ . The prior is

$$dQ(\xi) = \left(\prod_{j=1}^{K} \frac{a^{b} \varsigma_{j}^{b-1} e^{-a\varsigma_{j}}}{\Gamma(b)}\right) d\varsigma_{j} \frac{e^{-\langle v, \Sigma_{n}^{-2} v \rangle/2}}{(2\pi)^{N/2} (\det \Sigma_{n}^{2})} dv$$

Before specifying the concentrated entropy function, we study the matrix A<sub>1</sub> defined as  $A_1 = (X_1 I) = \begin{pmatrix} B_1 \\ \overline{B}_1 \end{pmatrix} = \begin{pmatrix} D_1 I_1 \\ \overline{D}_1 \overline{I}_1 \end{pmatrix}$ . Note that  $\begin{pmatrix} D_1 \\ \overline{D}_1 \end{pmatrix}$  splits  $X_1$  and  $\begin{pmatrix} I_1 \\ \overline{I}_1 \end{pmatrix}$  splits the  $N \times N$  identity matrix to

match the splitting of X. The concentrated entropy function is

$$\ln Z_{1}\left(\lambda_{1},-\overline{\lambda}\right) = \sum_{j=1}^{K} b \ln \left(\frac{a}{a + \left(D_{1}^{t}\lambda_{1} + \overline{D}_{1}^{t}\overline{\lambda}\right)_{j}}\right) + \sigma_{1,n}^{2}\left(\lambda_{1}^{2} + \overline{\lambda}^{2}\right)$$

(Note that  $D_1^t \lambda_1 + \overline{D}_1^t \overline{\lambda} = X_1^t \mu_1$  and  $D_2^t \lambda_2 - \overline{D}_2^t \overline{\lambda} = X_2^t \mu_2$ .) A similar expression exists for  $\ln Z_2(\lambda_2, -\overline{\lambda})$ .

The problem (12) consists of minimizing

$$lnZ(\lambda_{1},-\lambda,\lambda_{2},\lambda) + \langle \lambda_{1} + \lambda_{2}, y_{2} \rangle$$
  
=  $\sum_{j=1}^{K_{1}} b \ln \left( \frac{a}{a + (D_{1}^{t}\lambda_{1} + \overline{D}_{1}^{t}\overline{\lambda})_{j}} \right) + \sum_{j=1}^{K_{2}} b \ln \left( \frac{a}{a + (D_{2}^{t}\lambda_{2} - \overline{D}_{2}^{t}\overline{\lambda})_{j}} \right)$   
+  $\sigma_{1,n}^{2} (\lambda_{1}^{2} + \overline{\lambda}^{2}) + \sigma_{2,n}^{2} (\lambda_{2}^{2} + \overline{\lambda}^{2}) + \langle \lambda_{1} + \lambda_{2}, y_{2} \rangle$ 

Once  $\lambda_1^*, \overline{\lambda}^*, \lambda_2^*$  are found, the optimal density is

$$dP_{i}^{*}(\xi_{i}) = \left(\prod_{j=1}^{K_{i}} \frac{\left(a + \left(X_{i}^{t} \mu_{i}\right)_{j}\right)^{b} \varsigma_{j}^{b-1} e^{-\left(\left(X_{i}^{t} \mu_{i}\right)_{j}+a\right) \varsigma_{j}}}{\Gamma(b)} d\varsigma_{j}\right) \frac{e^{-\left\|v + \sigma_{i,n}^{2} \mu_{i}^{*}\right\|^{2} / 2\sigma_{i,n}^{2}}}{\left(2\pi\sigma_{i,n}^{2}\right)^{N_{i}/2}} dv$$

The estimated parameters are

$$(\beta_i^*)_j = E_{P_i^*}[(\varsigma_i)_j] = \frac{b}{(a + (X_i^t \mu_i)_j)}; \text{ for } j = 1, ..., K_i \text{ and } i = 1, 2$$

The realized residuals are

$$(\varepsilon_{i}^{*})_{l} = E_{P_{i}^{*}}[(v_{i})_{l}] = -\sigma_{i,n}^{2}(\mu_{i}^{*})_{l}; \text{ for } l = 1,...,N_{i} \text{ and } i = 1,2$$

## 4.3. Bernoulli Priors

This example represents another extreme case where it is assumed that the  $\beta$ 's are bounded. For simplicity, assume that we know that all  $\beta$ 's lie in the interval [a, b], which makes  $C_s = [a,b]^{K_i}$  the choice for all of the constraints on the signal space. For the noise component, we follow the previous formulation of normal priors.

With this background, the prior measure used is

$$dQ(\xi) = \left(\prod_{j=1}^{K} \frac{1}{2} \left( \delta_a \left( d\varsigma_j \right) + \delta_b \left( d\varsigma_j \right) \right) \right) \frac{e^{-\langle v, \Sigma_n^{-2} v \rangle/2}}{\left( 2\pi \right)^{N/2} \left( \det \Sigma_n^2 \right)} dv$$

The concentrated entropies are

$$\ln Z_{i}\left(\lambda_{i},\overline{\lambda}\right)_{i} = \sum_{j=1}^{K_{i}} \ln \frac{1}{2} \left(e^{-g_{i,j}a} + e^{-g_{i,j}b}\right) + \sigma_{n}^{2} \left\|\mu_{i}\right\|^{2}, i = 1, 2$$

where  $g_i = D_i^t \lambda_i + \overline{D}_i^t \overline{\lambda}_i$  and  $\mu_i = (\lambda_i \overline{\lambda}_i)^t$ . Recall that  $\overline{\lambda}_2 = -\overline{\lambda}_1$ . In this case, the function to be minimized is

$$lnZ\left(\lambda_{1},\overline{\lambda},\lambda_{2},-\overline{\lambda}\right)+\left\langle\lambda_{1}+\lambda_{2},y_{2}\right\rangle$$
$$=\sum_{j=1}^{K_{1}}\ln\frac{1}{2}\left(e^{-g_{1,j}a}+e^{-g_{1,j}b}\right)+\sum_{j=1}^{K_{2}}\ln\frac{1}{2}\left(e^{-g_{2,j}a}+e^{-g_{2,j}b}\right)$$
$$+\sigma_{1,n}^{2}\left(\lambda_{1}^{2}+\overline{\lambda}^{2}\right)+\sigma_{2,n}^{2}\left(\lambda_{2}^{2}+\overline{\lambda}^{2}\right)+\left\langle\lambda_{1}+\lambda_{2},y_{2}\right\rangle$$

which is minimized over the region described in (12). Again, the optimal solutions (minimizing  $\lambda_1^*, \overline{\lambda}^*, \lambda_2^*$ ) is to be found numerically. The estimated post-data is

$$dP_{i}^{*}(\xi) = \left(\prod_{j=1}^{K} \left(p_{i,j}\delta_{a}(d\varsigma_{j}) + q_{i,j}\delta_{b}(d\varsigma_{j})\right)\right) \frac{e^{-\|v + \sigma_{i,n}^{2}\mu_{i}^{*}\|^{2}/2\sigma_{i,n}^{2}}}{\left(2\pi\sigma_{i,n}^{2}\right)^{N_{i}/2}}dv$$

for i = 1,2 and where

$$p_{i,j} = \frac{e^{-ag_{i,j}}}{e^{-ag_{i,j}} + e^{-bg_{i,j}}} = 1 - q_{i,j}$$

from which the estimated parameters and residuals are given by

$$\beta_{i,j}^{*} = \frac{ae^{-ag_{i,j}} + be^{-bg_{i,j}}}{e^{-ag_{i,j}} + e^{-bg_{i,j}}} \quad \text{and} \quad \varepsilon_{i}^{*} = -\sigma_{i,n}^{2}\mu_{i}^{*}$$

### **5.** Empirical Example

We illustrate the applicability of our approach using an empirical application consisting of a small data set. The objective here is to demonstrate that our IT estimator is easy to apply and can be used for many different priors. The small sample performance of the IT-GME version of that estimator (uniform discrete priors) and detailed comparisons with other competing estimators is already shown in GMP and it falls outside the objectives of this note. The empirical example is based on one of the examples analyzed in GMP with data drawn from the March 1996 Current Population Survey. We estimated the wage-participation model for the subset of respondents in the labor market. Workers who are self-employed are excluded from the sample. Since the normal maximum likelihood estimator did not converge for that data [15], only results for the OLS, Heckman two-step, a semi-parametric estimator with a nonparametric selection mechanism due to [5], AP, and the different IT models developed here are reported [16]. To make our results comparable across the IT estimators, we use the empirical

standard deviations in all three cases and use supports between -100 and 100 for the IT-GME (uniform discrete priors) and the IT-Bernoulli case. In both the IT-Normal and IT-Bernoulli the priors used for the noise components are normal (as is shown in Section 4). Under these very similar specifications, we would expect all three IT examples to yield comparable estimates. Naturally, there are many other priors to choose from, but the objective here is just to show the flexibility and applicability of our approach.

We analyze a sample of 151 Native American females, of whom 65 are in the labor force. The wage equation covariates include years of education, a dummy for currently enrolled in school, potential experience (age - education - 6) and potential experience squared, a dummy for rural location, and a dummy for central city location. The covariates in the selection equation include all the variables in the wage equation and the amount of welfare payments received in the previous year, a dummy equal one for married, and the number of children. We use the three exclusion restrictions to identify the wage equation in the parametric and nonparametric two-step approaches.

|                    | OLS    | 2-Step  | AP     | IT-GME | IT-Normal | IT-Bernoulli |
|--------------------|--------|---------|--------|--------|-----------|--------------|
| Constant           | 1.073  | 1.771   | NA     | 1.038  | 1.049     | 1.068        |
| Education          | 0.055  | 0.043   | 0.044  | 0.054  | 0.056     | 0.055        |
| Experience         | 0.038  | 0.023   | 0.038  | 0.038  | 0.038     | 0.038        |
| Experience Squared | -0.001 | -0.0005 | -0.001 | -0.001 | -0.001    | -0.001       |
| Rural              | 0.214  | 0.268   | 0.332  | 0.210  | 0.215     | 0.214        |
| Central City       | -0.170 | -0.091  | -0.171 | -0.186 | -0.166    | -0.169       |
| Enrolled in School | -0.290 | -0.471  | -0.190 | -0.301 | -0.283    | -0.288       |
| λ                  |        | -0.461  |        |        |           |              |
| $\mathbf{R}^2$     | 0.355  | 0.376   | NA     | 0.343  | 0.355     | 0.354        |
| MSPE               | 0.157  | 0.135   | NA     | 0.147  | 0.144     | 0.144        |

Table 1. Estimates of the Native American wage equation (151 individuals; 65 in labor force).

Notes: Bold numbers reflect significantly different than zero at the 10% level

Table 1 presents the estimated coefficients for the wage equation. The  $R^2$  and Mean Squared Prediction Error (MSPE) for each model are presented as well. All IT estimators outperform the other estimators in terms of predicting selection [17]. The estimated return to education is about 5% across all estimation methods, but only statistically significantly different from 0 for the OLS and the IT estimators. Though, all estimators have estimated parameters of the same magnitude and sign, only the OLS and the three reported IT estimates are statistically significantly different from zero in most cases.

# 6. Conclusion

In this short paper we develop a simple to apply, information-theoretic, method for analyzing nonlinear data with sample selection problem. Rather than using a likelihood approach or a semi-parametric approach we generalized further the IT-GME model of Golan, Moretti and Perloff (2004). Our model (i) allows for bounded and unbounded supports on all the unknown parameters, (ii) allows us to use a whole class of priors (continuous or discrete), (iii) is specified as a nonlinear concentrated entropy model, and (iv) is easy to apply. Like GMP our model works well even with

small data. This is shown in our empirical example. The extensions developed here mark a significant improvement on the GMP model and other IT, generalized entropy models.

A detailed set of sampling experiments comparing our IT method with all other competitors, under different data processes, will be done in future work.

# Acknowledgement

We thank Enrico Moretti and Jeff Perloff for their comments on earlier versions of this work.

# **References and Notes**

- 1. Heckman, J. Sample selection bias as a specification error. *Econometrica* **1979**, 47,153-161.
- 2. Manski, C.F. Semiparametric analysis of discrete response: Asymptotic properties of the maximum score estimator. *J. Econom.* **1985**, *27*, 313-333.
- 3. Cosslett, S.R. Distribution-free maximum likelihood estimator of the binary choice model. *Econometrica* **1981**, *51*, 765-782.
- 4. Han, A.K. Non-parametric analysis of a generalized regression model: The maximum rank correlation estimator. *J. Econom.* **1987**, *35*, 303-316.
- 5. Ahn, H.; Powell, J.L. Semiparametric estimation of censored selection models with a nonparametric selection mechanism. *J. Econom.* **1993**, *58*, 3-29.
- 6. Golan, A.; Moretti E.; Perloff, J.M. A small sample estimation of the sample selection model. *Econom. Rev.* 2004, 23, 71-91.
- 7. Golan, A.; Judge, G.; Miller, D. *Maximum Entropy Econometrics: Robust Estimation with Limited Data*; John Wiley & Sons: New York, NY, USA, 1996.
- 8. Golan, A.; Judge, G.; Perloff, J.M. Recovering information from censored and ordered multinomial response data. *J. Econom.* **1997**, *79*, 23-51.
- 9. Golan A. An information theoretic approach for estimating nonlinear dynamic models. *Nonlinear Dynamics Econom.* **2003**, *7*, 2.
- 10. Maddala, G.S. *Limited-Dependent and Qualitative Variables in Econometrics*; Cambridge University Press: Cambridge, UK, 1983.
- 11. Kullback, S. Information Theory and Statistics; John Wiley & Sons: New York, NY, USA, 1959.
- 12. Kullback, S.; Leibler, R.A. On information and sufficiency. Ann. Math. Statist. 1951, 22, 79-86.
- 13. Gokhale, D.V.; Kullback, S. *The Information in Contingency Tables*; Marcel Dekker: New York, NY, USA, 1978.
- 14. Golan, A.; Gzyl, H. An information theoretic estimator for the linear model. Working paper, 2008.
- 15. Due to the small size of the data and the large proportion of censored observations, none of the maximum likelihood methods would converge with all standard software.
- 16. For discussion of the data, detailed analyses and discussion of the different estimators as well as a detailed discussion of the AP [5] application, see GMP [6]. The 2-step and AP estimates are taken from that paper, Table 8.
- 17. To keep the Table simple and since these specific results are not of interested here, they are not presented.

# Appendix

# **Proof of Theorem 2.1.**

Proof: Recall (11)

$$Sup\{Sup\{S(P_{1},Q_{1}) + S(P_{2},Q_{2}) | y_{2} = B_{2}E_{P_{2}}[\xi_{2}], \eta_{1} = B_{1}E_{P_{1}}[\xi_{1}]; \\ \overline{\eta}_{2} = \overline{B}_{2}E_{P_{2}}[\xi_{2}]; \overline{\eta}_{1} = \overline{B}_{1}E_{P_{1}}[\xi_{1}]\} | y_{2} > \eta_{1}, \overline{\eta}_{2} \le \overline{\eta}_{1}\}$$
(A.1)

First, note that the inner problem in (A.1) is equivalent to

$$\ell\left(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},\overline{\lambda}_{2}\right) = \ln Z\left(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},\overline{\lambda}_{2}\right) + \left\langle\lambda_{1},\eta_{1}\right\rangle + \left\langle\overline{\lambda}_{1},\overline{\eta}_{1}\right\rangle + \left\langle\lambda_{2},y_{2}\right\rangle + \left\langle\overline{\lambda}_{2},\overline{\eta}_{2}\right\rangle$$

where  $\lambda_i$  and  $\overline{\lambda_i}$  (*i*=1, 2) are the Lagrange multipliers associated with the data (9) and Z is the normalization factor of  $dP_1 dP_2$ :

$$Z\left(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},\overline{\lambda}_{2}\right) = \iint_{C_{1}\times C_{2}} e^{-\langle\lambda_{1},B_{1}\xi_{1}\rangle} e^{-\langle\overline{\lambda}_{1},\overline{B}_{1}\xi_{1}\rangle} e^{-\langle\overline{\lambda}_{2},B_{2}\xi_{2}\rangle} e^{-\langle\overline{\lambda}_{2},\overline{B}_{2}\xi_{2}\rangle} dQ_{1}\left(\xi_{1}\right) dQ_{2}\left(\xi_{2}\right)$$
$$= \int_{C_{1}} e^{-\langle\lambda_{1},B_{1}\xi_{1}\rangle} e^{-\langle\overline{\lambda}_{1},\overline{B}_{1}\xi_{1}\rangle} dQ_{1}\left(\xi_{1}\right) \int_{C_{2}} e^{-\langle\lambda_{2},B_{2}\xi_{2}\rangle} e^{-\langle\overline{\lambda}_{2},\overline{B}_{2}\xi_{2}\rangle} dQ_{2}\left(\xi_{2}\right) = Z_{1}\left(\lambda_{1},\overline{\lambda}_{1}\right) Z_{2}\left(\lambda_{2},\overline{\lambda}_{2}\right)$$

Note that the inner *sup* is over  $(P_1, P_2)$ , and the outer *sup* is over the  $\eta$ 's in the region indicated within the  $\{\cdot\}$ . The basic idea here is to replace the inequalities appearing in problem (10) with equalities. Next, the dual-unconstrained model of this inner primal problem is the solution to

$$\inf_{\boldsymbol{\lambda}} \left\{ \ell \left( \lambda_{i}, \overline{\lambda_{i}} \right) \mid \lambda_{i} \in \mathbb{R}^{|J|}, \text{ and } \overline{\lambda_{i}} \in \mathbb{R}^{N-|J|} \text{ for } i = 1, 2 \right\}$$

where |J| is the number of observations where  $y_{2i} > y_{1i}$ . With this step, the equivalent dual model of the primal problem (A.1) is

$$Sup\{Inf\{lnZ(\lambda_{1},\overline{\lambda}_{1},\lambda_{2},\overline{\lambda}_{2})+\langle\lambda_{1},\eta_{1}\rangle+\langle\overline{\lambda}_{1},\overline{\eta}_{1}\rangle+\langle\lambda_{2},y_{2}\rangle+\langle\overline{\lambda}_{2},\overline{\eta}_{2}\rangle| \\ \lambda_{i} \in \mathbb{R}^{|J|}, \text{ and } \overline{\lambda_{i}} \in \mathbb{R}^{|N-|J|}\}|y_{2} > \eta_{1},\overline{\eta}_{2} \leq \overline{\eta}_{1} \quad for \ i=1,2\}$$
(A.2)

Next, we rewrite the constraints for the outer problem. The constraint  $y_2 > \eta_1$  is rewritten as  $\eta_1 = y_2 - \zeta$ ,  $\zeta > 0$ , and the constraint  $\overline{\eta}_2 \leq \overline{\eta}_1$  is written as  $\overline{\eta}_2 \in \mathbb{R}^{N-|J|}$ ,  $\overline{\eta}_1 = \overline{\eta}_2 + \overline{\zeta}$ ,  $\overline{\zeta} \geq 0$ . Model (A.2) becomes

$$Sup\{Inf\{InZ(\lambda_{1},\overline{\lambda_{1}},\lambda_{2},\overline{\lambda_{2}})+\langle\lambda_{1}+\lambda_{2},y_{2}\rangle-\langle\lambda_{1},\zeta\rangle+\langle\overline{\lambda_{1}}+\overline{\lambda_{2}},\overline{\eta_{1}}\rangle+\langle\overline{\lambda_{2}},\overline{\zeta}\rangle|\\\lambda_{i}\in\mathbb{R}^{|J|}, \text{ and } \overline{\lambda_{i}}\in\mathbb{R}^{N-|J|}\}|\zeta>0,\overline{\eta_{1}},\overline{\zeta}\geq0\}$$

Next, exchanging the sup and the inf operations we get

$$Inf \{ lnZ(\lambda_{1}, \overline{\lambda_{1}}, \lambda_{2}, \overline{\lambda_{2}}) + \langle \lambda_{1} + \lambda_{2}, y_{2} \rangle + Sup \{ -\langle \lambda_{1}, \zeta \rangle + \langle \overline{\lambda_{1}} + \overline{\lambda_{2}}, \overline{\eta_{1}} \rangle + \langle \overline{\lambda_{2}}, \overline{\zeta} \rangle | \\ \zeta > 0, \overline{\eta_{1}}, \overline{\zeta} \ge 0 \} | \lambda_{i} \in \mathbb{R}^{|J|}, \text{ and } \overline{\lambda_{i}} \in \mathbb{R}^{N-|J|} \}$$

To compute the inner supremum, note that

$$Sup\{\langle \lambda, \zeta \rangle | \zeta > 0\} = \begin{cases} 0 & \text{if } -\lambda \in \mathbb{R}^{d}_{+} \\ \infty & \text{otherwise} \end{cases}$$

$$Sup\{\langle \lambda, \zeta \rangle | \overline{\zeta} \ge 0\} = \begin{cases} 0 & \text{if } -\lambda \in \mathbb{R}^{d}_{+} \\ \infty & \text{otherwise} \end{cases}$$

where  $\mathbb{R}^{d}_{+}$  denotes the non-negative orthant in  $\mathbb{R}^{d}$ , the right hand side of the last identity is usually written as  $I_{\mathbb{R}^{d}_{+}}(-\lambda)$  and  $I_{\mathbb{A}}(x)$  is defined as  $I_{\mathbb{A}}(x) = \begin{cases} 0 & \text{if } x \in A \\ \infty & \text{if } x \notin A \end{cases}$ . The difference between the first and second problem is that in the first the supremum is reached only when  $\lambda = 0$ , whereas in the second it is reached at the boundary of  $\mathbb{R}^{d}_{+}$ . Similarly,

$$\operatorname{Max}\left\{\left\langle \lambda,\zeta\right\rangle \mid \zeta\in\mathbb{R}^{d}\right\}=I_{\left\{0\right\}}\left(\lambda\right)$$

Noting that  $\overline{\lambda}_1 = -\overline{\lambda}_2 \equiv \overline{\lambda}$ , our MinMax problem (A.2) reduces to finding

$$\operatorname{Inf}\left\{\ln Z\left(\lambda_{1},-\overline{\lambda},\lambda_{2},\overline{\lambda}\right)+\left\langle\lambda_{1}+\lambda_{2},y_{2}\right\rangle+I_{\mathbb{R}^{d}_{+}}\left(\lambda_{1}\right)+I_{\mathbb{R}^{d}_{+}}\left(\overline{\lambda}_{1}\right)\mid\lambda_{1}\in\mathbb{R}^{|J|},\lambda_{2}\in\mathbb{R}^{|J|}\text{ and }\overline{\lambda}\in\mathbb{R}^{|J-J|}\right\}$$

which simplified to

$$Inf\left\{lnZ\left(\lambda_{1},-\overline{\lambda},\lambda_{2},\overline{\lambda}\right)+\left\langle\lambda_{1}+\lambda_{2},y_{2}\right\rangle\mid\lambda_{1}\in\mathbb{R}^{|J|}_{+},\lambda_{2}\in\mathbb{R}^{|J|}\text{ and }\overline{\lambda_{1}}\in\mathbb{R}^{|J-|J|}_{+}\right\}$$

which is (12).

© 2010 by the authors; licensee MDPI, Basel, Switzerland. This article is an Open Access article distributed under the terms and conditions of the Creative Commons Attribution license (http://creativecommons.org/licenses/by/3.0/).