

# **A Small-Sample Estimator for the Sample-Selection Model**

by

Amos Golan, Enrico Moretti, and Jeffrey M. Perloff

August 2002

## **ABSTRACT**

A semiparametric estimator for evaluating the parameters of data generated under a sample selection process is developed. This estimator is based on the generalized maximum entropy estimator and performs well for small and ill-posed samples. Theoretical and sampling comparisons with parametric and semiparametric estimators are given. This method and standard ones are applied to three small-sample empirical applications of the wage-participation model for female teenage heads of households, immigrants, and Native Americans.

Key Words: Maximum Entropy, Sample Selection, Monte Carlo Experiments.

Contact:

Amos Golan

Department of Economics

American University

Roper 200

4400 Massachusetts Ave., N.W.

Washington, D.C. 20016-8029

agolan@american.edu

(202)885-3783

We are grateful to Christopher Mckelvey for outstanding research assistance. We thank James Powell for very helpful suggestions. Two anonymous referees made helpful, constructive comments. Moretti and Perloff thank the Institute for Industrial Relations at Berkeley and the Giannini Foundation for support.

## Table of Contents

1. Introduction	1
2. The Model	2
3. Estimation Approach	3
3.1 Review of Maximum Entropy and Generalized Maximum Entropy Estimators	3
3.2 A GME Sample-Selection Estimator	7
4. Sampling Experiments	11
4.1 Experimental Designs	11
4.2 Alternative Estimators	13
4.3 Discussion of Results	16
5. Empirical Applications	19
5.1 Teenage Heads of Households	20
5.2 Recent Immigrants	21
5.3 Native Americans	22
5.4 Summary	22
6. Conclusion	23
References	24

## **A Small-Sample Estimator for the Sample-Selection Model**

### 1. INTRODUCTION

The problem of sample selection arises frequently in econometric studies of individuals' wages or labor supply and other topics. When sample sizes are small, existing parametric (full and limited information maximum likelihood, Heckman 1976, 1979) and semiparametric estimators (Manski, 1975, 1985; Cosslett, 1981; Han, 1987; Ahn and Powell, 1993) have difficulties.

We have three objectives. First, we develop a semiparametric estimator for the sample-selection problem that performs well when the sample is small. This estimator has its roots in information theory and is based on the generalized maximum entropy (GME) approach of Golan, Judge, and Miller (1996) and Golan, Judge, and Perloff (1997). Second, we use Monte Carlo experiments to compare and contrast the small-sample behavior of our GME estimator with other parametric and semiparametric estimators. Third, we apply this method to examine the wage-participation of several groups of females where the data sets are small.

Section 2 discusses the sample-selection model. Section 3 reviews GME and develops a GME sample-selection estimator with the relevant inferential statistics. Section 4 lays out the experimental design and discusses the sampling results. Section 5 applies the various methods to the wage-participation model for female teenage heads of households, immigrants, and Native Americans, all of which involve relatively small samples. Section 6 briefly summarizes the results.

## 2. THE MODEL

Many sample-selection models exist. For specificity, we consider a common labor model (see, e.g., Maddala, 1983). Suppose the  $i$ th person values staying home (working in the home) at  $y_{1i}^*$  and can earn  $y_{2i}^*$  in the marketplace. If  $y_{2i}^* > y_{1i}^*$ , the individual chooses to work in the marketplace,  $y_{1i} = 1$ , and we observe the market value,  $y_{2i} = y_{2i}^*$ . Otherwise,  $y_{1i} = 0$  and  $y_{2i} = 0$ .

The individual's value at home and in the marketplace depends on education, experience, and other demographic characteristics:

$$y_{1i}^* = \underline{x}'_1 \underline{\beta}_1 + e_{1i} \quad (2.1)$$

$$y_{2i}^* = \underline{x}'_2 \underline{\beta}_2 + e_{2i}, \quad (2.2)$$

where  $\underline{x}_1 = (1, x_{12}, \dots, x_{1L})'$ ,  $\underline{x}_2 = (1, x_{22}, \dots, x_{2K})'$ ,  $\underline{\beta}_1$  and  $\underline{\beta}_2$  are  $L$  and  $K$ -dimensional vectors of unknowns. We observe

$$y_{1i} = \begin{cases} 1 & \text{if } y_{2i}^* > y_{1i}^* \\ 0 & \text{if } y_{2i}^* \leq y_{1i}^*. \end{cases} \quad (2.3)$$

$$y_{2i} = \begin{cases} \underline{x}'_2 \underline{\beta}_2 + e_{2i} & \text{if } y_{2i}^* > y_{1i}^* \\ 0 & \text{if } y_{2i}^* \leq y_{1i}^*. \end{cases} \quad (2.4)$$

Our objective is to estimate  $\underline{\beta}_1$  and  $\underline{\beta}_2$ . Typically in these types of studies, the researcher is interested primarily in  $\underline{\beta}_2$ .

### 3. ESTIMATION APPROACH

We use a GME approach to estimate the sample-selection model. We start by providing some background as to how the generalized maximum entropy approach works, and then develop the GME sample-selection estimator.

#### 3.1 Review of Maximum Entropy and Generalized Maximum Entropy Estimators

The GME estimator is based on the classic maximum entropy (ME) approach of Jaynes (1957a, 1957b, 1984), which uses the entropy-information measure of Shannon (1948) to recover the unknown probability distribution of underdetermined problems. Shannon's (1948) entropy measure reflects the uncertainty (state of knowledge) we have about the occurrence of a collection of events. Letting  $x$  be a random variable with possible outcomes  $x_s$ ,  $s = 1, 2, \dots, n$ , with probabilities  $p_s$  such that  $\sum_s p_s = 1$ , Shannon (1948) developed the information criterion, *entropy*, as

$$H \equiv -\sum_s p_s \ln p_s, \quad (3.1)$$

with  $x \ln(x) \rightarrow 0$  as  $x \rightarrow 0$ . The function  $H$ , is a unique function of the distribution  $\underline{p} = (p_1, p_2, \dots, p_n)'$  and measures the amount of information in  $\underline{x}$ . It reaches a maximum of  $\ln(n)$  when  $p_1 = p_2 = \dots = p_n = 1/n$ , and it is zero when  $p_s = 1$  for one value of  $s$ . To recover the unknown probabilities  $\underline{p}$  that characterize a given data set, Jaynes (1957a, 1957b) proposed maximizing entropy, subject to available sample-moment information and adding up constraints on the probabilities. For an axiomatic development of ME and a review of its proper-

ties, see Kullback (1959), Levine (1980), Shore and Johnson (1980), Skilling (1989), and Csiszár (1991).

Unlike the ME method, where the observed moments are assumed to be exact, under the GME approach, the observed sample moments, or similarly each data point, are assumed to be noisy. As such, the GME is a generalization of the ME approach and uses a dual-objective (precision and prediction) function. To illustrate the GME approach, we examine the linear model

$$\underline{y} = X\underline{\beta} + \underline{e}, \quad (3.2)$$

where  $\underline{\beta}$  is a  $K$ -dimensional vector of unobserved parameters,  $\underline{y} = (y_1, y_2, \dots, y_T)'$  is a  $T$ -dimensional vector of observed data, and  $X$  is a  $(T \times K)$  design matrix. Our objective is to recover the unknown vector  $\underline{\beta}$  using as few assumption as possible. Consistent with our goal, we impose no distributional assumptions and no assumptions regarding the exact relationship between sample and population moments. That is, our objective is to simultaneously recover the signal  $\underline{\beta}$  and the noise (unknown error distribution)  $\underline{e}$  where both are unknown.

In order to relax the distributional assumptions, we need to provide some bounds on the solutions. In this way, the GME approach (developed by Golan, Judge, and Miller, 1996) is related to regularization methods. As a first step to setting these bounds, the parameters of Equation 3.2 are reparameterized as

$$\underline{y} = X\underline{\beta} + \underline{e} = XZ\underline{p} + V\underline{w}. \quad (3.3)$$

In this reformulation, the coefficient  $\beta_k$  on the  $k^{\text{th}}$  variable in  $X$  is defined as

$$\beta_k \equiv \sum_m z_{km} p_{km} \quad (3.4)$$

where

$$\sum_m p_{km} = 1 \quad (3.5)$$

and  $\underline{z}_k = (z_{k1}, z_{k2}, \dots, z_{kM})'$  is an  $M \geq 2$  dimensional vector. This vector serves as a discrete support space for each one of the  $K$  unknowns and is specified to span the possible (unknown) values of  $\beta_k$ . This reformulation converts the unknown elements of the vector  $\underline{\beta}$  from the real line to the  $[0, 1]$  interval with the properties of  $K$  proper probability distributions  $\underline{p}_k$  defined over the  $K$  support spaces  $\underline{z}_k$ .

How should we specify  $Z$ ? If we possess no knowledge as to the possible values of  $\beta_k$ , we specify  $\underline{z}_k$  to be symmetric around zero, with large negative and positive boundaries. For example,  $z_{k1} = -z_{kM} = -10^6$ . Often, however, we have some knowledge regarding the possible values of  $\underline{\beta}$  and use that information to specify  $Z$ .<sup>1</sup>

Similarly, we transform each  $e_t$  into  $T$  proper probability distributions. Define a discrete support space  $\underline{v}$  of dimension  $J \geq 2$ , *equally spaced and symmetric around zero* and associate with it a set of weights  $w_t$  such that

$$e_t \equiv \sum_j v_j w_{tj} \quad (3.6)$$

and

$$\sum_j w_{tj} = 1, \quad (3.7)$$

and  $V$  is a  $T \times J$  matrix of the  $T$  identical vectors  $v$ . The end points,  $v_1$  and  $v_J$ , are chosen to be  $-3\sigma_y$  and  $3\sigma_y$  where  $\sigma_y$  is the empirical standard deviation of  $\underline{v}$ .

---

<sup>1</sup> We can specify  $\underline{z}_k$  to be continuous and without bounds (see Golan and Gzyl, 2002), but here we use a discrete support space.

Having converted the two sets of unknowns into probability distributions, the estimation problem is to maximize a dual-loss function where emphasis is placed on both prediction and precision (smoothness) of the estimates:

$$\max_{\underline{p}, \underline{w}} H(\underline{p}, \underline{w}) = -\sum_k \sum_m p_{km} \ln p_{km} - \sum_t \sum_j w_{tj} \ln w_{tj} \quad (3.8)$$

subject to

$$y_t = \sum_k \sum_m x_{tk} z_{km} p_{km} + \sum_j v_j w_{tj} \quad (3.9)$$

$$\sum_m p_{km} = 1 \quad (3.10)$$

$$\sum_j w_{tj} = 1. \quad (3.11)$$

This estimator shrinks *all* unknown parameters to the center of the support given the data.

The ME estimator is a special case of the GME, in which no weight is placed on the noise component and the  $T$  observations (3.4) are represented as  $K$  zero moments.

Letting  $\hat{\lambda}_t$  be the estimate of the Lagrange multiplier associate with constraints (3.9), the optimal solution is

$$\hat{p}_{km} = \frac{\exp\left(-\sum_t \hat{\lambda}_t x_{tk} z_{km}\right)}{\sum_m \exp\left(-\sum_t \hat{\lambda}_t x_{tk} z_{km}\right)} \equiv \frac{\exp\left(-\sum_t \hat{\lambda}_t x_{tk} z_{km}\right)}{\Omega_k(\hat{\lambda})} \quad (3.12)$$

and



$$\hat{w}_{tj} = \frac{\exp(-\hat{\lambda}_t v_j)}{\sum_j \exp(-\hat{\lambda}_t v_j)} \equiv \frac{\exp(-\hat{\lambda}_t v_j)}{\Psi_t(\hat{\lambda})}. \quad (3.13)$$

The resulting point estimates are  $\hat{\beta} \equiv \sum_m z_{km} \hat{p}_{km}$  and  $\hat{\underline{\lambda}}_t \equiv \sum_j v_j \hat{w}_{tj}$ .

Golan, Judge, and Miller (1996) show that a dual, concentrated model can be constructed. Substituting into the first element of the right-hand side of the Lagrangean corresponding to Equations 3.8 - 3.9, the post-data  $\hat{p}$  and  $\hat{w}$ , yields

$$L(\underline{\lambda}) = \sum_t y_t \lambda_t + \sum_k \ln \Omega_k(\lambda) + \sum_t \ln \Psi_t(\lambda). \quad (3.14)$$

Setting the derivative of  $L(\underline{\lambda})$  with respect to  $\underline{\lambda}$  equal to zero yields  $\hat{\underline{\lambda}}$ , from which we can derive  $\hat{\underline{\beta}}$  and  $\hat{\underline{\lambda}}$ .

### 3.2 A GME Sample-Selection Estimator

We now apply the same approach to the sample-selection problem. We start by reparameterizing the signals  $\underline{\beta}_1$  and  $\underline{\beta}_2$  to be proper probability distributions that are defined over some support. We start by choosing a support space with  $M \geq 2$  of discrete points  $\underline{z}_{11} = [z_{111}, z_{112}, \dots, z_{11M}]'$  and  $\underline{z}_{2k} = [z_{2k1}, z_{2k2}, \dots, z_{2kM}]'$  that span the possible range of the unknowns  $\underline{\beta}_1$  and  $\underline{\beta}_2$ .

Then we let

$$\beta_{1l} = \sum_m p_{1lm} z_{1lm}, \quad l = 1, \dots, L, \quad (3.15a)$$

and

$$\beta_{2k} = \sum_m p_{2km} z_{2km}, \quad k = 1, \dots, K, \quad (3.15b)$$

where  $\underline{p}_{1l}$  and  $\underline{p}_{2k}$  are proper probability vectors that correspond to the  $M$ -dimensional support vectors of weights. As the dimension of  $M$  increases, we recover more moments for each point estimate  $\underline{\beta}_1$  and  $\underline{\beta}_2$ . However, for all practical purposes, the results are not sensitive to  $M$  as long as  $M$  is at least of dimension 3.

Similarly, we treat the errors as unknowns and use the following parametrization. Let each  $\underline{e}_1$  and  $\underline{e}_2$  be specified as

$$e_{1i} = \sum_j w_{1ij} v_j \quad \text{and} \quad e_{2i} = \sum_j w_{2ij} v_j \quad (3.16)$$

where  $\underline{w}_1$  and  $\underline{w}_2$  are proper probability vectors and  $\underline{v}$  is a support space of dimension greater than or equal to two that is symmetric about zero. The lower and upper bounds of the support space  $\underline{v}$  are  $-3\sigma_{y_2}$  and  $3\sigma_{y_2}$  respectively where  $\sigma_{y_2}$  is the empirical standard deviation of the observed non-zero left-hand side variables.

Having parameterized the unknowns, we now wish to maximize the dual loss (objective) function, which is the sum of the joint entropies of the signal and noise in the system, subject to the sample-selection model, Equations 2.3 and 2.4:

$$\begin{aligned} \max_{\underline{p}_1, \underline{p}_2, \underline{w}_1, \underline{w}_2} H(\underline{p}_1, \underline{p}_2, \underline{w}_1, \underline{w}_2) = & -\sum_l \sum_m p_{1lm} \ln p_{1lm} - \sum_k \sum_m p_{2km} \ln p_{2km} \\ & - \sum_i \sum_j w_{1ij} \ln w_{1ij} - \sum_i \sum_j w_{2ij} \ln w_{2ij} \end{aligned} \quad (3.17)$$

subject to

$$\sum_k \sum_m x_{2ik} z_{2km} p_{2km} + \sum_j v_j w_{2ij} = y_{2i}, \quad \text{if } y_{2i} > 0, \quad (3.18)$$

$$\sum_k \sum_m x_{2ik} z_{2km} p_{2km} + \sum_j v_j w_{2ij} > \sum_l \sum_m x_{1il} z_{1lm} p_{1lm} + \sum_j v_j w_{1ij}, \quad (3.19)$$

if  $y_{2i}^* > 0$ ,

$$\sum_k \sum_m x_{2ik} z_{2km} p_{2km} + \sum_j v_j w_{2ij} \leq \sum_l \sum_m x_{1il} z_{1lm} p_{1lm} + \sum_j v_j w_{1ij}, \quad (3.20)$$

if  $y_{2i} = 0$ ,

$$\sum_m p_{1lm} = 1 ; \quad \sum_m p_{2km} = 1 ; \quad (3.21)$$

$$\sum_j w_{1ij} = 1 ; \quad \sum_j w_{2ij} = 1. \quad (3.22)$$

The optimization yields estimates of  $\hat{\underline{\rho}}_1$ ,  $\hat{\underline{\rho}}_2$ ,  $\hat{\underline{w}}_1$ , and  $\hat{\underline{w}}_2$ , from which we obtain estimates  $\hat{\underline{\beta}}_1$ ,  $\hat{\underline{\beta}}_2$ ,  $\hat{\underline{\epsilon}}_1$ , and  $\hat{\underline{\epsilon}}_2$  using Equations 3.15 and 3.16.

The following conditions ensure consistency and asymptotic normality of our sample-selection GME estimator:

- (i) The errors' supports  $\underline{v}$  for each equation are symmetric around zero.
- (ii) The support space  $Z$  spans the true values for each one of the unknown parameters  $\underline{\beta} = (\underline{\beta}'_1, \underline{\beta}'_2)'$  and has finite lower and upper bounds ( $z_{111}$  and  $z_{11M}$  for  $\underline{\beta}_1$  and  $z_{2k1}$  and  $z_{2kM}$  for  $\underline{\beta}_2$ ).
- (iii) The errors are independently and identically distributed. [Note: this assumption *does not* restrict the errors to be uncorrelated across equations.]
- (iv)  $\text{Plim } (1/T)X'X$  exists and is nonsingular, where  $X$  is a block diagonal matrix consisting of  $X_1$  and  $X_2$ .

The proofs of consistency and asymptotic normality follow immediately from those in Golan, Judge, and Miller (1996), Golan, Judge, and Perloff (1997), Golan, Perloff, and Shen

(2001), and Mittelhammer and Cardell (1996). These asymptotic properties can also be established via the empirical likelihood approach (Owen 1990, 1991, Qin and Lawless 1994, and Golan and Judge 1996).

In general given the four conditions

$$\sqrt{T}(\underline{\beta} - \beta) \xrightarrow{d} N(\underline{0}, Q_{GME}),$$

where

$$Q_{GME} = \text{plim } T^{-1} \left( X' (\hat{\Sigma}^{-1} \otimes I_T) X \right)^{-1}$$

is the asymptotic covariance matrix for the GME. Since  $\underline{\beta}$  is a continuous function of  $\underline{\lambda}$  (the Lagrange multipliers), this statement follows immediately from Qin and Lawless (1994, Lemma 1 and Theorem 1). The asymptotic variances are

$$\hat{\sigma}_\delta^2 = \frac{1}{T - K_\delta} \sum_t \hat{e}_{\delta i}^2,$$

for  $\delta = 1, 2$ , where  $\hat{e}_{\delta i} \equiv \sum_j v_j \hat{w}_{\delta ij}$  and  $K_\delta = L$  where  $\delta = 1$  or  $K$  where  $\delta = 2$ .

The solution to the GME estimation problem is unique and therefore allows us to identify all the estimates. By using the GME model (3.17)-(3.22), we convert an ill-posed or under-determined problem into a well-posed problem where the number of unknowns equals the number of first-order equations. Further, by using equations (3.19)-(3.20), we are able to specify the economic model in a more natural way and gain efficiency (as we demonstrate below).

Unlike the likelihood approaches, the GME method requires no assumption about the distribution or the covariance between the two equations. The GME differs from other models in how it handles identification. Most previous approaches use an exclusion restriction to identify the outcome equation in the sample selection model. On the other hand, the GME approach achieves identification from the inequality structure of Equations 3.19-3.20, which allows the covariance elements to be nonzero.

GME's only disadvantage is that the computation time increases markedly as the number of observations increases. However, if one is analyzing a single data set (rather than running simulations), the increase in time is not a major consideration.<sup>2</sup>

#### 4. SAMPLING EXPERIMENTS

In recent years, there have been several Monte Carlo studies of sample selection estimators for relatively large data sets. These studies include Hay, Leu, and Rohrer (1987), Manning, Duan, and Rogers (1987), Hartman (1991), and Leung and Yu (1996). Their results differ because of differences in their experimental designs.

##### *4.1 Experimental Designs*

Leung and Yu (1996) argue that several of the earlier studies that found superior performance of ordinary least squares (OLS) over maximum likelihood (ME) sample-selection estimators was due to unusual experimental designs. In particular, they argue that studies

---

<sup>2</sup> For example with a Pentium 4 1.7 GHz PC with 256 MB of RAM, it takes 0.109 seconds of execution time to solve the system of equations with  $K = L = 4$  and 100 observations, 0.938 seconds to solve a system with  $K = L = 7$  and 500 observations, and 3.25 seconds to solve  $K = L = 7$  and 1,000 observations.

such as Manning, Duan, and Rogers (1987) got their results because they drew regressors from a uniform distribution with a range of  $[0, 3]$ . Because this range is narrow, the covariates are highly collinear and the Mills' ratio term used in two-stage estimators is highly correlated with the regressor. Leung and Yu find that the ME sample-selection estimators perform better than OLS when they draw regressors from a larger range,  $[0, 10]$ . In order to give maximum likelihood estimators the greatest possible advantage, we use Leung and Yu's larger range for the right-hand-side variables.

Following Leung and Yu, most of our experiments involve only a single regressor (in addition to the intercept) in both the choice and level equations, so  $L = K = 2$ . In all designs,  $\beta_{12} = \beta_{22}$ . We vary  $\beta_{11}$ ,  $\beta_{21}$ , and the intercepts to control the level of censoring. The support spaces for  $\underline{z}_1$  and  $\underline{z}_2$  are all specified to be symmetric about zero with large negative and positive bands,  $\underline{z}_{1m} = \underline{z}_{2m} = (-100, -50, 0, 50, 100)'$  for all the unknown  $\underline{\beta}_1$  and  $\underline{\beta}_2$ . These supports reflect our state of ignorance relative to the unknown  $\underline{\beta}$ 's in the range  $[-100, 100]$ . The support spaces for the errors  $\underline{e}_1$  and  $\underline{e}_2$  are based on the *empirical* standard deviations of the observed  $y_{2i}$ ,  $\sigma_2^*$ , such that  $\underline{v}_1 = \underline{v}_2 = (-3\sigma_2^*, 0, 3\sigma_2^*)'$  for all  $i = 1, 2, \dots, T$ .

We used the computer program GAMS to generate the data. We repeated each experiment 1,000 times. To show the robustness of the GME estimator, we repeated the experiments for different right-hand-side variables, different number of observations, different number of regressors, normal and non-normal distributions, and for correlations between  $\underline{e}_1$  and  $\underline{e}_2$  of  $\rho = 0$  and  $\rho = 0.5$ . Table 1 describes the various designs.

We use the performance criteria of Leung and Yu (1996) to summarize the performance of each experiment. The first measure is the mean square error (MSE) for the relevant

equation. We also use the slope parameter bias and its mean square error, where  $\text{Bias}(\hat{\beta}_{22}) \equiv \hat{\beta}_{22} - \beta_{22}$  (and the subscript "22" indicates the second coefficient in the  $\underline{\beta}_2$  vector from the wage equation). The final criterion is the mean square prediction error (MSPE) for the second equation where

$$MSPE \equiv \frac{1}{1000} \sum_{i=1}^{1000} [\hat{E}(y_{2i}) - E(y_{2i})]^2. \quad (4.1)$$

#### 4.2 Alternative Estimators

We compare our new estimator to alternative parametric and semiparametric estimators. The alternative estimators include OLS, Heckman's two-step approach (2-Step) method, full-information maximum likelihood (FIML), and a semiparametric estimator with a nonparametric selection mechanism (AP) due to Ahn and Powell (1993).

The simplest alternative is to estimate the second equation using ordinary least squares, ignoring the sample-selection problem. We used GAMS to estimate both the GME and OLS models.

The two most commonly used likelihood-based, parametric approaches are the Heckman two-step and maximum likelihood estimators. We estimated these models using the computer program Limdep. Because of the relatively small sample sizes, these parametric estimation methods often failed to produce plausible estimates. Indeed, as Nawata and Nagase (1996) show, the FIML estimator may not converge or may converge to a local optimum. They use Monte Carlo experiments to show that FIML may not be a proper estimator when there is a high degree of multicollinearity between the estimated indicator

value of the first equation and the right-hand-side variable in the second equation. For these two estimators, we reject an estimate if Limdep reports a failure to converge or if the estimated correlation coefficient between the two errors does not lie within the range  $(-1, 1)$ . In our experiments, these failures are due primarily to small sample sizes. This failure virtually disappears with samples of 500 or 1,000 observations.

As the GME estimator can be viewed as an estimator from the class of semiparametric estimators, we also compare the sampling experiments with the Ahn and Powell (1993) estimator. The AP approach is designed to deal with a well-known problem of the parametric likelihood estimator, which assumes that the errors in the two equations are jointly normally distributed. If the joint distribution of the error terms is misspecified, these parametric estimators are inconsistent. Ahn and Powell (1993) propose a two-step approach where both the joint distribution of the error term and the functional form of the selection equation is unknown.

The AP estimator is robust to misspecification of the distribution of residuals and the form of the selection equation. When the distribution of residuals is not normal and the sample size is large, we expect the AP estimator to perform better than FIML and 2-Step estimators. However, when the sample size is small, it is not clear whether AP would dominate FIML and 2-Step estimators, as the large sample size requirement for the AP estimator is not met. So far as we know, no previous study has examined the small-sample performance of the AP estimator.<sup>3</sup>

---

<sup>3</sup> We use Matlab to obtain the AP estimates in our experiments, where the kernel functions are taken to be normal density functions. Following Ahn and Powell (1993), in the first-step kernel regression, the data were first linearly transformed so that the components of



Although it is possible to estimate the structural form of the selection equation using Heckman's method, the AP method only identifies the outcome equation. Moreover, since the outcome equation is estimated by regressing differences in dependent variables on corresponding differences in explanatory variables in the AP method, the intercept term is not identified. Hence, not all statistics reported for the other estimators are reported for the AP method in our tables.

In the GME model, the structural form of the selection model and  $\beta_1$  are estimated and there is no need for any weighing procedure.<sup>4</sup> The unknown Lagrangean multipliers (one for each observation) in the GME are the implicit and natural weight of each observation (Golan and Judge, 1996).

---

the vector of exogenous variables were orthogonal in the sample, with variances that were equal one for each component. The first-step bandwidth parameter  $h_1$  was selected in each iteration by least-square cross validation over a crude grid of possible values. The choice of the second-step kernel bandwidths,  $h_2$ , is less straightforward. Cross validation does not necessarily produces the best bandwidths (Powell and Stoker, 1996). We set  $h_2 = 0.7$ . We experimented with different values of  $h_2$  between 0.0005 and 1. For each of these values between 0.0005 to 1, the point estimates were equal to the ones presented here up to the third decimal point. When selection depends on one variable only, the first step is not needed. The second step weights can be obtained by conditioning on the selection equation regressor. The one-step estimator is asymptotically identical to the two-step estimator. For those experiment designs where selection depends on only one variable, we calculated both one- and two-step estimators and got similar results. To be consistent across designs, we report the two-step estimates in all tables.

<sup>4</sup> The coefficient  $\beta_1$  can be identified in the AP method if the symmetry restriction of the GME approach is imposed along with a restriction on the selection equation of a single index. For example, Chen (1999) assumes that the joint distribution of the residuals is symmetric around the origin and depends via a linear index on the exogenous variables of the model. Chen shows that under these symmetry and single index assumptions the intercept term for the outcome equation can be estimated. Thus, one could estimate a single index first stage (e.g., Heckman's probit) and then estimates the second stage using AP with a symmetry assumption. Unlike the original AP estimator, this one is not semiparametric in all stages.

### 4.3 Discussion of Results

In all experiments, the 2-Step and FIML approaches failed to estimate a large proportion of the samples while the OLS, AP, and GME models *always* produced estimates.<sup>5</sup>

When we compare various measures of fit in the following discussion, we discuss only the "plausible" 2-Step and FIML repetitions, which favors these two approaches substantially.

The summary statistics for the other approaches include the difficult samples that the 2-Step and FIML approaches could not handle.

Tables 2-5 report results for the five experimental designs. The first column shows the technique. The number in the parentheses following the FIML or 2-Step label is the percent of the 1,000 repetition where that estimator failed to converge and produce plausible values (the estimated correlation lies within the plausible range:  $|\rho| < 1$ ). The next two columns show the number of observations and the proportion of censored observations. The next two columns report the sum of the mean square errors for all the coefficients in the second equation including the constant,  $MSE(\hat{\beta}_2)$ , and for just the second coefficient in the second equation,  $MSE(\hat{\beta}_{22})$ . The following column shows the bias for this coefficient. The last column shows the mean squared prediction error, MSPE. We cannot report the  $MSE(\hat{\beta}_2)$  and MSPE for the AP approach because it does not produce an estimate of the constant term.

Because of the problem of obtaining convergence with the Heckman two-step and full-information estimators, we tried Nawata and Nagase's (1996) alternative method (which we estimated using Matlab). This estimator is labeled "NN" in Table 2. With this technique, we

---

<sup>5</sup> The problem of lack of convergence with the FIML approach is well known. The same problem occurs with the 2-Step procedure where there are relatively few observations with inadequate variation.

were more likely to find convergence or "plausible" estimates, though many estimates have  $|\rho| = .99$  for small samples. Moreover, the bias increased. Presumably the reported bias increased because we calculated the bias of the traditional techniques conditioned on dropping "implausible" estimates. Because the NN approach does not solve the problems with the likelihood approach for small samples, we do not report the estimates for subsequent tables.

In all the sampling experiments reported in Tables 2-5, the GME strictly dominates the OLS (hence we do not discuss this comparison further). In virtually all of the sampling results, the GME has smaller  $MSE(\hat{\beta}_2)$  and  $MSE(\hat{\beta}_{22})$ , indicating the stability of the GME estimator relative to the other estimators for both the whole vector of estimated parameters and of the slope parameter  $\beta_{22}$  ( $\beta_{21}$  is the intercept). Further, the bias of the GME estimator is smaller than for the other estimators in many cases.

In general, the GME dominates the AP. The AP method is designed to provide robust estimates with large samples and has been shown to perform well with large samples. However, it performs relatively poorly in our small-sample experiments, presumably because it imposes very little structure on the data (particularly linearity of the selection equation function). The AP bias is lower than OLS bias in all designs, but does not otherwise provide significantly better results than OLS.

The objective of the 2-Step and FIML estimators is to maximize prediction within the sample. It is, therefore, not surprising that the likelihood methods produce the best results in terms of the MSPE in most experiments (where we compare just the successful likelihood estimates to all the estimates for the alternative approaches). The fraction of repetitions for which the 2-Step and FIML estimators fail to provide plausible estimates is very large,

ranging from 11% to 99%. In samples of 20 and 50 observations, both the 2-Step and FIML estimators fail to provide plausible estimates in more than half of the repetitions. As the sample size increases, the percentage of implausible estimates decreases.

We now discuss each of the different experiments in more detail. Table 2 reports the effect of sample size on the estimates in Experimental Design 1. The GME approach dominates the AP method on all criteria except bias in the single case where  $T = 20$ . For the smallest sample size,  $T = 20$ , the GME estimator is superior to the likelihood methods based on all criteria. For  $T = 50$ , GME is superior on all criteria except MSPE. For sample sizes of  $T = 75$  and larger, the GME is superior in terms of MSE while the likelihood estimators (where they work at all) have smaller bias and better MSPE.

Table 3 reports the effect of the level of censoring on the estimates. In general, the performance of the estimators, as measured by MSPE, gets worse as the proportion of censored observations increases. The results are similar to Table 1 where the GME always has the lowest MSE, thus exhibiting the highest level of stability from sample to sample. We view this result as a strong one because the GME is superior even for a small proportion of censored observations.

In Table 4, we investigate the robustness of the estimators to various distribution or ill-posed specifications. The first row reports results based on the  $\chi^2_{(4)}$  distribution normalized to have a unit of variance. Again the GME is the most stable estimator while the two likelihood estimators have the smallest MSPE. Surprisingly, where they work, the likelihood approaches often perform better than the AP method in terms of bias (though not in terms of the MSE) even when the distribution is misspecified. This result may be due to the small

sample size. In small samples, the first-step, AP nonparametric regression estimator is likely to be imprecise. Because they impose "less structure" on data, nonparametric estimators typically need many observations to achieve good precision levels (Silverman, 1986). Hartman (1991), using a different experiment with a sample of a thousand, found that maximum likelihood performed badly with a misspecified error distribution, particularly with respect to the MSE.

The second row of Table 4 shows results based on  $\underline{x}_1 = \underline{x}_2$ . In this case the GME is superior to all other estimators under all the different statistics reported. Further, both the FIML and 2-Step estimators "work" only for a very small proportion of the samples. This last result is not surprising because the likelihood estimators are not identified (or are identified only by the nonlinearity of the inverse Mills ratio). As the AP estimator does not impose any restriction on the form of the selection equation, it is not identified in the limit. We report results for AP for completeness.

In Table 5, we compare estimators where  $K = 3$ ,  $\rho = 0$  or  $0.5$ , and  $T = 50$  or  $100$ . The GME dominates the other methods in terms of mean square error (except in the  $\rho = 0$ ,  $T = 100$  case), while the maximum likelihood (in the relatively few repetition where it produces plausible estimates) is superior in terms of prediction.

## 5. EMPIRICAL APPLICATIONS

We used each approach on three empirical applications with small data sets drawn from the March 1996 Current Population Survey (CPS). In each application, we estimated the wage-participation model (Equation 2.3 and 2.4) for the subset of respondents in the labor

market. In all three application, we exclude from the sample workers who are self-employed.<sup>6</sup> In no case did the normal maximum likelihood estimator converge, so we report results for only the OLS, Heckman two-step, AP, and GME models.

### *5.1 Teenage Heads of Households*

We first examine the labor market behavior of female teenagers who are heads of households. Recently, political debates concern the relationship between the increasing number of teen-age pregnancies and the generosity of welfare payments. Knowing which teenagers choose to work and how much they earn may help inform such a debate.

The wage equation covariates include years of education, a dummy for currently enrolled in school, potential experience ( $\text{age} - \text{education} - 6$ ) and potential experience squared, a dummy for Black, a dummy for rural location, a dummy for central city location, and a dummy for U.S. citizenship. 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 teenagers, and the number of children. The March 1996 CPS has 43 female teenagers who are head of an household for whom all the relevant variables are available. Of these, 29 are employees.

---

<sup>6</sup> There are 24,740 observations in the original March 1996 CPS sample. For Tables 6 through 8, we dropped (in order) 1,886 self-employed people and the 10,447 men. Then, to obtain the sample for Table 6, we also dropped (in order) 11,485 people who were 18 or younger, 866 who were not heads of households, and 13 people for whom relevant variables were missing. To obtain the sample for Table 7, we also dropped (in order) 12,298 people who had not immigrated to the United States within the last 5 years and 2 people who had missing relevant variables. To obtain the sample for Table 8, we also dropped (in order) 12,252 non-Native Americans and 5 people who had missing relevant variables.

Table 6 shows the estimated wage equation coefficients and asymptotic standard errors and an selection outcome table representing the probability of correct prediction. Except for the GME, in all models, an individual is assigned to a category if the fitted probability of being in that category is greater than 0.5. With the GME, we determine the category directly from the inequalities. GME predicts selection substantially better than the probit model and AP. As was discussed above, the intercept is not identified for AP (consequently, the MSPE is not identified either). Heckman's two-step estimator failed to yield an estimated correlation coefficient,  $\rho$ , between -1 and 1, so the table reports a  $\rho$  that is censored at 1.

None of the estimators finds a positive return to education that is statistically significantly larger than zero using the 0.05 criterion. Only the AP and GME methods find statistically significant experience effects. All find a large positive, statistically significant central city effect. The coefficient on Black is positive and surprisingly large, however only the GME estimate is statistically significantly different than zero at the 0.05 level. Only the AP and GME methods find a statistically significant positive return to U.S. citizenship.

## 5.2 *Recent Immigrants*

Next we analyze a sample of 107 female immigrants who entered the United States in the five years preceding the interview (27 of whom are in the labor force). Although there is now a significant literature on labor market performance of recent immigrants, most of the research has been conducted on men rather than women. Table 7 reports estimates for the same model as for the teenagers. Return to education are positive (though only the GME estimates is statistically significant) and smaller than the 8 to 10% returns usually reported in

the literature for U.S.-born workers. Once again, the GME methods predicts selection better than do the Heckman two-step and AP models.

### *5.3 Native Americans*

Finally, we analyze a sample of 151 Native American females, of whom 65 are in the labor force. We are unaware of any previous study of wages and participation by Native American females. We dropped the (irrelevant) race and U.S. citizenship dummy variables. The estimated return to education is around 5% across estimation methods, but only statistically significantly different from 0 for the OLS and GME estimators. All estimators show a positive rural effect, a negative center city effect, and a negative enrolled in school effect, but only the GME estimates are statistically significant. Again, the GME does a superior job of predicting selection.

### *5.4 Summary*

In all three of these applications based on small samples, we obtain fairly similar coefficients estimates (though the GME estimates tend to be slightly closer to the OLS estimates than to the other two sets), the GME does a better job of predicting labor force participation, and we cannot estimate Heckman's full-information maximum likelihood model. In one of these cases, Heckman's two-step model fails to produce a plausible estimate of the correlation coefficient, which brings the entire estimate into question. In each case, the GME's estimated asymptotic standard errors are much smaller than those of the other methods (followed by those of the OLS and AP).



## 6. CONCLUSION

In a large number of empirical economic analysis, the data sets are relatively small and there is a need for a stable and consistent estimator that converges to an optimal solution and performs well under these circumstances. Our new generalized maximum entropy (GME) estimator meets this objective. For small samples, the GME approach has smaller mean square error measures than other well-known estimators such as ordinary least squares, Heckman's 2-step method, full-information maximum likelihood, and Ahn and Powell's method.

We compared GME to these alternative estimators in small sample experiments. All but one of our experimental designs uses a normal distribution, which favors the likelihood approaches. In these small samples, the OLS, Ahn and Powell, and GME methods always work, but the 2-Step and FIML methods frequently fail to converge or provide estimates of the correlation coefficient that do not lie within the plausible range.

Under all scenarios, the GME proved to be the most stable estimator (had the lowest variance and mean square errors), while the likelihood approaches predicted within the sample better when it worked at all (except for small sample sizes where the GME out-performed the other estimators under all criteria). The GME approach performed better than the OLS in all cases and better than the AP estimator in most cases. Finally, the GME works where the right-hand-side variables are identical in both equations, a situation where the likelihood methods cannot work at all and the AP method does not perform as well. Thus, if precision and stability of the estimates of a sample-selection model based on a relatively small data set are the objective, the GME estimator appears to be the appropriate choice.

References

- Ahn, H. and J. L. Powell, "Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism," *Journal of Econometrics*, 58, 1993:3-29.
- Chen, S., "Distribution-Free Estimation of the Random Coefficient Dummy Endogenous Variable Model," *Journal of Econometrics*, 91, 1999:171-199.
- Cosslett, S. R., "Distribution-free Maximum Likelihood Estimator of the Binary Choice Model," *Econometrica*, 51, 1981:765-782.
- Csiszár, I., "Why Least Squares and Maximum Entropy? An Axiomatic Approach to Inference for Linear Inverse Problems," *The Annals of Statistics*, 19, 1991:2032-2066.
- Golan, A., and H. Gzyl, "A Generalized Maxentropic Inversion Procedure for Noisy Data," *Applied Mathematics and Computation*, 127, 2002:249-260.
- Golan, A., and G. Judge, "A Maximum Entropy Approach to Empirical Likelihood Estimation and Inference," Working paper, University of California, Berkeley, 1996.
- Golan, A., G. Judge, and D. Miller, *Maximum Entropy Econometrics: Robust Estimation With Limited Data*, New York: John Wiley & Sons, 1996.
- Golan, A., G. Judge, J. M. Perloff, "Recovering Information from Censored and Ordered Multinomial Response Data," *Journal of Econometrics*, 79, 1997:23-51.
- Golan, A., J. M. Perloff, E. Z. Shen, "Estimating a Demand System with Nonnegativity Constraints: Mexican Meat Demand," *Review of Economics and Statistics*, 83, 2001:541-550.
- Han, A. K., "Non-parametric analysis of a generalized regression model: The maximum rank correlation estimator," *Journal of Econometrics*, 35, 1987:303-316

- Hartman, R. S., "A Monte Carlo Analysis of Alternative Estimators in Models Involving Selectivity," *Journal of Business & Economic Statistics*, 9, 1991:41-9.
- Hay, J., R. Leu, and P. Rohrer, "Ordinary least squares and sample-selection models of health-care demand," *Journal of Business & Economic Statistics*, 5, 1987:499-506.
- Jaynes, E. T., "Information Theory and Statistical Mechanics," *Physics Review*, 106, 1957a:620-630.
- Jaynes, E. T., "Information Theory and Statistical Mechanics, II," *Physics Review*, 108, 1957b:171-190.
- Jaynes, E. T., "Prior Information and Ambiguity in Inverse Problems," *Inverse Problems*, D. W. McLaughlin, ed., Providence, Rhode Island: American Mathematical Society, 1984.
- Kullback, J., *Information Theory and Statistics*, New York: John Wiley & Sons, 1959.
- Leung, S. F., and S. Yu, "On the choice between sample selection and two-part models," *Journal of Econometrics*, 72, 1996:197-229.
- Levine, R. D., "An Information Theoretical Approach to Inversion Problems," *Journal of Physics*, A, 13, 1980:91-108.
- Maddala, G.S. (1983): *Limited-Dependent and Qualitative Variables in Econometrics*, Cambridge: Cambridge University Press, 1983.
- Manning, W., N. Duan, and W. Rogers, "Monte Carlo Evidence on the Choice Between Sample Selection and Two-part Models," *Journal of Econometrics*, 35, 1987:59-82.
- Mittelhammer, R., and S. Cardell, "On the Consistency and Asymptotic Normality of the Data Constrained GME Estimator of the GML," working paper, Washington State University, Pullman, WA.

- Nawata, K., and N. Nagase, "Estimation of Sample Selection Bias Models," *Econometric Review*, 15, 1996:387-400.
- Owen, A.B., "Empirical Likelihood Ratio Confidence Regions," *Annals of Statistics*, 18, 1990:90-120.
- Owen, A.B., "Empirical Likelihood for Linear Models," *The Annals of Statistics*, 19, 1991: 1725-45.
- Powell, J. L., and T. M. Stoker, "Optimal bandwidth choice for Choice for Density-Weighted Averages," *Journal of Econometrics*, 75, 291-316, 1996.
- Qin, J., and J. Lawless, "Empirical Likelihood and General Estimating Equations," *The Annals of Statistics*, 22, 1994:300-325.
- Shannon, C. E., "A Mathematical Theory of Communication," *Bell System Technical Journal*, 27, 1948:379-423.
- Shore, J. E., and R. W. Johnson, "Axiomatic Derivation of the Principle of Maximum Entropy and the Principle of Minimum Cross-Entropy," *IEEE Transactions on Information Theory*, IT-26, 1980:26-37.
- Silverman, B. W., *Density Estimation for Statistics and Data Analysis*, Chapman and Hall, 1986.
- Skilling, J., "The Axioms of Maximum Entropy," in J. Skilling, ed., *Maximum Entropy and Bayesian Methods in Science and Engineering*, Dordrecht: Kluwer Academic, 1989:173-87.

**Table 1**  
**Experimental Design**

<i>Design</i>	<i>Regressors</i>	<i>T</i>	<i>K</i>	<i>Approximate Percent Censored</i>	<i>Error Distribution</i>	$\rho$
1	$x_1 \sim U(0,10), x_2 \sim U(0,10)$	20, 50, 75, 100, 125	2	50	bivariate $N$	0.5
2	$x_1 \sim U(0,10), x_2 \sim U(0,10)$	100	2	25, 50, 75	bivariate $N$	0.5
3	$x_1 \sim U(0,10), x_2 \sim U(0,10)$	100	2	50	bivariate $\chi^2$	0.5
4	$x_1 = x_2 \sim U(0,10)$	100	2	50	bivariate $N$	0.5
5	$x_1 \sim U(0,10), x_2 \sim U(0,10)$	50, 100	3	50	bivariate $N$	0, 0.5

**Table 2**  
**Sample Results for Experimental Design 1**

<i>Estimation Method</i>	<i>Number of Obs.</i>	<i>Proportion of Censored Obs.</i>	<i>MSE(<math>\hat{\beta}_2</math>)</i>	<i>MSE(<math>\hat{\beta}_{22}</math>)</i>	<i>Bias(<math>\hat{\beta}_{22}</math>)</i>	<i>MSPE</i>
FIML (0.956)*	20	44.86	1.624	0.056	-0.075	1.518
2-Step (0.960)*			3.515	0.711	-0.212	1.562
NN (0.089/0.223)#			0.670	0.019	-0.018	1.140
OLS			0.594	0.120	-0.026	1.009
AP				0.012	-0.011	
GME			0.558	0.011	-0.013	1.001
FIML (0.538)*			50	42.20	0.359	0.0072
2-Step (0.305)*	0.361	0.0072			-0.013	0.642
NN (0.002/0.043)#	0.373	0.0120			-0.047	1.064
OLS	0.339	0.0069			-0.026	1.021
AP		0.0069			-0.025	
GME	0.297	0.0063			-0.013	1.008
FIML (0.411)*	75	51.41			0.462	0.0085
2-Step (0.224)*			0.453	0.0083	-0.012	0.545
NN (0.003/0.003)#			0.474	0.0131	-0.064	1.089
OLS			0.454	0.0082	-0.036	1.076
AP				0.0080	-0.034	
GME			0.373	0.0071	-0.020	1.054
FIML (0.208)*			100	51.37	0.294	0.0047
2-Step (0.117)*	0.266	0.0045			-0.007	0.367
NN (0.003/0.002)#	0.290	0.0071			-0.053	1.042
OLS	0.309	0.0049			-0.036	1.055
AP		0.0049			-0.035	
GME	0.230	0.0040			-0.022	1.030
FIML (0.164)*	125	51.41			0.205	0.0035
2-Step (0.092)*			0.196	0.0034	-0.004	0.306
NN (0.000/0.000)#			0.241	0.0081	-0.056	1.077
OLS			0.228	0.0038	-0.023	1.036
AP				0.0037	-0.022	
GME			0.193	0.0035	-0.015	1.023

\* The fraction in parentheses indicates the share of repetitions for which this estimator failed to converge and produce "plausible" results.

# The first fraction is the share of repetition that did not converge, the section fraction is the share for which  $|\rho| = 0.99$ .

**Table 3**  
**Sample Results for Experimental Design 2**

<i>Estimation Method</i>	<i>Number of Obs.</i>	<i>Proportion of Censored Obs.</i>	$MSE(\hat{\beta}_2)$	$MSE(\hat{\beta}_{22})$	$Bias(\hat{\beta}_{22})$	$MSPE$
FIML (0.685)* 2-Step (0.830)* OLS AP GME	100	26.06	0.0932 0.0883 0.0939 0.0808	0.00194 0.00186 0.00197 0.00190 0.00178	-0.0071 -0.0059 -0.0179 -0.0172 -0.0112	0.2848 0.2860 0.9877 0.9839
FIML (0.792)* 2-Step (0.883)* OLS AP GME	100	51.37	0.2939 0.2663 0.3091 0.2301	0.00472 0.00445 0.00487 0.00487 0.00397	-0.0052 -0.0065 -0.0361 -0.0348 -0.0215	0.3662 0.3670 1.0550 1.0300
FIML (0.661)* 2-Step (0.823)* OLS AP GME	100	75.11	1.1321 1.1422 1.1177 0.9500	0.01679 0.01669 0.01611 0.01613 0.01440	-0.0073 -0.0062 -0.0356 -0.0351 -0.0177	0.3312 0.3309 1.3045 1.2480

\* The fraction in parentheses indicates the share of repetitions for which this estimator converged and produced "plausible" results.

**Table 4**  
**Sample Results for Experimental Designs 3 and 4**

<i>Design</i>	<i>Estimation Method</i>	<i>Number of Obs.</i>	<i>Proportion of Censored Obs.</i>	$MSE(\hat{\beta}_2)$	$MSE(\hat{\beta}_{22})$	$Bias(\hat{\beta}_{22})$	$MSPE$
3 $\chi^2_{(4)}$	FIML (0.759)*	100	51.29%	0.377	0.0051	-0.0077	0.3698
	2-Step (0.898)*			0.283	0.0042	-0.0088	0.3699
	OLS			0.252	0.0039	-0.0273	1.0476
	AP				0.0039	-0.0266	
	GME			0.192	0.0033	-0.0161	1.0275
4	FIML (0.123)*	100	49.92	9217.590	9193.80	30.9090	116.370
	2-Step (0.015)*			0.777	0.00401	0.0303	0.6106
	OLS			0.234	0.00227	-0.0019	1.1448
	AP				0.00233	-0.0019	
	GME			0.127	0.00234	0.0003	1.0419

\* The fraction in parentheses indicates the share of repetitions for which this estimator converged and produced "plausible" results.



**Table 5**  
**Sample Results for Experimental Design 5**

<i>Estimation Method</i>	<i>Correlation</i>	<i>Number of Obs.</i>	<i>Proportion of Censored Obs.</i>	$MSE(\hat{\beta}_2)$	$MSE(\hat{\beta}_{22})$	$Bias(\hat{\beta}_{22})$	$MSPE$
FIML (0.313)* 2-Step (0.538)* OLS AP GME	$\rho = 0$	50	49.79	0.43298 0.46207 0.44895 0.41531	0.00626 0.00727 0.00716 0.00710 0.00680	-0.0022 -0.0108 -0.0156 -0.0149 -0.0128	0.19302 0.19694 1.11194 1.07460
FIML (0.252)* 2-Step (0.411)* OLS AP GME	$\rho = .5$	50	49.52	0.45369 0.42840 0.43692 0.40531	0.01187 0.00649 0.00671 0.00673 0.00635	-0.0072 -0.0005 -0.0023 -0.0018 0.0004	0.37473 0.18853 1.05764 1.04195
FIML (0.583)* 2-Step (0.789)* OLS AP GME	$\rho = 0$	100	49.40	0.19493 0.19040 0.24461 0.19992	0.00333 0.00322 0.00442 0.00434 0.00371	-0.0093 -0.0124 -0.0383 -0.0370 -0.0300	0.09365 0.09459 1.09071 1.05963
FIML (0.569)* 2-Step (0.793)* OLS AP GME	$\rho = .5$	100	49.40	0.20378 0.20240 0.20345 0.17967	0.00348 0.00339 0.00357 0.00357 0.00319	-0.0077 -0.0097 -0.0218 -0.0210 -0.0155	0.09351 0.09379 1.03300 1.02113

\* The fraction in parentheses indicates the share of repetitions for which this estimator converged and produced "plausible" results.

**Table 6**  
**Wage Equation for Female Teen Heads of Households**  
 (N = 43; 29 in the labor force)

	<i>OLS</i>	<i>2-Step</i>	<i>AP</i>	<i>GME</i>
Constant	1.181 (1.070)	-0.165 (1.828)	NA NA	1.091 (0.002)
Education	0.017 (0.085)	0.115 (0.143)	0.016 (0.044)	0.0187 (0.0119)
Black	0.203 (0.186)	0.256 (0.241)	0.202 (0.135)	0.192 (0.074)
Experience	0.206 (0.141)	0.103 (0.207)	0.198 (0.123)	0.180 (0.068)
Experience Squared	-0.091 (0.041)	-0.062 (0.059)	-0.089 (0.032)	-0.085 (0.017)
Rural	0.085 (0.114)	0.149 (0.168)	0.080 (0.082)	0.080 (0.067)
Central City	0.324 (0.117)	0.411 (0.176)	0.323 (0.079)	0.332 (0.062)
Enrolled in School	-0.033 (0.128)	-0.197 (0.225)	-0.037 (0.112)	-0.043 (0.060)
U.S. Citizen	0.264 (0.186)	0.304 (0.259)	0.260 (0.086)	0.277 (0.109)
$\lambda$		0.443 (0.409)		
$\rho$		1		
R <sup>2</sup>	0.458	0.504		0.454
MSPE	0.058	0.043		0.039

<i>Actual</i>	<i>Predicted</i>					
	<i>2-Step's Probit</i>		<i>AP</i>		<i>GME</i>	
	0	1	0	1	0	1
0	9	5	0	14	12	2
1	2	27	0	29	3	26

**Table 7**  
**Wage Equation for Female Immigrants**  
 (N = 107; 27 in the labor force)

	<i>OLS</i>	<i>2-Step</i>	<i>AP</i>	<i>GME</i>
Constant	1.131 (0.648)	0.750 (1.350)	- -	1.190 (0.001)
Education	0.052 (0.034)	0.061 (0.042)	0.058 (0.056)	0.045 (0.010)
Black	0.051 (0.452)	0.210 (0.644)	0.080 (0.239)	0.051 (0.010)
Experience	0.014 (0.058)	0.025 (0.060)	0.027 (0.080)	0.021 (0.018)
Experience Squared	-0.0005 (0.001)	-0.001 (0.002)	-0.0009 (0.003)	-0.0008 (0.0005)
Rural	0.364 (0.709)	0.045 (0.272)	0.419 (0.301)	0.365 (0.292)
Central City	0.051 (0.319)	0.061 (0.272)	0.070 (0.356)	0.059 (0.049)
Enrolled in School	-0.080 (0.626)	-0.026 (0.559)	0.014 (0.545)	-0.065 (0.112)
$\lambda$		0.166 (0.536)		
$\rho$		0.290		
R <sup>2</sup>	0.165	0.168		0.155
MSPE	0.300	0.228		0.313

		<i>Predicted</i>					
		<i>2-Step's Probit</i>		<i>AP</i>		<i>GME</i>	
<i>Actual</i>		0	1	0	1	0	1
	0		75	5	80	0	75
1		19	8	27	0	0	27

**Table 8**  
**Wage Equation for Native Americans Females**  
 (N = 151; 65 in the labor force)

	<i>OLS</i>	<i>2-Step</i>	<i>AP</i>	<i>GME</i>
Constant	1.073 (0.394)	1.771 (0.677)	NA NA	1.128 (0.0003)
Education	0.055 (0.028)	0.043 (0.031)	0.044 (0.043)	0.052 (0.008)
Experience	0.038 (0.016)	0.023 (0.020)	0.038 (0.027)	0.037 (0.009)
Experience Squared	-0.001 (0.0003)	-0.0005 (0.0005)	-0.001 (0.001)	-0.001 (0.0002)
Rural	0.214 (0.130)	0.268 (0.196)	0.332 (0.175)	0.221 (0.072)
Central City	-0.170 (0.183)	-0.091 (0.200)	-0.171 (0.124)	-0.172 (0.052)
Enrolled in School	-0.290 (0.216)	-0.471 (0.279)	-0.190 (0.242)	-0.319 (0.085)
$\lambda$		-0.461 (0.344)		
$\rho$		-0.894		
R <sup>2</sup>	0.355	0.376		0.354
MSPE	0.157	0.135		0.144

<i>Actual</i>	<i>Predicted</i>					
	<i>2-Step's Probit</i>		<i>AP</i>		<i>GME</i>	
	0	1	0	1	0	1
0	66	20	68	18	60	26
1	24	41	23	42	0	65