

Partial Identification of Poverty Measures with Contaminated and Corrupted Data

Juan Carlos Chavez-Martin del Campo*
Department of Applied Economics and Management
Cornell University
jcc73@cornell.edu

September 8, 2004

Abstract

Much of the statistical analysis for poverty measurement regards the data employed to estimate poverty statistics as error-free observations. However, it is amply recognized that surveys responses are not perfectly reliable and that the quality of the data is often poor, especially for developing countries. Robust estimation addresses this problem by searching for poverty measures that are not highly sensitive to errors in the data. However, given the assumptions of robust estimation, the rationale for point estimation is not apparent. In the present study we tackle the problem by implementing a different strategy. Since a particular poverty measure is not point identified under the assumptions of robust estimation and some outcomes that are possible *ex ante* are ruled out *ex post*, we apply a fully non-parametric method to show that for the family of additively separable poverty measures it is possible to find identification regions under very mild assumptions. We investigate identifiability for the class of P_α poverty measures, showing that there exists an " α -ordering" for the size of the identification region. We apply two conceptually different confidence intervals for partially identified poverty measures: the first type of confidence interval covers the entire identification region, while the other covers each element of the identification region with fixed probability. The methodology developed in the paper is applied to analyze rural poverty in Mexico.

JEL Classification: C14, I32.

Keywords: *Poverty Measurement, Bounds, Partial Identification, Contamination Model, Identification Regions, Confidence Intervals.*

*I am grateful to Francesca Molinari and participants of the 2004 Latin American Meeting of the Econometric Society held in Santiago de Chile for helpful comments and suggestions. My research has been supported by El Consejo Nacional de Ciencia y Tecnologia (CONACYT) and the Ford/MacArthur/Hewlett Program of Graduate Fellowships in the Social Sciences. All remaining errors are my own.

1 Introduction

Since the appearance of Sen's seminal paper [21], research on poverty measurement has focused on the theoretical properties of aggregate poverty measures. Much of the statistical analysis of poverty measurement regards the data employed to estimate a specific poverty measure as error-free observations, implicitly assuming that the real problem to be concerned about is sample size. However, it is amply recognized that surveys responses are not perfectly reliable. Financial and technological constraints may affect the quality of the data, something that is particularly relevant for developing countries, making "truth" very difficult to measure. [1, 24, 22]

Measurement error has several dimensions for poverty estimation. For example, the poverty line is set for heterogenous groups of people without considering idiosyncratic differences in the cost of basic needs, arbitrary imputations are made when missing and zero outcomes appear in the sample, and the variable of interest is misreported by an important subset of survey respondents. [25]

Often the methodologies applied to solve these problems are arbitrary; at the same time, the results are highly sensitive to such adjustments. For instance, Szekely, Lustig, Cumpa and Mejia [25] applied several techniques to adjust for misreporting. In the case of Mexico, they found that, depending on the method for performing the adjustment, either 14 percent or 76.6 percent of the population is below the poverty line (in absolute terms it implies a difference of 57 million individuals). This has important policy implications since, depending on which of these numbers is used as a reference, the amount of resources directed to social programs can be considered either appropriate or totally insufficient.

Several approaches have been developed in order to analyze the effects of measurement error on poverty measurement. For instance, Chesher and Schluter [2] consider multiplicative measurement error distributed continuously and independently of true income to investigate the sensitivity of welfare measures to alternative amounts of measurement error. Ravallion [18] considers additive random errors when estimating individual-specific poverty lines, finding that heterogeneity in error distributions generates ambiguous poverty rank-

ings. An alternative approach, robust estimation, aims at developing point estimators that are not highly sensitive to errors in the data.¹ The objective is to guard against worse-case scenarios that errors in the data could conceivably produce. In that sense it takes an ex-ante perspective of the problem. Cowell and Victoria-Feser [3] apply this approach to poverty measurement by using the concept of the influence function to assess the influence of an infinitesimal amount of contamination upon the value of a poverty statistic [7]. They found that poverty measures that take as their primitive concept poverty gaps rather incomes of the poor are in general robust under this criterion. In particular, they proved that the class P_α of poverty measures developed by Foster, Greer and Thorbecke [6] is robust under data contamination.

In the present study we tackle the problem by implementing a different strategy. First, our approach does not consider classical measurement error, that is to say, we do not assume the existence of chronic errors affecting every observation. Instead of assuming that the error distributions have no mass point at zero, we consider the impact of intermittent errors by setting an upper bound to the proportion of gross errors within the data. Second, since the population parameters of interest are not point identified under the assumptions of robust estimation and some outcomes that are possible ex ante are ruled out ex post, we follow Horowitz and Manski [9] and apply a partial identification approach for poverty measurement². By using a fully non-parametric method, we show that for the family of additively separable poverty measures it is possible to find identification regions under very mild assumptions. In the case of additively poverty measures those identification regions are intervals whose lower and upper bounds are likely to be estimated from sample data. By using some results from the literature on order statistics, we construct confidence intervals asymptotically covering the entire identification region with fixed probability. Imbens and Manski have proposed a conceptually different type of confidence interval, one that covers the true value of the parameter with fixed probability. We extend their idea to the present

¹See Hampel et al [8] and Huber [10] for a comprehensive treatment of robust inference.

²Examples of applications of this approach in other settings are Molinari [17] and Dominitz and Sherman [5]. See Manski [16] for an overview of this literature

setting by imposing some pertinent assumptions.

The paper is organized as follows. Section 2 introduces some important concepts for poverty measurement. Section 3 states the problem formally, presenting both the contaminated and corrupted sampling models within the context of poverty measurement. Section 4 investigates the identification region for poverty measures belonging to the additively separable class. It is shown that, by using some stochastic dominance properties, we can find upper and lower bounds for poverty measures within that class. In section 5, we investigate identifiability for the class of P_α poverty measures, showing that there exists an " α -ordering" for the size of the identification region. Section 6 develops two conceptually different kinds of confidence intervals for partially identified poverty measures. Section 7 provides an empirical illustration by applying the methodology to the measurement of rural poverty in Mexico. Section 8 concludes. Most of the mathematical details are in the Appendix.

2 Poverty Measurement: Conceptual Framework

Let \mathcal{A} denote the σ -algebra of Lebesgue measurable sets on \mathfrak{R} . Let \mathcal{P} denote the set of all probability distributions on $(\mathfrak{R}, \mathcal{A})$. Thus for any $P \in \mathcal{P}$ the triple $(\mathfrak{R}, \mathcal{A}, P)$ is a probability space. Let $z \in \mathfrak{R}_{++}$ be the poverty line.

A person is said to be in poverty if her income, $y \in \mathfrak{R}$ or any other measure of her economic status is strictly below z . An aggregate poverty index is defined as a functional of the distribution $P \in \mathcal{P}$. Formally:

Definition 1 *A Poverty Index is a functional $\Pi(P; z) : \mathcal{P} \times \mathfrak{R}_{++} \rightarrow \mathfrak{R}$ that indicates the degree of poverty when a particular variable has distribution P and the poverty line is z .*

An important type of poverty measures is the *Additively Separable Poverty* class³, which is defined as follows:

$$\Pi(P; z) = \int \pi(y; z) dP \tag{1}$$

³Members of this class are the FGT, the Watts, and the Clark, Hemming and Ulph poverty measures. See Seidl [20] for a survey of poverty measures.

Where $\pi(y; z) : \mathfrak{R}_{++} \times \mathfrak{R} \rightarrow \mathfrak{R}$, is the poverty evaluation function for an individual, indicating the severity of poverty for a person with income y when the poverty line is fixed at z .

Since the axiomatic approach to poverty measurement proposed by Sen [21], most economists interested in the phenomenon of poverty have quantified poverty in a manner consistent with those principles. One of those principles, the *focus axiom*, requires a poverty measure to be independent of the income distribution of the non poor. The *monotonicity axiom* says that, everything else equal, a reduction in the income of a poor individual must increase the poverty measure; the *transfer axiom* emphasizes the positive effect of a regressive transfer on the poverty measure, that is to say, given other things, a pure transfer of income from a poor individual to any other individual that is richer must increase the poverty measure. Finally, Kakwani [13] has proposed a 4th property that prioritizes transfers taking place down in the distribution, other things being equal; this transfer sensitivity axiom argues that if a transfer $t > 0$ of income takes place from a poor individual with income y to a poor individual with income $y + \delta$ ($\delta > 0$), then the magnitude of the increase in poverty must be smaller for larger y .

3 Statement of the Problem

Let each member j of population J be characterized by the pair of welfare indicators (y_1^j, y_0^j) in the space $\mathfrak{R} \times \mathfrak{R}$ where y_1^j is the outcome of interest denoting the "true" equivalent income (or expenditure) for a given poverty line z . Let the random variable $(y_1, y_0) : J \rightarrow \mathfrak{R} \times \mathfrak{R}$ have distribution $P(y_1, y_0)$. Let a random sample be drawn from $P(y_1, y_0)$. Let's assume that instead of observing y_1 , one observes a random variable y defined by:

$$y \equiv wy_1 + (1 - w)y_0 \tag{2}$$

Realizations of y with $w = 0$ are said to be data errors, those with $w = 1$ are error-free, and y itself is a contaminated version of y_1 . Let $Q(y)$ denote the distribution of the observable y . Let $P_i = P_i(y_i)$ denote the marginal distribution of y_i . Let $P_{ij} =$

$P_{ij}(y_i | w = j)$ denote the distribution of y conditional on the event $w = j$ for $i = 0, 1$ and $j = 0, 1$. Let $p = P(w = 0)$ be the marginal probability of a data error. With data errors, the sampling process does not identify P_1 (the object of interest) but only $Q(y)$, the distribution of the observable y . By the law of total probability, these two distributions can be decomposed as follows:

$$P_1 = (1 - p)P_{11} + pP_{10} \tag{3}$$

$$Q(y) = (1 - p)P_{11} + pP_{00} \tag{4}$$

This problem can be approached from different perspectives. In robust estimation P_1 is held fixed and $Q(y)$ is allowed to range over all distributions consistent with both equations. In the context of poverty measurement, the objective would be to estimate the maximum possible distance between $\Pi(Q; z)$ and $\Pi(P_1; z)$. In contrast, the present analysis holds $Q(y)$ fixed because it is identified by the data, and P_1 is allowed to range over all distributions consistent with (3) and (4). This approach recognizes that the parameter of interest might not be point identified, but it can often be bounded.

The sampling process reveals only the distribution $Q(y)$. However, informative identification regions emerge if knowledge of the empirical distribution is combined with a non-trivial upper bound, λ , on p .

This investigations analyzes two different cases of data errors. In the first case, we will assume that the occurrence of data errors is independent of the sample realizations from the population of interest; formally:

$$P_1 = P_{11} \tag{5}$$

This particular model of data errors is known as "contaminated data" or "contaminated sampling" model. [10] In the other case, (5) does not hold and it is only assumed that there exists a non-trivial upper bound on the error probability. Horowitz and Manski [9] refer to this case as "corrupted sampling".

Define the sets

$$\mathcal{P}_1(p) \equiv \mathcal{P} \cap \{(1-p)\phi_{11} + p\phi_{10} : (\phi_{11}, \phi_{10}) \in \mathcal{P}_{11}(p) \times \mathcal{P}\} \quad (6)$$

$$\mathcal{P}_{11}(p) \equiv \mathcal{P} \cap \left\{ \frac{Q - p\phi_{00}}{1-p} : \phi_{00} \in \mathcal{P} \right\} \quad (7)$$

If there exists a non-trivial upper bound, λ , on the probability of data errors, then it can be proved that P_{11} and P_1 belong to the sets $\mathcal{P}_{11}(\lambda)$ and $\mathcal{P}_1(\lambda)$ respectively, where $\mathcal{P}_{11}(\lambda) \subset \mathcal{P}_1(\lambda)$. These restrictions are sharp in the sense that they exhaust all the available information, given the maintained assumptions [9].

4 Partial Identification of Poverty Measures

Suppose now that a proportion $p < 1$ of the data is erroneous. Furthermore, assume there exists a non-trivial upper bound, λ , for p , so $p \leq \lambda < 1$.⁴ From the analysis above, we know that the distribution of interest P_1 is not identified: i.e. $\mathcal{P}_1(\lambda)$ is not a singleton.

Even though P_1 is not identified, it is partially identified in the sense that it belongs to the identification region $\mathcal{P}_1(\lambda)$. There is a mapping from this set into the set of values in \mathfrak{R} of a given poverty measure. So the natural question is if there is a way to bound such values. As we will see below, it is possible to do it for the class of additively separable poverty indices for which the poverty evaluation function is decreasing by ordering the distributions in \mathcal{P}_λ according to a stochastic dominance criterion. Such criterion is defined as follows:

Definition 2 *Let $F, G \in \mathcal{P}$. Distribution F First Order Stochastically dominates (FOD) distribution G if*

$$F((-\infty, x]) \leq G((-\infty, x])$$

for all $x \in \mathfrak{R}$.

⁴In practice, upper bounds on the probability of data errors can be estimated from a validation data set or by the proportion of imputed data in the sample. See Kreider and Pepper [15] for an application of a validation model.

In the case of monotone functions, there is a well-known equivalent result for FOD that will be helpful to obtain the identification regions for the poverty measures:

Lemma 1 *The Distribution F first-order stochastically dominates the distribution G if and only if, for every non decreasing function $\varphi : \mathfrak{R} \rightarrow \mathfrak{R}$, we have*

$$\int \varphi(x)dF(x) \geq \int \varphi(x)dG(x) \quad (8)$$

Finally, let me introduce a basic concept that is a building block for identification regions.

Definition 3 *For $\alpha \in (0, 1]$, the α -quantile of $Q(y)$ is $r(\alpha) = \inf\{t : Q((-\infty, t]) \geq \alpha\}$.*

Now we can state the main result of this section. Following the approach of Horowitz and Manski [9] to find sharp bounds on parameters that respect stochastic dominance ⁵ we can construct identification regions for ASP measures.

Proposition 1 *Let it be known that $p \leq \lambda < 1$. Define probability distributions L_λ and U_λ on \mathfrak{R} as follows:*

$$L_\lambda = \begin{cases} \frac{Q(y \leq t)}{1-\lambda} & \text{for } t < r(1-\lambda) \\ 1 & \text{otherwise} \end{cases}$$

$$U_\lambda = \begin{cases} 0 & \text{for } t < r(\lambda) \\ \frac{Q(y \leq t) - \lambda}{1-\lambda} & \text{otherwise} \end{cases}$$

If $\Pi(P; z)$ belongs to the family of Additively Separable Poverty Measures and the poverty evaluation function is non-increasing in y , then identification regions for $\Pi(P_{11}; z)$ and $\Pi(P_1; z)$ are given by:

$$\mathbf{H}[\Pi(P_{11}; z)] = [\Pi_l(U_\lambda; z), \Pi_u(L_\lambda; z)] \quad (9)$$

and

$$\mathbf{H}[\Pi(P_1; z)] = [(1-\lambda)\Pi_l(U_\lambda; z) + \lambda\psi_0, (1-\lambda)\Pi_u(L_\lambda; z) + \lambda\psi_1] \quad (10)$$

⁵A parameter $\delta(\cdot)$ respects stochastic dominance if $\delta(F) \geq \delta(G)$ whenever F FOD G .

where $\psi_0 = \inf_y \pi(y; z)$ and $\psi_1 = \sup_y \pi(y; z)$.

Proof: See appendix.

These results are quite intuitive. The smallest feasible value of $\Pi(P_{11}; z)$ occurs when the subpopulation of persons with $w = 1$ have the smallest values of y . That is, when $P_{11} = L_\lambda$. Analogous reasoning gives the largest feasible value of $\Pi(P_{11}; z)$.

Example 1 Assume $P_1 = P_{11}$. Let $Q(y) = U[0, 1]$, $0 < p < \lambda < z < 1 - \lambda$. Let the poverty measure be given by $\varphi = \int_0^\infty 1(y < z) d\phi$. Then, $\varphi(P_1; z) \in [\frac{z-\lambda}{1-\lambda}, \frac{z}{1-\lambda}]$. If $P_1 \neq P_{11}$ then $\varphi(P_1; z) \in [z - \lambda, z + \lambda]$. Notice that $\varphi(Q; z)$ belongs to both intervals.

5 Identifiability of a Poverty Measure: An α -Ordering

Following Sen [21], there has been a widely use of distributive-sensitive poverty measures. This trend is epitomized by the class P_α of poverty measures developed by Foster et al [6], which is not only a member of the class of ASP poverty measures, but also one of the most widely-used in applied work.

Define $\Gamma = \{P \in \mathcal{P} : P((-\infty, y]) = 0, \forall y < 0\}$, i.e. the support of y is on \mathfrak{R}_+ . The P_α measure is given by

$$P_\alpha(F; z) = \int 1(y < z) \left(\frac{z-y}{z} \right)^\alpha dP \quad (11)$$

Where $\alpha \geq 0$ can be viewed as a measure of poverty aversion: The larger α , the greater the relative importance of the poorest individuals. Since P_α belongs to the class of ASP measures and its evaluation function is non-increasing in y , we can find its identification region in presence of contamination by Proposition 1. Define $g^\alpha = 1(y < z) \left(\frac{z-y}{z} \right)^\alpha$, $P_{\alpha\lambda}^L = \int g^\alpha dU_\lambda$, and $P_{\alpha\lambda}^U = \int g^\alpha dL_\lambda$. From Proposition 1, the identification region for P_α when the data is contaminated is given by:

$$\mathbf{H}[P_\alpha] = [P_{\alpha\lambda}^L, P_{\alpha\lambda}^U] \quad (12)$$

In the case of corrupted data the identification region is defined by the interval

$$\mathbf{H}[P_\alpha] = [(1 - \lambda)P_{\alpha\lambda}^L, (1 - \lambda)P_{\alpha\lambda}^U + \lambda] \quad (13)$$

We investigate the effects of contaminated and corrupted data on the identifiability of this poverty measure for different values of α . Let \mathcal{W} be the class of additively separable poverty measures with bounded evaluation functions and image given by the set $[\psi_0, \psi_1]$, where ψ_1 and ψ_0 are defined as above. For this class of poverty measures and given a distribution function, $F(y) \in \Gamma$, and a poverty line, z , we define a measure of identifiability as follows:

$$\chi = 1 - \frac{m(\mathbf{H}[\Pi])}{\psi_1 - \psi_0} \quad (14)$$

Where $\mathbf{H}[\Pi]$ is the identification region for some poverty measure $\Pi \in \mathcal{W}$, and $m : \mathcal{B} \rightarrow [0, \infty]$ is the Lebesgue measure on the Borel sets, \mathcal{B} , of \mathfrak{R} . Therefore, in the case of connected identification regions we have $m(\mathbf{H}[\Pi]) = \Pi_u - \Pi_l$.

Notice that the family of P_α poverty measures belong the \mathcal{W} class since $\psi_1 - \psi_0 = 1$ for all $\alpha \geq 0$ and $F(y) \in \Gamma$. For this particular case, our identifiability measure is given by:

$$\chi_\alpha = 1 - m_\alpha \quad (15)$$

Where $m_\alpha = m(\mathbf{H}[P_\alpha])$. When the data is either corrupted or contaminated we have the following result:

Proposition 2 *Let $Q \in \Gamma$ and $z \leq \max\{r(\lambda), r(1 - \lambda)\}$. If the data is either corrupted or contaminated, then $\chi_\alpha \geq \chi_\beta$ whenever $\alpha > \beta$.*

Proof: Define $\delta_{\beta-\alpha} = g^\beta - g^\alpha$. Since $m_\alpha = P_{\alpha\lambda}^u - P_{\alpha\lambda}^l$ when the data are contaminated we

have:

$$\begin{aligned}
m_\beta &= \int g^\beta dL_\lambda - \int g^\beta dU_\lambda \\
&= \int g^\alpha dL_\lambda + \int \delta_{\beta-\alpha} dL_\lambda - \int g^\alpha dU_\lambda - \int \delta_{\beta-\alpha} dU_\lambda \\
&= m_\alpha + \frac{1}{1-\lambda} \int \delta_{\beta-\alpha} 1(y \leq r(1-\lambda)) dQ - \frac{1}{1-\lambda} \int \delta_{\beta-\alpha} 1(y \geq r(\lambda)) dQ \\
&= m_\alpha + \frac{1}{1-\lambda} \int \zeta(y) dQ \\
&= m_\alpha + \frac{1}{1-\lambda} E_Q(\zeta(y))
\end{aligned}$$

Where $\zeta(y) = \delta_{\beta-\alpha}[1(y \leq r(1-\lambda)) - 1(y \geq r(\lambda))]$. It is easy to show that $\zeta(y) \leq 0$ for all $y \in \text{supp}Q$. Therefore we have $E_Q(\zeta(y)) \leq 0$, and the result follows in the case of contamination. When the data is corrupted, notice that $m(H(P_\alpha)) = (1-\lambda)m_\alpha + \lambda$, so $m(H(P_\beta)) - m(H(P_\alpha)) \leq 0$. \square

Since the poverty evaluation function is equal to zero for all $y \geq z$ and all α , the identifiability of the poverty measure is a monotone function of α down in the income distribution whenever $z \leq \max\{r(\lambda), r(1-\lambda)\}$. In other words, for this model of errors, there exists a positive relationship between the sensitivity of the poverty measure to the income of the poorest individuals (represented by the parameter α) and the identifiability of the same measure, χ_α . Therefore, a connection between the identifiability of the poverty measure and the axiomatic approach can be established. For example, it can be shown that for $\alpha > 0$, P_α satisfies the Monotonicity Axiom, the Transfer Axiom for $\alpha > 1$, and the Transfer Sensitivity axiom for $\alpha > 2$.

6 Confidence intervals for Partially Identified Poverty Measures

Let $(\mathfrak{R}, \mathcal{A}, Q)$ be a probability space, and let \mathcal{P} be a space of probability distributions. The distribution Q is not known, but a random sample y_1, y_2, \dots, y_N is available.

In the point identified case ($\lambda = 0$), a consistent estimator of the class of ASP measures

is given by

$$\hat{\Pi} = \frac{1}{N} \sum_{i=1}^N 1(y_i < z) \pi(y_i; z) \quad (16)$$

where $\pi(y; z)$ is a measurable function. By applying The Central Limit Theorem, the standard $100 \cdot \gamma\%$ confidence interval for $\Pi(P; z)$ is given by:

$$CI_{\gamma}^{\Pi} = \left[\hat{\Pi} - z_{\frac{\gamma+1}{2}} \frac{\sigma}{\sqrt{N}}, \hat{\Pi} + z_{\frac{\gamma+1}{2}} \frac{\hat{\sigma}}{\sqrt{N}} \right] \quad (17)$$

where z_{τ} is the τ quantile of the standard normal distribution.⁶

I will apply two conceptually different methodologies to estimate confidence intervals when data is contaminated. The first methodology considers symmetric confidence intervals for the entire identification region $\mathbf{H}[\Pi(P_1; z)]$. The second type of confidence interval, developed by Imbens and Manski [11], rather than cover the entire identification region with fixed probability γ , asymptotically covers the true value of the parameter with this probability. Besides, this type of confidence interval ensures that its exact coverage probability does converge uniformly to its nominal values. By doing so, one is able to avoid the problem of having wider confidence intervals when the parameter is point identified that when is set-identified.

For the first class of confidence intervals, I will make use of a result on L-statistics due to Stigler [23], who explores the asymptotic behavior of trimmed means. Define the confidence interval $CI_{\gamma}^{[\Pi_l, \Pi_u]}$ as

$$CI_{\gamma}^{[\Pi_l, \Pi_u]} = \left[\hat{\Pi}_l - z_{\frac{\gamma+1}{2}} \frac{\hat{\sigma}_l}{\sqrt{n}}, \hat{\Pi}_u + z_{\frac{\gamma+1}{2}} \frac{\hat{\sigma}_u}{\sqrt{n}} \right] \quad (18)$$

Where $\hat{\sigma}_l^2$ and $\hat{\sigma}_u^2$ are, respectively, consistent estimators for

$$\sigma_l^2 = \frac{Var_{U_{\lambda}}(\pi(y; z)) + (\pi(r(1-\lambda)) - \Pi_l)\lambda}{1-\lambda} \quad (19)$$

$$\sigma_u^2 = \frac{Var_{L_{\lambda}}(\pi(y; z)) + (\pi(r(\lambda)) - \Pi_u)\lambda}{1-\lambda} \quad (20)$$

⁶Kakwani [14] describes this methodology for ASP indices

Proposition 3 *Let $0 < \lambda < 1$ be known. Assume $E(\pi(y; z)^2) < \infty$. Let $r(1-\lambda)$ and $r(\lambda)$ be continuity points of $Q(y)$. Let the poverty evaluation function, $\pi(y; z)$, be a non-increasing function that is continuous at $r(1-\lambda)$ and $r(\lambda)$. Then*

$$\lim_{n \rightarrow \infty} Pr([\Pi_l, \Pi_u] \subset CI_\gamma^{[\Pi_l, \Pi_u]}) \geq \gamma \quad (21)$$

Proof: See appendix.

For the second type of confidence interval, define $\Delta = \Pi_U - \Pi_L$ and $\hat{\Delta} = \hat{\Pi}_U - \hat{\Pi}_L$ and consider the following set of assumptions:

Assumption 1 $Q(y) \in \mathcal{F}$, where \mathcal{F} is the set of distribution functions for which $E(|\pi(y; z)|^3) < \infty$, F'' is bounded in the neighborhoods of $r(\lambda)$ and $r(1-\lambda)$ while $F'(r(\lambda)) > 0$ and $F'(r(1-\lambda)) > 0$.

Assumption 2 $\underline{\sigma}^2 \leq \sigma_l^2, \sigma_u^2 \leq \bar{\sigma}^2$ for some positive and finite $\underline{\sigma}^2$ and $\bar{\sigma}^2$.

Assumption 3 $\Pi_u - \Pi_l \leq \bar{\Delta} < \infty$

Assumption 4 For all $\epsilon > 0$ there are $\nu > 0$, K and N_0 such that $N \geq N_0$ implies $Pr\left(\sqrt{N} |\hat{\Delta} - \Delta| > K\Delta^\nu\right) < \epsilon$, uniformly in $Q \in \mathcal{F}$.

Define the confidence interval \overline{CI}_γ^Π as:

$$\overline{CI}_\gamma^\Pi = \left[\hat{\Pi}_l - \frac{\overline{C}_N \hat{\sigma}_l}{\sqrt{N}}, \hat{\Pi}_u + \frac{\overline{C}_N \hat{\sigma}_u}{\sqrt{N}} \right] \quad (22)$$

where \overline{C}_N satisfies

$$\Phi\left(\overline{C}_N + \sqrt{N} \frac{\hat{\Delta}}{\max(\hat{\sigma}_l, \hat{\sigma}_u)}\right) - \Phi(-\overline{C}_N) = \gamma \quad (23)$$

Proposition 4 *Let $0 < \lambda < 1$. Let $r(1-\lambda)$ and $r(\lambda)$ be continuity points of $Q(y)$. Let the poverty evaluation function, $\pi(y; z)$, be a non-increasing function that is continuous at*

$r(1 - \lambda)$ and $r(\lambda)$. Suppose assumptions 1,2,3 and 4 hold. Then

$$\lim_{N \rightarrow \infty} \inf_{P \in \mathcal{P}} Pr \left(\Pi \in \overline{CI}_{\gamma}^{\Pi} \right) \geq \gamma \quad (24)$$

Proof: See appendix.

7 An Application to Rural Poverty in Mexico

The methodology developed in this paper is applied to the data obtained from the 2002 *Encuesta Nacional de Ingreso y Gasto de los Hogares* (ENIGH) held by INEGI [12]. This household income and expenditure survey is one of a series of surveys that are carried out under the same days of each year using identical sampling techniques.

The households are divided into zones of high and low population density. Low density population zones are those areas with fewer than 2500 inhabitants. It is common to identify these areas as rural ones. The rest of the zones (those with more than 2500 inhabitants) are identified as urban areas. The sample is representative for both urban and rural areas and at the national level. For the purposes of this study, we will just concentrate on the rural sub-sample which includes 6753 observations.

We have considered the extreme poverty line for rural areas constructed by INEGI-CEPAL for the 1992 ENIGH, following the methodology applied by SEDESOL [19] to inflate both the poverty line and all of the data into August 2000 prices. The rural poverty line is equal to 494.77 monthly 2002 pesos. In this paper we have used per capita current disposable income as indicator of economic welfare⁷. It is divided into monetary and non-monetary income. The monetary sources include wages and salaries, entrepreneurial rents, incomes from cooperatives, transfers and other monetary sources. Non-monetary incomes include gifts, autoconsumption, imputed rents and payments in kind.

The identification regions and the three different 95% confidence intervals for the class of FGT poverty measures are presented for both the contamination and the corruption models

⁷Due to lack of information, a final transformation of the original data was required: we will assume that each household member obtains the same proportion of total income as the others.

in Figures 1 and 2 respectively⁸. The contamination model applies if the occurrence of events that produces data errors is statistically independent of y_1 , the outcome of interest. The corruption model applies if the occurrence of those events is not statistically independent of y_1 . The first confidence interval corresponds to the point identified case ($\lambda = 0$). It is based on the point estimator ± 1.96 times its standard error. The second confidence interval is equal to the estimator of the lower bound minus 1.96, and the estimator of the upper bound plus 1.96 times their standard errors. The third confidence interval is the adjusted interval for the parameter \overline{C}_N .

Notice that for this particular data set $Q(z) + \lambda < 1$ for all $\lambda \leq .70$, so proposition 2 can be applied in this case. It is clear from the empirical analysis how the identifiability of the poverty measure, χ_α , is an increasing function of α , confirming the result stated in that proposition.

We found that there is almost no difference between the last two types of confidence intervals, that is to say, between the confidence interval covering the entire identification region and the one that provides the appropriate coverage for the parameter of interest.

It is clear from the empirical exercise that only considering random sampling errors without paying attention to the nature of measurement errors is very likely to produce considerable bias in our poverty estimation. This supports the view of the partial identification approach for which identifiability is not a matter of sample size, that is to say, the exact knowledge of a poverty measure cannot be inferred from any finite number of observations when the data is either corrupted or contaminated.

⁸We have no estimate of the frequency of data errors in the sample, so we present a sensitivity analysis using different values of λ .

Figure 1: Identification regions and confidence intervals for FGT poverty measures under contamination model: Rural Mexico, 2002

λ	$P_{\alpha\lambda}^L$	$P_{\alpha\lambda}^U$	χ_α	$CI_{0.95}^\Pi$	$CI_{0.95}^{[\Pi_L, \Pi_U]}$	$\overline{CI}_{0.95}^\Pi$
$\alpha = 0$						
0	0.287	0.287	1.000	[0.276, 0.298]		
0.01	0.282	0.289	0.993		[0.271, 0.300]	[0.272, 0.299]
0.02	0.275	0.292	0.983		[0.265, 0.304]	[0.266, 0.302]
0.03	0.268	0.294	0.974		[0.257, 0.306]	[0.259, 0.304]
0.05	0.252	0.299	0.953		[0.241, 0.311]	[0.243, 0.309]
0.07	0.234	0.304	0.930		[0.223, 0.316]	[0.225, 0.314]
0.10	0.209	0.312	0.897		[0.198, 0.325]	[0.200, 0.323]
$\alpha = 1$						
0	0.093	0.093	1.000	[0.089, 0.098]		
0.01	0.088	0.094	0.994		[0.084, 0.099]	[0.085, 0.098]
0.02	0.083	0.095	0.988		[0.079, 0.100]	[0.080, 0.099]
0.03	0.077	0.096	0.981		[0.074, 0.101]	[0.074, 0.100]
0.05	0.066	0.097	0.969		[0.062, 0.103]	[0.063, 0.102]
0.07	0.055	0.099	0.956		[0.052, 0.106]	[0.053, 0.105]
0.10	0.042	0.101	0.941		[0.039, 0.109]	[0.040, 0.108]
$\alpha = 2$						
0	0.042	0.042	1.000	[0.040, 0.045]		
0.01	0.038	0.043	0.995		[0.036, 0.046]	[0.036, 0.045]
0.02	0.034	0.043	0.991		[0.032, 0.047]	[0.033, 0.046]
0.03	0.031	0.043	0.988		[0.029, 0.048]	[0.029, 0.047]
0.05	0.024	0.044	0.980		[0.022, 0.049]	[0.022, 0.048]
0.07	0.018	0.045	0.973		[0.016, 0.050]	[0.017, 0.050]
0.10	0.011	0.046	0.965		[0.010, 0.053]	[0.011, 0.052]

Figure 2: Identification regions and confidence intervals for FGT poverty measures under corruption model: Rural Mexico, 2002

λ	$P_{\alpha\lambda}^L$	$P_{\alpha\lambda}^U$	χ_α	$CI_{0.95}^\Pi$	$CI_{0.95}^{[\Pi_L, \Pi_U]}$	$\overline{CI}_{0.95}^\Pi$
$\alpha = 0$						
0	0.287	0.287	1.000	[0.276, 0.298]		
0.01	0.279	0.296	0.983		[0.268, 0.307]	[0.270, 0.306]
0.02	0.270	0.307	0.963		[0.259, 0.318]	[0.261, 0.316]
0.03	0.260	0.316	0.944		[0.250, 0.327]	[0.251, 0.325]
0.05	0.239	0.334	0.905		[0.229, 0.345]	[0.231, 0.344]
0.07	0.218	0.352	0.866		[0.208, 0.364]	[0.209, 0.362]
0.10	0.188	0.381	0.807		[0.179, 0.393]	[0.180, 0.391]
$\alpha = 1$						
0	0.093	0.093	1.000	[0.089, 0.098]		
0.01	0.087	0.103	0.984		[0.083, 0.108]	[0.084, 0.107]
0.02	0.081	0.113	0.968		[0.077, 0.118]	[0.078, 0.117]
0.03	0.075	0.123	0.952		[0.071, 0.128]	[0.072, 0.127]
0.05	0.063	0.142	0.921		[0.059, 0.148]	[0.060, 0.147]
0.07	0.051	0.162	0.889		[0.048, 0.168]	[0.049, 0.167]
0.10	0.038	0.191	0.847		[0.035, 0.198]	[0.036, 0.197]
$\alpha = 2$						
0	0.042	0.042	1.000	[0.040, 0.045]		
0.01	0.038	0.052	0.986		[0.036, 0.055]	[0.036, 0.055]
0.02	0.034	0.062	0.972		[0.032, 0.066]	[0.032, 0.065]
0.03	0.030	0.072	0.958		[0.028, 0.076]	[0.028, 0.075]
0.05	0.022	0.092	0.930		[0.021, 0.097]	[0.021, 0.096]
0.07	0.016	0.112	0.904		[0.015, 0.117]	[0.015, 0.116]
0.10	0.010	0.141	0.869		[0.009, 0.147]	[0.010, 0.146]

8 Conclusions

In the last decade a growing body of research has studied inference in settings where parameters of interest are not point identified. The main contribution of this paper has been to bring about this literature in the context of poverty measurement.

When the observed data is corrupted or contaminated and without making parametric assumptions on the distribution from which the data are drawn, a particular poverty measure is not point identified. By applying the work on contaminated and corrupted samples developed by Horowitz and Manski [9], and using some properties common to an important subset of poverty measures, we have been able to identify bounds for the class of addi-

tively separable poverty indices. Moreover, we have shown that, for the class of P_α poverty measures, the more distributive-sensitive a poverty measure is, the smaller the size of the identification region under very plausible conditions.

We have extended two different confidence intervals to the setting of partially identified poverty measures. The first type of confidence interval provides coverage for the entire identification region, while the second one asymptotically covers the true value of the the poverty measure with fixed probability. We have illustrated the methodology developed in the paper with an application to rural poverty in Mexico. It is clear from both the theoretical and the empirical analysis that only considering random sampling errors without paying attention to the effects of measurement errors on poverty estimation is very likely to produce considerable bias in our inferences.

In future work, we plan to address questions about the identifying power of validation and covariate data, and monotonicity restrictions among other factors.

9 Appendix

Proof of Proposition 1: We need to show that $\Pi(U_\lambda; z) \leq \Pi(P; z)$ and $\Pi(L_\lambda; z) \geq \Pi(P; z)$ for all $P \in P_\lambda$. Set $\psi(y; z) = -\pi(y; z)$, so $\psi(y; z)$ is a non-decreasing function. By lemma 1, it is enough to prove that U_λ stochastically dominates every member of P_λ and L_λ is stochastically dominated by every member of that set. The rest of the proof is identical to proposition 4 in Horowitz and Manski [9]□

Before proving Propositions 3 and 4, we present a number of preliminary results. Let Y_1, Y_2, \dots, Y_n be *i.i.d.* random variables with distribution function $F(y)$. Let $Y_{(1)}, Y_{(2)}, \dots, Y_{(n)}$ denote the order statistics of the sample. Consider the trimmed mean given by

$$S_n = \frac{1}{[(\beta - \alpha)n]} \sum_{i=[\alpha n]+1}^{[\beta n]} Y_{(i)} \quad (25)$$

where $0 \leq \alpha < \beta \leq 1$ are any fixed numbers and $[\cdot]$ represents the greatest integer function. Let $r(\alpha)$ and $r(\beta)$ be continuity points of $F(y)$. Further, define

$$G(y) = \begin{cases} 0 & \text{if } y < r(\alpha) \\ \frac{F(y) - \alpha}{\beta - \alpha} & \text{if } r(\alpha) \leq y < r(\beta) \\ 1 & \text{otherwise} \end{cases}$$

and set

$$\mu = \int_{-\infty}^{\infty} y dG(y) \quad (26)$$

$$\sigma^2 = \int_{-\infty}^{\infty} y^2 dG(y) - \mu^2 \quad (27)$$

Lemma 2 (Stigler [23]) *Let Y_1, Y_2, \dots, Y_n be i.i.d. random variables with distribution function $F(y)$. then*

$$\begin{aligned} n^{\frac{1}{2}}(S_n - \mu) &\xrightarrow{d} N(0, (1 - \alpha)^{-2}((1 - \alpha)\sigma^2 + (r(\alpha) - \mu)^2\alpha(1 - \alpha))) && \text{if } \beta = 1 \text{ and } \int_0^{\infty} y^2 dF(y) < \infty \\ n^{\frac{1}{2}}(S_n - \mu) &\xrightarrow{d} N(0, (\beta)^{-2}((\beta)\sigma^2 + (r(\beta) - \mu)^2\beta(1 - \beta))) && \text{if } \alpha = 0 \text{ and } \int_{-\infty}^0 y^2 dF(y) < \infty \end{aligned}$$

Lemma 3 (Berry-Esseen for trimmed means, de Wet [4]) *Let Y_1, Y_2, \dots, Y_n be i.i.d. random variables with distribution function $F(y) \in \mathcal{F}$. Then*

$$\begin{aligned} \sup \left| Pr \left(\sqrt{N} \frac{(S_n - \mu)}{\sigma} < x \right) - \Phi(x) \right| &\longrightarrow 0 && \text{if } \beta = 1 \text{ and } \int_{r(\alpha)}^{\infty} |y|^3 dF(y) < \infty \\ \sup \left| Pr \left(\sqrt{N} \frac{(S_n - \mu)}{\sigma} < x \right) - \Phi(x) \right| &\longrightarrow 0 && \text{if } \alpha = 0 \text{ and } \int_{-\infty}^{r(\beta)} |y|^3 dF(y) < \infty \end{aligned}$$

For Lemma 4 define $\Delta = \theta_u - \theta_l$ and let $\hat{\theta}_u$ and $\hat{\theta}_l$ and $\hat{\Delta} = \hat{\theta}_u - \hat{\theta}_l$ be estimators for θ_l , θ_u and Δ and consider the following set of assumptions:

- i) There are estimators for the lower and upper bound $\hat{\theta}_l$ and $\hat{\theta}_u$ that satisfy: $\sqrt{N}(\hat{\theta}_l - \theta_l) \xrightarrow{d} \mathcal{N}(0, \sigma_l^2)$, and $\sqrt{N}(\hat{\theta}_u - \theta_u) \xrightarrow{d} \mathcal{N}(0, \sigma_u^2)$, uniformly in $P \in \mathcal{P}$ and there are estimators for σ_l^2 and σ_u^2 that converge to the true values uniformly in $P \in \mathcal{P}$.
- ii) For all $P \in \mathcal{P}$, $\sigma_l^2 \leq \bar{\sigma}_l^2$, $\sigma_u^2 \leq \bar{\sigma}_u^2$ for some positive and finite $\bar{\sigma}_l^2$ and $\bar{\sigma}_u^2$, $\theta_u - \theta_l \leq \bar{\Delta} < \infty$.
- iii) For all $\epsilon > 0$ there are $\nu > 0$, K and N_0 such that $N \geq N_0$ implies $Pr \left(\sqrt{N} |\hat{\Delta} - \Delta| > K\Delta^\nu \right) < \epsilon$, uniformly in $P \in \mathcal{P}$.

Define the confidence interval $\overline{CI}_\gamma^\theta$ as:

$$\overline{CI}_\gamma^\theta = \left[\hat{\theta}_l - \frac{\overline{C}_N \sigma_l}{\sqrt{N}}, \hat{\theta}_u + \frac{\overline{C}_N \sigma_u}{\sqrt{N}} \right] \quad (28)$$

where \bar{C}_N satisfies

$$\Phi(\bar{C}_N + \sqrt{N} \frac{\hat{\Delta}}{\max(\hat{\sigma}_l, \hat{\sigma}_u)}) - \Phi(-\bar{C}_N) = \gamma \quad (29)$$

Lemma 4 (Imbens and Manski, 2004) *Suppose assumptions i), ii), and iii) hold. Then*

$$\lim_{N \rightarrow \infty} \inf_{P \in \mathcal{P}} Pr(\theta \in \bar{CI}_\gamma^\theta) \geq \gamma \quad (30)$$

Proof of Proposition 3:

Define the events

$$A_n = \left\{ \Pi_l : \Pi_l \geq \hat{\Pi}_l - z_{\frac{\gamma+1}{2}} \frac{\hat{\sigma}_l}{\sqrt{n}} \right\}$$

$$B_n = \left\{ \Pi_u : \Pi_u \leq \hat{\Pi}_u + z_{\frac{\gamma+1}{2}} \frac{\hat{\sigma}_u}{\sqrt{n}} \right\}$$

From the definition of the confidence interval, $CI_\gamma^{[P_L, P_U]}$

$$Pr([\Pi_l, \Pi_u] \subset CI_\gamma^{[\Pi_l, \Pi_u]}) = Pr(A \cap B)$$

By Bonferroni's inequality:

$$Pr(A_n \cap B_n) \geq Pr(A_n) + Pr(B_n) - 1 \quad (31)$$

By lemma 3, $i = u, l$

$$\frac{\sqrt{n}(\hat{\Pi}_i - \Pi_i)}{\hat{\sigma}_i} \xrightarrow{d} \mathcal{N}(0, 1)$$

Therefore:

$$Pr([\Pi_l, \Pi_u] \subset CI_\gamma^{[\Pi_l, P_u]}) \geq 2\left(\frac{\gamma+1}{2}\right) - 1 = \gamma$$

asymptotically. \square

Proof of Proposition 4:

The result is a direct consequence of lemmas 3 and 4. \square

References

- [1] J. Bound, C. Brown, and Nancy Mathiowetz. Measurement error in survey data. In J. Heckman and E. Leamer, editors, *Handbook of Econometrics*, vol. 5. Elsevier Science, 2001.
- [2] Andrew Chesher and Christian Schluter. Welfare measurement and measurement error. *Review of Economic Studies*, 69:357–378, 2002.
- [3] Frank A. Cowell and Maria-Pia Victoria-Feser. Poverty measurement with contaminated data: A robust approach. *European Economic Review*, 40:1761–1771, 1996.
- [4] T. de Wet. Berry-esseen results for the trimmed mean. *S. Afr. Statist. J.*, 10:77–96, 1976.
- [5] J. Dominitz and R. Sherman. Nonparametric analysis of mixture models with verification. Mimeo, California Institute of Technology, 2003.
- [6] James Foster, Joel Greer, and Erick Thorbecke. A class of decomposable poverty measures. *Econometrica*, 52:761–766, 1984.
- [7] F. Hampel. The influence curve and its role in robust estimation. *Journal of the American Statistical Association*, 69:383–393, 1974.
- [8] F. Hampel, E. Ronchetti, P. Rousseeuw, and W. Stahel. *Robust Statistics*. Wiley, New York, 1986.
- [9] Joel E. Horowitz and Charles F. Manski. Identification and robustness with contaminated and corrupted data. *Econometrica*, 63:281–302, 1995.
- [10] P. Huber. *Robust Statistics*. Wiley, NY, 1981.
- [11] G. Imbens and C. Manski (forthcoming). Confidence intervals for partially identified parameters. *Econometrica*.

- [12] INEGI. *Encuesta Nacional de Ingresos y Gastos de los Hogares de 2002*. Instituto Nacional de Estadística, Geografía, e Informática, Mexico, 2002.
- [13] N. Kakwani. On a class of poverty measures. *Econometrica*, 1980.
- [14] N. Kakwani. Statistical inference in the measurement of poverty. *The Review of Economics and Statistics*, pages 632–639, 1993.
- [15] B. Kreider and J. Pepper. Inferring disability status from corrupt data. Mimeo, 2004.
- [16] Charles Manski. *Partial Identification of Probability Distributions*. Springer-Verlag, first edition, 2003.
- [17] Francesca Molinari. Identification of probability distributions with misclassified data. Mimeo, Cornell University, 2003.
- [18] M. Ravallion. Poverty rankings using noisy data on living standards. *Economics Letters*, 45:481–485, 1994.
- [19] SEDESOL. Nota técnica para la medición de la pobreza con base en los resultados de la encuesta nacional de ingresos y gastos de los hogares, 2002. Secretaría de Desarrollo Social, Mexico, 2002.
- [20] C. Seidl. Poverty measurement: A survey. In D. Boss, M. Rose, and C. Seidl, editors, *Welfare and Efficiency in Public Economics*. Springer-Verlag, 1998.
- [21] A. Sen. Poverty: An ordinal approach to measurement. *Econometrica*, 54, 1976.
- [22] T. Srinivassan. Data base for development analysis: An overview. *Journal of Development Economics*, 44:3–27, 1994.
- [23] S. Stigler. The asymptotic distribution of the trimmed mean. *The Annals of Statistics*, 1:472–477, 1973.
- [24] J. Strauss and D. Thomas. Measurement and mismeasurement of social indicators. *American Economic Review*, 86:30–34, 1996.

- [25] M. Szekely, N. Lustig, M. Cumpa, and J. Mejia. Do we know how much poverty there is. Mimeo, Inter-American Development Bank, 2000.