

# Measuring Poverty Dynamics with Synthetic Panels Based on Cross-Sections

*Hai-Anh Dang*  
*Peter Lanjouw*

The World Bank  
Development Research Group  
Poverty and Inequality Team  
June 2013



## Abstract

Panel data conventionally underpin the analysis of poverty mobility over time. However, such data are not readily available for most developing countries. Far more common are the “snap-shots” of welfare captured by cross-section surveys. This paper proposes a method to construct synthetic panel data from cross sections which can provide point estimates of poverty mobility. In contrast to traditional pseudo-panel methods that require multiple rounds of cross-sectional data to study poverty at the cohort level, the proposed method can

be applied to settings with as few as two survey rounds and also permits investigation at the more disaggregated household level. The procedure is implemented using cross-section survey data from several countries, spanning different income levels and geographical regions. Estimates fall within the 95 percent confidence interval—or even one standard error in many cases—of those based on actual panel data. The method is not only restricted to studying poverty mobility but can also accommodate investigation of other welfare outcome dynamics.

---

This paper is a product of the Poverty and Inequality Team, Development Research Group. It is part of a larger effort by the World Bank to provide open access to its research and make a contribution to development policy discussions around the world. Policy Research Working Papers are also posted on the Web at <http://econ.worldbank.org>. The authors may be contacted at [hdang@worldbank.org](mailto:hdang@worldbank.org) and [planjouw@worldbank.org](mailto:planjouw@worldbank.org).

*The Policy Research Working Paper Series disseminates the findings of work in progress to encourage the exchange of ideas about development issues. An objective of the series is to get the findings out quickly, even if the presentations are less than fully polished. The papers carry the names of the authors and should be cited accordingly. The findings, interpretations, and conclusions expressed in this paper are entirely those of the authors. They do not necessarily represent the views of the International Bank for Reconstruction and Development/World Bank and its affiliated organizations, or those of the Executive Directors of the World Bank or the governments they represent.*

# Measuring Poverty Dynamics with Synthetic Panels Based on Cross-Sections

Hai-Anh H. Dang and Peter F. Lanjouw\*

**Keywords:** transitory and chronic poverty, mobility, welfare, synthetic panels

JEL Codes: C53, D31, I32, O15

---

\*Dang ([hdang@worldbank.org](mailto:hdang@worldbank.org); corresponding author) and Lanjouw ([p.f.lanjouw@vu.nl](mailto:p.f.lanjouw@vu.nl)) are respectively economist and professor with the Poverty and Inequality Unit, Development Research Group, World Bank, and VU University Amsterdam. We thank Francois Bourguignon, Alan Dorfman, Chris Elbers, Francisco Ferreira, Gary Fields, Paul Glewwe, Bill Greene, Bo Honore, Dean Jolliffe, Aart Kraay, Christoph Lakner, Yue Man Lee, Michael Lokshin, David McKenzie, Reema Nayar, Tuoc Van Phan, Sergiy Radyakin, Carolina Sanchez-Paramo, Renos Vakis, Roy van der Weide, Nobuo Yoshida, and participants at meetings of the Econometric Society, International Association for Applied Econometrics, International Conference on Panel Data, and seminars at IFPRI, University of New South Wales, and World Bank for helpful discussions on earlier versions of this paper. We further thank Renos Vakis and Leonardo Lucchetti for their help with the Peruvian data.

## **I. Introduction**

To effectively reduce poverty, we want to understand the factors that help households escape poverty as well as those that induce them to remain in, or fall back into, poverty. It is commonly claimed that we need panel data to answer these questions, especially at the household or individual level. However, for a variety of reasons, cross-sectional data are far more common, in most developing countries, than panel data. Panel data collection can be very costly, for example, and can also pose a number of logistical and capacity-related challenges. For whatever reason, the scarcity of panel data has rendered the analysis of welfare dynamics difficult, if not impossible, in many developing country settings.

To overcome the non-availability of panel data, there have been a number of studies that develop pseudo-panels (or synthetic panels) out of multiple rounds of cross-sectional data. Following the seminal contributions of Deaton (1985), synthetic panels based on cohorts have been widely used to track income and consumption outcomes over time (e.g., Deaton and Paxson, 1994; Banks, Blundell, and Brugiavini, 2001; Pencavel, 2007). Other outcomes have also been analyzed with synthetic panels; these include, for example, labor responses to tax reforms (Blundell, Duncan, and Meghir, 1988), the returns to academic and vocational qualifications (Mcintosh, 2006), and household demand for private medical insurance (Propper, Rees, and Green, 2001). Notably, since cross-section samples are typically refreshed each time that the surveys are fielded, synthetic panels are possibly less exposed to the concerns surrounding attrition and measurement error that are often leveled

at actual panel data.<sup>1</sup> Thus, unsurprisingly, the econometrics of pseudo-panel data is a rapidly growing field of research.<sup>2</sup>

Perhaps because of their emphasis on cohorts rather than the household or individual, pseudo panel methods have not been widely applied to the study of poverty dynamics. Two notable exceptions are Bourguignon, Goh and Kim (2004) and Güell and Hu (2006) who construct synthetic panels at the household level. However, these two approaches require certain assumptions that may not always be easily satisfied in available cross sections: the former requires at least three rounds of cross section data and assumes a first-order autoregression (AR (1)) process through which past household or individual incomes (earnings) can affect present outcomes; the latter is exclusively restricted to duration analysis.

Against this background, a recent paper by Dang, Lanjouw, Luoto, and McKenzie (2014) proposes both parametric and non-parametric approaches to construct synthetic panels at the household level from two rounds of cross sections with rather parsimonious assumptions.<sup>3</sup> These synthetic panels are then used to predict lower-bound and upper-bound estimates of household transitions into and out of poverty in Vietnam and Indonesia. Drawing on both cross sectional data and genuine panel data, the authors compare mobility estimates based on synthetic panels to those that would be obtained from actual panel data and find that the “true” estimate of mobility (as revealed by the actual panel data) is generally sandwiched between the upper bounds and lower bounds based on the synthetic

---

<sup>1</sup> See, for example, Glewwe and Jacoby (2000) and Kalton (2009) for recent overviews of the advantages and disadvantages of cross sections and panel data in both developing and richer country contexts.

<sup>2</sup> See, for example, Inoue (2008) for a recent development and a brief review of this literature.

<sup>3</sup> This method focuses on temporal imputation or survey-to-survey imputation over time. For related studies that investigate spatial imputation or survey-to-census imputation, see, e.g., Elbers, Lanjouw, and Lanjouw (2003), some recent economic applications of which include Agostini and Brown (2010), Elbers et al. (2007), and Demombynes and Ozler (2005). More broadly, this method is related to the literature on identifying the bounds on the joint distribution for outcomes in different samples (see, e.g., Cross and Manski, 2002) and the statistical literature on imputing missing data (see, e.g., Little and Rubin, 2002). See also Ridder and Moffitt (2007) for a recent review on the econometrics of data combination.

panels.<sup>4</sup> In particular, the interval between the bound estimates can be narrowed if an appropriate range for the correlation between the error terms can be postulated. Dang et al. (2014) suggest that panel data, where available from other sources such as other countries with similar characteristics, might be scrutinized to identify such a narrower range. Despite its infancy, applications of Dang et al.’s method in various settings have been yielding encouraging results.<sup>5</sup>

In this paper we generalize the method introduced by Dang et al. in several important aspects. First, by introducing a method to approximate the appropriate correlation term and its theoretical upper bound using each country’s *own* cross sectional surveys, we overcome the limitations of having to draw on external estimates from actual panels in similar settings. Notably, such “similar” actual panels might neither be available, nor offer very convincing estimates (see further discussion in section III below). Thus with the more general framework offered here, we can apply our method to study poverty dynamics using *just only* two rounds of cross section data that satisfy our (rather standard) assumptions. The approximation of point values for the correlation term allows us to move beyond the *bound* estimates proposed by Dang et al. (2014) to actual *point* estimates of poverty mobility. This represents a major advance in terms of interpretation, and potential application in practice. In particular, we can investigate different measures of poverty dynamics, such as the population shares in different poverty status categories in both survey periods considered together (i.e., unconditional or joint probabilities) or the population shares in different poverty status categories in one period given their welfare status in the other period (i.e.,

---

<sup>4</sup> For presentational convenience, we generally refer to household movements between poverty or consumption categories as mobility.

<sup>5</sup> For example, recent applications/ validations of Dang et al.’s method against true panel data include Bierbaum and Gassmann (2012) for the Kyrgyz Republic, Cruces et al., (forthcoming) for Chile, Nicaragua, and Peru, and Martinez et al. (2013) for the Philippines.

conditional probabilities). The former measure provides one way to define chronic poverty,<sup>6</sup> and the latter permits the study of poverty transitions. We also provide the formulae for the standard errors on point estimates, which is not available with the bound estimates offered by Dang et al (2014).

Second, we generalize the construction of these synthetic panels to settings where more than two rounds of data are available. This offers useful results since few actual panel datasets in developing countries span more than two periods; and in those cases where they do, they may be suffering heavily from attrition problems. By considering poverty mobility in more than two periods, we can investigate richer inter-temporal profiles of movement into and out of poverty.

Third, our proposed framework further extends Dang et al. (2014)'s investigation of household transitions into and out of poverty to a much more general setup of household movements among different consumption groups. A typical description of the former is usually shown in a 2x2 poverty transition matrix—where household movement is tracked between a poor category or non-poor category from one period to the next—while an example of the latter is a 5x5 consumption transition matrix, where household movement is tracked not just for two poverty categories but for five consumption quintiles.

We validate our estimates both theoretically with a Monte Carlo simulation exercise and empirically by using cross sectional and actual panel survey data from several high-income and developing countries including Bosnia-Herzegovina, Lao PDR, Peru, the United States, and Vietnam. We find that our synthetic panel estimates are close to—and mostly lying within the 95 percent confidence intervals or even one standard error in many cases—those

---

<sup>6</sup> Also see Calvo and Dercon (2009) and Foster (2009) for more discussion on different definitions of chronic poverty. We restrict our discussion in this paper to a money-metric measure of poverty, for a multidimensional measure see Alkire and Foster (2011).

of actual panel data. Assuming our estimation model is valid, the standard errors on our model-based synthetic panel estimates are also found to be smaller than the sampling-based standard errors of actual panel data (an advantage derived due to the fact that sample sizes of cross-section survey data generally exceed those of the actual panel data).

This paper consists of seven sections. We start in the next section with a brief description of the method introduced in Dang et al., which is followed by our generalization of this method and estimation procedures in Section III; the Monte Carlo simulation exercise is provided in Section IV. We describe the data in Section V before applying our method to investigate poverty dynamics in Section VI, and Section VII concludes.

## II. Bound Estimates on Poverty Mobility

Let  $x_{ij}$  be a vector of household characteristics observed in survey round  $j$ ,  $j= 1$  or  $2$ , that are also observed in the other survey round for household  $i$ ,  $i= 1, \dots, N$ . A key requirement is that these household characteristics are time variant. Subject to data availability, these household characteristics can include such clearly time-invariant variables as sex, ethnicity, religion, language, place of birth, and parental education as well as variables that can be converted into time invariant versions based, for example, on information about household heads' age and education. The vector  $x_{ij}$  can also include time-varying household characteristics if retrospective questions about the round-1 values of such characteristics are asked in the second round survey.

Let  $y_{ij}$  then represent household consumption or income in survey round  $j$ . The linear projection of household consumption (or income) on household characteristics for each survey round is given by

$$y_{i1} = \beta_1' x_{i1} + \varepsilon_{i1} \tag{1}$$



$$y_{i2} = \beta_2' x_{i2} + \varepsilon_{i2} \quad (2)$$

Let  $z_j$  be the poverty line in period  $j$ . We are interested in knowing such quantities as

$$P(y_{i1} < z_1 \text{ and } y_{i2} > z_2) \quad (3a)$$

which represents the percentage of households that are poor in the first period but nonpoor in the second period (considered together for two periods), or

$$P(y_{i2} > z_2 \mid y_{i1} < z_1) \quad (3b)$$

which represents the percentage of poor households in the first period that escape poverty in the second period.

If panel data are available, we can easily estimate the quantities in (3a) and (3b); otherwise, we can use synthetic panels for this purpose. Assume that the underlying population being sampled in survey rounds 1 and 2 are the same, or more specifically,  $x_{i1} = x_{i2}$ , and  $y_{i1} \mid x_{i1}$  and  $y_{i2} \mid x_{i2}$  have identical distributions (Assumption 1 in Dang et al, 2014). We can rely on the time-invariant variables  $x_{ij}$  that are collected in both survey rounds to predict the consumptions in period 1 for households interviewed in period 2, and vice versa.<sup>7</sup> In particular, assume that  $\varepsilon_{i1}$  and  $\varepsilon_{i2}$  have a bivariate normal distribution with correlation coefficient  $\rho$  and standard deviations  $\sigma_{\varepsilon_1}$  and  $\sigma_{\varepsilon_2}$  respectively (Assumption 2').

If  $\rho$  is known, Dang et al. propose to estimate quantity (3a) by

$$P(y_{i1} < z_1 \text{ and } y_{i2} > z_2) = \Phi_2 \left( \frac{z_1 - \beta_1' x_{i2}}{\sigma_{\varepsilon_1}}, \frac{z_2 - \beta_2' x_{i2}}{\sigma_{\varepsilon_2}}, -\rho \right) \quad (4)$$

where  $\Phi_2(\cdot)$  stands for the standard bivariate normal cumulative distribution function (cdf)

(and  $\phi_2(\cdot)$  stands for the standard bivariate normal probability density function (pdf)).

---

<sup>7</sup> In other words, this assumption implies that households in period 2 that have similar characteristics to those of households in period 1 would have achieved the same consumption levels in period 1 or vice versa.

However, since  $\rho$  is usually unknown in most contexts, assume also that it is bounded by the interval  $[0, 1]$  (Assumption 2'').<sup>8</sup> Since for any  $x, y$ , and  $\rho$ ,  $\frac{\partial \Phi_2(x, y, \rho)}{\partial \rho} = \phi_2(x, y, \rho) > 0$  (Sungur, 1990), equation (4) indicates that a lower (higher) value of  $\rho$  means a higher (lower) probability of being poor in the first period but non-poor in the second period. Thus the lower bound and upper bound estimates of mobility can be established by identifying the appropriate range of values for the correlation term  $\rho$ . Absent any other information, Dang et al. (2014) suggest one can start by assuming that  $\rho$  is either 0 or 1. However, by examining empirical estimates from actual panel data for other countries, they propose a narrower range for Indonesia during 1997-2000 of  $[0.3, 0.7]$ . This in turn yields a lower bound and upper bound interval of  $[8.1, 11.8]$  for the proportion of households that were poor in 1997 but non-poor in 2000, which encompasses the “true” proportion rate of 10.1 percent based on actual panel data (Dang et al, 2014).

### III. Point Estimates on Poverty Mobility in a Generalized Framework

Despite its relevance for identifying bound estimates on poverty dynamics at the household level, the method introduced in Dang et al. (2014) suffers from several drawbacks. First, unless one is content to work with the extreme case of  $\rho$  in the  $[0, 1]$  range, a group of countries with actual panel data that is comparable to the country under investigation must be found so that a more reasonable empirical range of values for  $\rho$  can be identified. This task would require a certain degree of homogeneity for these countries, since  $\rho$  may vary depending on a host of factors caused by any difference ranging from

---

<sup>8</sup> Dang et al. (2014) provide several reasons why this assumption can be expected to hold and show this is the case using household survey data from several countries. Also see Dang et al. for more discussion on the implications for these assumptions and proofs for the bound estimates.

economic structures to modeling methods and survey designs. Complicating this issue even further, for the same country,  $\rho$  is likely to be different for different household welfare outcomes; thus an appropriate range of correlation term for, say, household food consumption must be estimated separately from that for household non-food consumption. Over time,  $\rho$  might also change.

Second, this method provides at best bound estimates, rather than point estimates. While bound estimates are certainly more useful than no estimates at all in the absence of true panel data, they also leave room for greater precision. There always exists a tradeoff between precision and encompassment with the bound approach: a larger bound interval is more likely to encompass the true rates but will be less precise, and vice versa.

In this section we generalize the method introduced in Dang et al. (2014) by presenting a method to estimate  $\rho$  based only on a country's own cross sectional surveys. We discuss the asymptotic properties of this new point estimator. We then extend the method to settings with three survey rounds or more. We maintain the Dang et al. (2014) framework described in the previous section.

### III.1. Theoretical Estimates for $\rho$

We offer the following proposition to obtain the simple correlation coefficient between household consumption in two survey rounds  $\rho_{y_{i1}y_{i2}}$ , which is closely related to  $\rho$ .

**Proposition 1- Approximate estimation of  $\rho_{y_{i1}y_{i2}}$**

*Assume household consumption follows a simple linear dynamic data-generating process given by  $y_{i2} = \alpha + \delta' y_{i1} + \eta_{i2}$  (\*), where  $\eta_{i2}$  is the random error term. Also assume that the sample size of each household survey round is large enough (or  $N \rightarrow \infty$ ), the number of cohorts ( $C$ ) constructed from the survey data is fixed, and the cohort dummy variables satisfy the relevance and exogeneity criteria for instrumental variables for  $y_{i1}$  in (\*). The simple correlation coefficient  $\rho_{y_{i1}y_{i2}}$  can then be approximated with the synthetic panel*

*cohort-level simple correlation coefficient  $\rho_{y_{c1}y_{c2}}$ , where  $c$  indexes the cohorts constructed from the household survey data.*

**Proof**

See Appendix 1.

In the absence of true panel data, we do not observe  $y_{i1}$  for the same household with household consumption in period 2, but we can predict it by projecting household consumption in period 1 on the cohort dummy variables; this process is equivalent to an instrumental variables (IV) estimation where the instrumental variables are the cohort dummy variables, and practically results in taking period-by-period sample averages within the head of household  $i$ 's cohort (Verbeek, 2008). Thus a consistent estimate for  $\rho_{y_{i1}y_{i2}}$  (and  $\delta$ ) requires that these instruments are relevant (i.e., statistically significant in the regression of household consumption on themselves) and exogenous (i.e., being uncorrelated with the error term  $\eta_{i2}$ ).

While the former assumption can be easily checked given the available data, the latter needs to be assumed and depends on our assumptions or on prior knowledge about the specific data under consideration. For example, while we expect the latter assumption holds in most contexts, it would be violated if there are cohort effects in the random error term  $\eta_{i2}$  (see, for example, Moffitt (1993) or McKenzie (2004)).<sup>9</sup> Furthermore, in addition to the relevance and exogeneity conditions, good instruments also need to be strong for unbiased estimates (Stock and Yogo (2005)), which effectively requires cohort dummy variables to

---

<sup>9</sup> On a related note, surveys that focus on a particular age group of the population (e.g., youth surveys) or with particular designs (e.g., oversampling certain age cohorts) are unlikely to provide consistent estimates for the whole population.

be strongly correlated with household consumption; fortunately, as with the relevance condition, this additional condition can be easily checked using the cross sections.<sup>10</sup>

Cohorts can be constructed from age or some combination of age and other time-invariant characteristics such as gender or ethnicity, as long as the cell size for each cohort is large enough. The assumptions stated in Proposition 1 on large cohort sizes are standard in the traditional pseudo-panel literature (e.g., type 1 asymptotics in Verbeek (2008)) and helps preclude measurement errors with cohort means. While it appears reasonable to assume that  $N$  tending to infinity given the usually large number of households interviewed with current household surveys, there seems to be no current consensus in the literature on how large  $n_c$  should be.<sup>11</sup> Thus Proposition 1 provides an approximation of  $\rho_{y_{i1}y_{i2}}$  based on asymptotic theory, and how well this approximation turns out to be in practice is an empirical issue. We will see later in our empirical estimates using household surveys that estimation results are rather encouraging. Furthermore Corollary 2.1, to be discussed below, will provide a lower value for  $\rho_{y_{i1}y_{i2}}$  as a check on our cohort estimate.

Armed with an estimate for  $\rho_{y_{i1}y_{i2}}$ , we can then proceed to propose an estimate for the partial correlation coefficient  $\rho$ , which in turn helps provide the point estimate for poverty

---

<sup>10</sup> Cohorts with a weaker correlation with household consumption would capture less variation in the latter. Clearly, the implicit assumption underlying traditional pseudo-panel analysis is that cohort dummy variables have a strong relationship with household consumption. In the extreme case, consumption (or poverty) mobility can happen entirely within cohorts, but we would expect this phenomenon to be rare in practice. Furthermore, this case would be easily detected since it will result in the synthetic panel cohort-level simple correlation coefficient  $\rho_{y_{e1}y_{e2}}$  being equal to 0 (i.e., since cohort means remain unchanged over time). We would like to thank David McKenzie for suggesting this point to us.

<sup>11</sup> Monte Carlo simulations by Verbeek and Nijman (1992) suggests that cohort sizes of 100 to 200 are sufficient, (but also see, e.g., Devereux (2007) for a related simulation). Given a typical sample size of 5000 for most household surveys, if the household head's age is used as the cohort variables, we can have approximately between 25 and 50 such cohorts of this sample size (depending on their age distribution).

mobility. We also provide an upper value on  $\rho$  as a robustness check on our estimates for this parameter.

**Proposition 2- Point estimate of  $\rho$**

Let  $R_j^2$ , for  $j= 1, 2$ , respectively represent the coefficients of determination obtained from estimating equations (1) and (2), and  $x_i$  represent the vector of household time-invariant characteristics. The partial correlation coefficient  $\rho$  can be estimated by

$$\rho = \frac{\rho_{y_{i1}y_{i2}} \sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})} - \beta_1' \text{var}(x_i) \beta_2}{\sigma_{\varepsilon_1} \sigma_{\varepsilon_2}} \quad (5)$$

**Corollary 2.1- Another approximation of  $\rho$**

If  $\beta_1 \approx \beta_2$ , the partial correlation coefficient  $\rho$  can also be estimated by

$$\rho = \frac{\rho_{y_{i1}y_{i2}} - \sqrt{R_1^2 R_2^2}}{\sqrt{1 - R_1^2} \sqrt{1 - R_2^2}} \quad (6)$$

**Corollary 2.2- Upper value of  $\rho$**

Assume that the error terms  $\varepsilon_{ij}$  in equations (1) and (2) follows the traditional household [random?] effects model and can be broken down as  $\varepsilon_{ij} = u_i + v_{ij}$  where conditional on the observed household characteristics, the unobserved household effects  $u_i$  has a normal distribution with mean 0 and variance  $\sigma_u^2$ , the idiosyncratic error terms  $v_{i1}$  and  $v_{i2}$  both have a normal distribution with means 0's and variance  $\sigma_v^2$ , and the covariance between  $v_{i1}$  and  $v_{i2}$  is 0. An upper value for the partial correlation coefficient  $\rho$  is given by the simple correlation coefficient  $\rho_{y_{i1}y_{i2}}$ .

**Corollary 2.3- Lower value of  $\rho_{y_{i1}y_{i2}}$**

The simple correlation coefficient  $\rho_{y_{i1}y_{i2}}$  for household consumption between the two survey rounds is greater than or equal to its lower value

$$i) \quad \frac{\beta_1' \text{var}(x_i) \beta_2}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} \text{ or} \quad (7)$$

$$ii) \quad \sqrt{R_1^2 R_2^2} \text{ if } \beta_1 \approx \beta_2 \quad (8)$$

with equality occurring when the estimation model fully captures all the variations in the dependent variable (i.e., all the error terms are zero).

**Proof**

See Appendix 1.

Several remarks are in order here. First, while the first way of estimating  $\rho$  given in (5) naturally follows from our framework provided by equations (1) and (2) and provides more accurate results, the second way of estimating  $\rho$  in (6) is somewhat neater and perhaps more amenable to interpretation. It suggests that, given the estimated parameters in equations (1) and (2) are close to each other, the partial correlation coefficient for household consumption can be interpreted as the simple correlation coefficient purged of (the geometric mean of) its multiple correlation with household (time-invariant) characteristics in the two survey rounds, and then reweighted by (the geometric mean of) the shares of the unexplained predicted errors. In our validation exercise to be discussed below, these two formulae give very similar estimates for  $\rho$ .<sup>12</sup>

Second, the variance-covariance matrix of the time-invariant household characteristics  $\text{var}(x_i)$  in expression (5) is the same for each round of true panel data, but can vary for the cross sectional surveys. It may thus be useful to check one's data to make sure these matrices are similar; otherwise, one may separately try the variance from each survey round to see if there is any difference in poverty estimates. In our empirical estimates discussed below, these variance-covariance matrices are very similar between survey rounds and make almost no difference to our estimate whether we use the one from the first survey round or from the second survey round.<sup>13</sup>

---

<sup>12</sup> Strictly speaking, we only require  $\beta_1\beta_2' \approx \beta_2\beta_1'$  instead of assuming  $\beta_1 \approx \beta_2$  for Corollary 2.1, but we make the above assumption for convenience. Note that for the three-variable case, the two formulae in (5) and (6) are identical. Also note that another way, still, to estimate  $\rho$  is using the recursion formula for partial correlation coefficients provided by Anderson (2003, p.41); however, this formula requires much more calculations than the given formulae above, thus we do not discuss it further.

<sup>13</sup> We abuse the notations  $\text{var}(x)$  and  $\text{var}(y)$  to refer to both the population true quantities and their sample estimates to keep the expressions simpler. Similarly, we subsequently use  $N$  to refer to both the total population and the sample survey.

Third, the assumption on the traditional household random effects in Corollary 2.2 is rather standard and should be satisfied as long as the variances of  $\varepsilon_{ij}$  are similar. As long as there are unobserved household characteristics that are not controlled for in the regression,  $\sigma_u^2$  will be positive. However, the more prediction power the model has due to inclusion of previously unobserved household time-invariant characteristics (e.g., through better data collection), the less variance these unobserved household characteristics have or the smaller  $\sigma_u^2$  is.

Fourth, we do not estimate  $\rho$  the same way we do with  $\rho_{y_1, y_2}$  as in Proposition 1, but need go through one more step with Proposition 2. The reason is straightforward once we recall that the cohort aggregation method in Proposition 1 is akin to an instrumental variable method where the cohort dummy (or age) variables work as the instruments (Moffitt, 1993). Thus, since the predicted error terms obtained from equations (1) and (2) are netted of age (and other time-invariant characteristics), when these error terms are aggregated by cohorts again, they would tend to zero.<sup>14</sup> On the other hand, we do not estimate  $\rho$  using the same procedures in Proposition 1, by, say, leaving out the age variables in estimating equations (1) and (2) since  $\rho$  obtained this way is different by construction from (and will overestimate) the partial correlation coefficient we are interested in.

Finally, it is rather straightforward to see that, given our assumption that  $\rho$  is non-negative, the numerators in equations (5) and (6) are also non-negative, thus leading to

---

<sup>14</sup> In fact, an informal check on a lower value for the partial correlation coefficient can be done by just implementing the same procedures in Proposition 1, where the predicted error terms are obtained from estimating equations (1) and (2) (including the age variable). However, as discussed above, while this estimated partial correlation can provide some value for checking purposes, it is likely to be not statistically different from 0.



Corollary 2.3.<sup>15</sup> Since the lower values of  $\rho_{y_{i1}y_{i2}}$  provided by Corollary 2.3 are derived in a different way, these represent a robustness check on our estimate based on the cohort analysis in Proposition 1 above. While these “tests” may not be powerful in the sense that they can give a tight estimate just under  $\rho_{y_{i1}y_{i2}}$ , at least they can provide some assurance that our cohort estimate of  $\rho_{y_{i1}y_{i2}}$  should satisfy a lower bound estimate and provide a positive estimate for  $\rho$ . Also a practical use of expression (8) is, given the  $R^2$ 's from two cross sections, we can use their geometric mean as a lower value to quickly gauge the strength of the simple correlation between household consumption in the two periods. A similar use of Corollary 2.1 when we know  $\rho_{y_{i1}y_{i2}}$  provides a shorthand calculation for the partial correlation  $\rho$ .

### III.2. Further Discussion on Point Estimates and Poverty Dynamics

To fully characterize the distribution for the point estimates in (4), we provide below its asymptotics.

#### **Proposition 3- Asymptotic results for point estimates for two periods**

*Assuming that household consumption can be explained by household characteristics as stated in equations (1) and (2) and all the standard regularity conditions are satisfied for each equation (i.e.,  $X'\varepsilon/N \xrightarrow{p} 0$  and  $X'X/N \xrightarrow{p} M$  finite and positive definite).<sup>16</sup> Let  $P$  represent household  $i$ 's ( $i=1, \dots, N$ ) quantity of poverty dynamics (e.g.,  $P = P(y_{i1} < z_1 \text{ and } y_{i2} > z_2)$ ),  $d_j$  an indicator function that equals 1 if the household is poor*

---

<sup>15</sup> Note that the seemingly intuitive result that the partial correlation coefficient should be less than or equal to the simple correlation coefficient (e.g., since the former results when all possible correlation with other household characteristics are removed from the latter) stated in Corollary 2.2 may only hold in our context of the multiple correlation between household consumption and household characteristics, but not in general. For example, where the  $R^2$ 's in expression (6) are not multiple correlation coefficients but just bivariate correlation coefficients, they can take on negative values that will invalidate this equality. This is the well-known suppression problem in the statistics literature; see, e.g., Friedman and Wall (2005) for a recent discussion.

<sup>16</sup> As is the usual practice, vectors of time-invariant characteristics  $x_i$ 's ( $k \times 1$ ) are transposed into row vectors and stacked on top of each other to form the matrix  $X$  ( $n \times k$ ), and the vectors of error terms  $\varepsilon$  ( $n \times 1$ ) are formed similarly from the scalars  $\varepsilon_i$ 's.

and equals -1 if the household is non-poor in period  $j$ ,  $j= 1, 2$ ,  $\rho_d = d_1 d_2 \rho$ , and  $\rho_{y_{i1}y_{i2},d} = d_1 d_2 \rho_{y_{i1}y_{i2}}$ , our point estimates are distributed as

$$\sqrt{n} \left[ P - \hat{\Phi}_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right) \right] \sim N(0, V) \quad (9)$$

The covariance-variance matrix  $V$  can be decomposed into two components, one due to sampling errors and the other due to model errors assuming these two errors are uncorrelated such that  $V = \Sigma_s + \Sigma_m$ .

The first component  $\Sigma_s$  is due to the sampling errors and can be estimated using the bootstrap method. The second component  $\Sigma_m$  is due to the model errors and can be

estimated as  $\sum_{m=1}^2 \nabla'_{\hat{\beta}_m} V(\hat{\beta}_m) \nabla_{\hat{\beta}_m} + \sum_{m=1}^2 \nabla'_{\hat{\sigma}_{\varepsilon_m}} V(\hat{\sigma}_{\varepsilon_m}) \nabla_{\hat{\sigma}_{\varepsilon_m}} + \nabla'_{\hat{\rho}_{y_{i1}y_{i2},d}} V(\hat{\rho}_{y_{i1}y_{i2},d}) \nabla_{\hat{\rho}_{y_{i1}y_{i2},d}}$

where

$$\nabla_{\hat{\beta}_m} = d_m \begin{pmatrix} -x_{ij} \\ \hat{\sigma}_{\varepsilon_m} \end{pmatrix} \phi \left( d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}} \right) \Phi \left( \frac{d_n \frac{z_n - \hat{\beta}_n' x_{ij}}{\hat{\sigma}_{\varepsilon_n}} - \hat{\rho}_d d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}}}{\sqrt{1 - \hat{\rho}_d^2}} \right) \quad (10)$$

$$- \frac{d_m d_n \text{var}(x_{ij}) \hat{\beta}_n}{\hat{\sigma}_{\varepsilon_m} \hat{\sigma}_{\varepsilon_n}} \phi_2 \left( d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}}, d_n \frac{z_n - \hat{\beta}_n' x_{ij}}{\hat{\sigma}_{\varepsilon_n}}, \hat{\rho}_d \right)$$

$$\nabla_{\hat{\sigma}_{\varepsilon_m}} = \begin{pmatrix} -d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}^2} \end{pmatrix} \phi \left( d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}} \right) \Phi \left( \frac{d_n \frac{z_n - \hat{\beta}_n' x_{ij}}{\hat{\sigma}_{\varepsilon_n}} - \hat{\rho}_d d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}}}{\sqrt{1 - \hat{\rho}_d^2}} \right) \quad (11)$$

$$- \left( d_m d_n \frac{\rho_{y_{im}y_{in}} \sqrt{\text{var}(y_{im}) \text{var}(y_{in})} - \beta_m' \text{var}(x_i) \beta_n}{\hat{\sigma}_{\varepsilon_m}^2 \hat{\sigma}_{\varepsilon_n}} \right) \phi_2 \left( d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}}, d_n \frac{z_n - \hat{\beta}_n' x_{ij}}{\hat{\sigma}_{\varepsilon_n}}, \hat{\rho}_d \right)$$

$$\nabla_{\hat{\rho}_{y_{i1}y_{i2},d}} = \frac{d_1 d_2 \sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}}{\hat{\sigma}_{\varepsilon_1} \hat{\sigma}_{\varepsilon_2}} \phi_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right) \quad (12)$$

with  $n = 3 - m$ ,  $V(\hat{\beta}_1)$  and  $V(\hat{\beta}_2)$  being respectively the estimated asymptotic covariance-variance matrix for the estimated coefficients obtained from equations (1) and (2),  $V(\hat{\sigma}_{\varepsilon_m})$

being approximated by  $\frac{(8N - 7)\hat{\sigma}_{\varepsilon_m}^2}{(4N - 3)^2}$ , and  $V(\hat{\rho}_{y_{i1}y_{i2},d})$  the estimated asymptotic variance obtained from Proposition 1.

## **Proof**

See Appendix 1.

A couple of remarks are in order. First, the model variance will be smaller the better fit we have for our regressions in equations (1) and (2). And the larger the sample sizes used for prediction, the further the sampling variance can be reduced; thus, this points to the advantages of cross sections over panel data when the former has much larger sample sizes than the latter. A natural extension of this would be to pool estimates from the two cross sections for a larger sample size to reduce the sampling errors even more.<sup>17</sup> Whether the model variance or the sampling variance is the dominant component would depend on the dynamics of the underlying regression relationship and the overall precision of our theoretical models. As shown later, our estimation results also indicate that the sampling variance is significantly larger than the model variance, which is consistent with findings in the small area estimation literature (see, e.g., Rao (2003, p. 35)). But more importantly, when the postulated estimation model is correct, it follows that the model-based variances (synthetic panel estimates in our case) are usually smaller than the design-based variances (weighted estimates based on actual panel data) (Matloff, 1981; Binder and Roberts, 2009).<sup>18</sup>

---

<sup>17</sup> Since we are mostly interested in estimating the means, assuming that the sample sizes of the cross sections are similar, we can simply use the corresponding population weight for each cross section. For estimates of other quantities such as totals, the population weights when pooling two cross sections can be adjusted by dividing by half; see, for example, Botman and Jack (1995) for a related discussion with the National Health Interview Surveys, and Kish (1999, 2002) for overviews on combining surveys.

<sup>18</sup> Notably, this decomposition of the error terms is similar in spirit to the decomposition in Elbers, Lanjouw, and Lanjouw (ELL) (2003), which is related to the familiar formula in sampling statistics that the mean squared error of the estimate is composed of the variance of the estimate plus its bias squared. However, the key difference between that study and ours is that ELL impute data from a survey into a census, but we are imputing data from a survey to a survey. Thus while ELL do not need to take into account the sampling errors, we have to model these sampling errors explicitly. An alternative to this analytical error would be the bootstrap error.

Second, the formulae in Proposition 3 are general and can be used to obtain our estimates using data either from the first or the second survey round, where the subscript  $j$  in  $x_{ij}$  should be adjusted accordingly to indicate data for the corresponding survey round. Estimates using either survey rounds are theoretically equivalent since the following identity always holds  $P(y_{i1} < z_1 \text{ and } y_{i2} > z_2) \equiv P(y_{i2} > z_2 \text{ and } y_{i1} < z_1)$ . Also note that while the number of households is the same for actual panel data across survey rounds, it can vary for cross sectional data, thus the number of observations ( $N$ ) in the variance approximation for  $V(\hat{\sigma}_{\varepsilon_m})$  should be adjusted accordingly.

Effective poverty reduction strategies require a good understanding of the proportions of the population that remains in certain poverty statuses in both periods, as well as the proportions of the population that move into or out of poverty in one period given their poverty status in the other period. Roughly speaking, the former outcomes include absolute numbers of poverty proportions (or joint probabilities) such as chronic poverty rates, while the latter outcomes include relative numbers of poverty proportions (or conditional probabilities) such as the proportion of the poor in the first period that exit poverty in the second period. The former outcomes thus provide a simultaneous view of poverty dynamics over time, but the latter outcomes emphasize its sequential nature, and both measures combined would provide a rich picture of poverty dynamics. We thus provide the following Corollary to Proposition 3 to provide the asymptotic results for such cases.

**Corollary 3.1- Asymptotic results for point estimates of relative quantities of poverty dynamics for two periods**

*Given the same assumptions in Proposition 3, let  $P_{i1}$  and  $P_{i,12}$  respectively represent household  $i$ 's ( $i=1, \dots, N$ ) quantities of poverty dynamics in period  $j$  ( $j= 1, 2$ ) and both periods (e.g.,  $P_{ij} = P(y_{ij} < z_j)$  and  $P_{i,12} = P(y_{i1} < z_1 \text{ and } y_{i2} > z_2)$ ),  $d_j$  an indicator function that equals 1 if the household is poor and equals -1 if the household is non-poor in period  $j$ ,*

and  $\rho_d = d_1 d_2 \rho$ . And let the sampled averaged estimated quantities of poverty dynamics

$$\text{represented by } \hat{\Phi}_2(\cdot) = \frac{1}{N} \sum_{i=1}^N \Phi_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right),$$

$$\hat{\Phi}(\cdot) = \frac{1}{N} \sum_{i=1}^N \hat{\Phi} \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right), \text{ our point estimates are distributed as}$$

$$\sqrt{n} \begin{bmatrix} \frac{P_{i,12}}{P_{ij}} \frac{\hat{\Phi}_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right)}{\hat{\Phi} \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right)} \end{bmatrix} \sim N(0, V_r) \quad (13)$$

where the covariance-variance matrix  $V_r$  can be estimated as

$$V_r = \left( \frac{\hat{\Phi}_2(\cdot)}{\hat{\Phi}(\cdot)} \right)^2 \left[ \frac{\text{Var}(\hat{\Phi}_2(\cdot))}{\left( \frac{\hat{\Phi}_2(\cdot)}{\hat{\Phi}(\cdot)} \right)^2} + \frac{\text{Var}(\hat{\Phi}(\cdot))}{\left( \frac{\hat{\Phi}(\cdot)}{\hat{\Phi}(\cdot)} \right)^2} - 2 \frac{\text{Cov}(\hat{\Phi}_2(\cdot), \hat{\Phi}(\cdot))}{\hat{\Phi}_2(\cdot) \hat{\Phi}(\cdot)} \right] \quad (14)$$

where  $\text{Var}(\hat{\Phi}(\cdot))$  can be decomposed into a model error  $\Sigma_{jm}$  and a sampling error  $\Sigma_{js}$ . The model error can be estimated as  $\Sigma_{jm} = \nabla'_{\hat{\beta}_j} V(\hat{\beta}_j) \nabla_{\hat{\beta}_j} + \nabla'_{\hat{\sigma}_{\varepsilon_j}} V(\hat{\sigma}_{\varepsilon_j}) \nabla_{\hat{\sigma}_{\varepsilon_j}}$  with

$$\nabla_{\hat{\beta}_j} = d_j \left( \frac{-x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right) \phi \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right) \text{ and } \nabla_{\hat{\sigma}_{\varepsilon_j}} = -d_j \left( \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}^2} \right) \phi \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right).$$

### Proof

See Appendix 1.

Note that we estimate  $\bar{\Phi}(\cdot)$  using equations (1) or (2) and its variance as discussed above, but do not use the corresponding sample-based statistics (i.e., poverty headcount ratio) to be consistent with the way we estimate  $\bar{\Phi}_2(\cdot)$ . If the model has good fit,  $\bar{\Phi}(\cdot)$  would be very similar to the sample-based poverty headcount ratio but has much smaller variance.<sup>19</sup> However, since we have to estimate both the numerators and denominators in the ratios (and their standard errors), this would reduce the accuracy of our estimates compared to those for the absolute quantities of poverty dynamics provided in Proposition 3.

---

<sup>19</sup> Another practical implication is that if we divide  $\hat{\Phi}_2(\cdot)$  by the sample poverty rate instead of  $\hat{\Phi}(\cdot)$ , this ratio can be larger than 100 percent when we consider estimates for certain subpopulation groups.

### III.3. Mobility for Three Periods or More

We now generalize this method to the general setting where there are three or even more rounds of survey data. More generally, assume that there are  $k$  rounds of survey data and household consumption levels can be explained by household characteristics for survey round  $j^{\text{th}}$  in the following equations

$$y_{ij} = \beta_j' x_{ij} + \varepsilon_{ij} \quad (15)$$

where  $j= 1, \dots, k$ . We are interested in knowing such quantities as

$$P(y_{i1} \sim z_1 \text{ and } y_{i2} \sim z_2, y_{i3} \sim z_3, \dots, y_{ik} \sim z_k) \quad (16)$$

where  $z_j$  is the poverty line in period  $j$  and the relation sign ( $\sim$ ) indicates either the larger sign ( $>$ ) or smaller sign ( $<$ ).

It is rather straightforward to see that the formula to calculate such quantities in (16) on data from the  $j^{\text{th}}$  survey round is the generalized version of (4) as follows

$$P(y_{i1} \sim z_1 \text{ and } y_{i2} \sim z_2, \dots, y_{ik} \sim z_k) = \Phi_k \left( d_1 \frac{z_1 - \beta_1' x_{ij}}{\sigma_{\varepsilon_{i1}}}, d_2 \frac{z_2 - \beta_2' x_{ij}}{\sigma_{\varepsilon_{i2}}}, \dots, d_k \frac{z_k - \beta_k' x_{ij}}{\sigma_{\varepsilon_{ik}}}, \Sigma_\rho \right) \quad (17)$$

where  $\Phi_k(\cdot)$  stands for the  $k$ -variate normal cumulative distribution function, and  $d_j$ , with  $j= 1, \dots, k$  an indicator variable that equals 1 when  $y_{ij}$  is smaller than the corresponding poverty line in the same period (i.e., household  $i$  is poor in period  $j$ ), and equals -1 otherwise.

Note that the matrix  $\Sigma_\rho$  of partial correlation coefficient is symmetric and is represented by

$$\Sigma_\rho = \begin{pmatrix} 1 & \cdot & \cdot & \cdot & \cdot & \cdot \\ \rho_{d12} & 1 & \cdot & \cdot & \cdot & \cdot \\ \rho_{d13} & \rho_{d23} & 1 & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \rho_{d1k} & \cdot & \cdot & \cdot & \cdot & 1 \end{pmatrix}_{k \times k}$$

where the subscripts  $jj$  indicates the particular two survey

rounds under consideration and all the elements on the diagonal are 1s.  $\rho_{djl}$  stands for the correlation coefficient between equation  $j$  and  $l$  and equals  $d_j * d_l * \rho_{jl}$ ; there are such  $\frac{k(k-1)}{2}$  correlation coefficients in the correlation coefficient matrix  $\Sigma_\rho$ .

However, compared to the previous case of two periods, the computation now becomes more involved since the number of integral dimensions corresponds to the number of survey rounds. Estimates will be likely to be less accurate for three periods, and longer periods in general, than those for two periods due to increased layers of (modeling and sampling) errors.

We can then generalize Proposition 3 to any setting with two periods or more. Estimates on the relative quantity of poverty dynamics can be obtained by extending the result in Corollary 3.1 to more periods, but again, note that estimates will be likely to be less accurate the more periods we consider.

**Proposition 4- Asymptotic results for point estimates for k periods**

*Assuming that household consumption can be explained by household characteristics as stated in the following equations*

$$y_{ij} = \beta_j' x_{ij} + \varepsilon_{ij} \tag{18}$$

*,  $j= 1, \dots, k$  and all the standard regularity conditions are satisfied for each equation (i.e.,  $X' \varepsilon / N \xrightarrow{p} 0$  and  $X' X / N \xrightarrow{p} M$  finite and positive definite). Let  $P$  represent household  $i$ 's ( $i=1, \dots, N$ ) quantity of poverty dynamics (e.g.,*

$P(y_{i1} \sim z_1, y_{i2} \sim z_2, y_{i3} \sim z_3, \dots, y_{ik} \sim z_k)$ ,  $d_j$  an indicator function that equals 1 if the household is poor and equals -1 if the household is non-poor in period  $j$ , and  $d_{jt} = d_j d_{1t}$ , our point estimates are distributed as

$$\sqrt{n} \left[ P - \Phi_k \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, \dots, d_k \frac{z_k - \hat{\beta}_k' x_{ij}}{\hat{\sigma}_{\varepsilon_k}}, \hat{\Sigma}_\rho \right) \right] \sim N(0, V) \quad (19)$$

The covariance-variance matrix  $V$  can be decomposed into two components, one due to sampling errors and the other due to model errors assuming these two errors are uncorrelated such that  $V = \Sigma_s + \Sigma_m$ .

The first component  $\Sigma_s$  is due to the sampling errors and can be estimated using the bootstrap method.

To make notations less cluttered, let  $\beta_{(jxl)}$  represent the matrix of estimated coefficients obtained from (18),  $\Phi_k(\cdot)$  the standard  $k$ -variate normal probability, and

$$\hat{a}_{dmj} = d_m \frac{z_m - \hat{\beta}_m' x_{ij}}{\hat{\sigma}_{\varepsilon_m}} \text{ and } \bar{a}_{dmnj} = \frac{(\hat{\rho}_{dmq} - \hat{\rho}_{dnq} \hat{\rho}_{dmn}) \hat{a}_{dmj} + (\hat{\rho}_{dnq} - \hat{\rho}_{dmq} \hat{\rho}_{dmn}) \hat{a}_{dnj}}{\sqrt{1 - \hat{\rho}_{dmn}^2}} \text{ for } m, n, q =$$

$1, \dots, k$ , and  $m \neq n \neq q$ . Also let  $\hat{\Sigma}_{\rho(-m)}$  be the  $(k-1) \times (k-1)$  partial correlation matrix given  $\hat{\beta}_m$

with the off-diagonal entries  $\hat{\rho}_{dst.m} = \frac{\hat{\rho}_{dst} - \hat{\rho}_{dsm} \hat{\rho}_{dtm}}{\sqrt{1 - \hat{\rho}_{dsm}^2} \sqrt{1 - \hat{\rho}_{dtm}^2}}$  for  $s, t = 1, \dots, k$  and  $s, t \neq m$ ;

similarly, let  $\hat{\Sigma}_{\rho(-m,-n)}$  be the  $(k-2) \times (k-2)$  partial correlation matrix given  $\hat{\beta}_m$  and  $\hat{\beta}_n$  with the

off-diagonal entries  $\hat{\rho}_{dst.mn} = \frac{\hat{\rho}_{dst.m} - \hat{\rho}_{dsn.m} \hat{\rho}_{dtn.m}}{\sqrt{1 - \hat{\rho}_{dsn.m}^2} \sqrt{1 - \hat{\rho}_{dtn.m}^2}}$  for  $s, t = 1, \dots, k$  and  $s, t \neq m, n$ . The

second component  $\Sigma_m$  is due to the model errors and can be estimated as

$$\sum_{m=1}^k \nabla'_{\hat{\beta}_m} V(\hat{\beta}_m) \nabla_{\hat{\beta}_m} + \sum_{m=1}^k \nabla'_{\hat{\sigma}_{\varepsilon_m}} V(\hat{\sigma}_{\varepsilon_m}) \nabla_{\hat{\sigma}_{\varepsilon_m}} + \sum_{m=1}^{k-1} \sum_{n=m+1}^k \nabla'_{\hat{\rho}_{y_{im}y_{in},d}} V(\hat{\rho}_{y_{im}y_{in},d}) \nabla_{\hat{\rho}_{y_{im}y_{in},d}}$$

where

$$\nabla_{\hat{\beta}_m} = d_m \left( \frac{-x_{ij}}{\hat{\sigma}_{\varepsilon_m}} \right) \phi(\hat{a}_{dmj}) \Phi_{k-1} \left( \frac{\hat{a}_{d1j} - \hat{\rho}_{dm1} \hat{a}_{dmj}}{\sqrt{1 - \hat{\rho}_{dm1}^2}}, \dots, \frac{\hat{a}_{dkj} - \hat{\rho}_{dmk} \hat{a}_{dmj}}{\sqrt{1 - \hat{\rho}_{dmk}^2}}, \hat{\Sigma}_{\rho(-m)} \right) + \quad (20)$$

$$+ \sum_{\substack{n=1 \\ n \neq m}}^k d_m d_n \frac{-\text{var}(x_{ij}) \hat{\beta}_n}{\hat{\sigma}_{\varepsilon_m} \hat{\sigma}_{\varepsilon_n}} \phi_2(\hat{a}_{dmj}, \hat{a}_{dnj}, \hat{\rho}_{dmn}) \Phi_{k-2}(\hat{a}_{d1j} - \bar{a}_{dm1j}, \dots, \hat{a}_{dkj} - \bar{a}_{dmkj}, \hat{\Sigma}_{\rho(-m,-n)})$$

$$\nabla_{\hat{\sigma}_{\varepsilon_m}} = \left( \frac{-\hat{a}_{dmj}}{\hat{\sigma}_{\varepsilon_m}} \right) \phi(\hat{a}_{dmj}) \Phi_{k-1} \left( \frac{\hat{a}_{d1j} - \hat{\rho}_{dm1} \hat{a}_{dmj}}{\sqrt{1 - \hat{\rho}_{dm1}^2}}, \dots, \frac{\hat{a}_{dkj} - \hat{\rho}_{dmk} \hat{a}_{dmj}}{\sqrt{1 - \hat{\rho}_{dmk}^2}}, \hat{\Sigma}_{\rho(-m)} \right) - \quad (21)$$

$$- \sum_{\substack{n=1 \\ n \neq m}}^k \left( d_m d_n \frac{\rho_{y_{im}y_{in}} \sqrt{\text{var}(y_{im}) \text{var}(y_{in})} - \beta_m' \text{var}(x_i) \beta_n}{\hat{\sigma}_{\varepsilon_m} \hat{\sigma}_{\varepsilon_n}} \right) \phi_2(\hat{a}_{dmj}, \hat{a}_{dnj}, \hat{\rho}_{dmn})^*$$

$$* \Phi_{k-2}(\hat{a}_{d1j} - \bar{a}_{dm1j}, \dots, \hat{a}_{dkj} - \bar{a}_{dmkj}, \hat{\Sigma}_{\rho(-m,-n)})$$



$$\nabla_{\hat{\rho}_{y_{im}y_{in},d}} = \frac{d_m d_n}{\sqrt{1-R_m^2} \sqrt{1-R_n^2}} * \phi_2(\hat{a}_{dmj}, \hat{a}_{dnj}, \hat{\rho}_{dmn}) * \Phi_{k-2}(\hat{a}_{d1j} - \bar{a}_{dm1j}, \dots, \hat{a}_{dkj} - \bar{a}_{dmkj}, \hat{\Sigma}_{\rho(-m,-n)}) \quad (22)$$

with  $V(\hat{\beta}_m)$  being the asymptotic covariance-variance matrix for the estimated coefficients obtained from the corresponding equation in (18) and  $V(\hat{\sigma}_{\varepsilon_m})$  being approximated by  $\frac{(8N-7)\hat{\sigma}_{\varepsilon_m}^2}{(4N-3)^2}$ .

### III.4. Mobility between Different Consumption Groups

The results in Proposition 3 can be straightforwardly extended to settings with more than two consumption groups as follows.

#### Proposition 5- Asymptotic results for point estimates for mobility between different groups for two periods

Assuming that household consumption can be explained by household characteristics as stated in equations (1) and (2) and all the standard regularity conditions are satisfied for each equation (i.e.,  $X'\varepsilon/N \xrightarrow{p} 0$  and  $X'X/N \xrightarrow{p} M$  finite and positive definite). Let  $P^{lm}$  represent household  $i$ 's ( $i=1, \dots, N$ ) probability of moving from consumption group  $l$  in period 1 to consumption group  $m$  in period 2, that is  $P^{lm} = P(z_1^{l-1} < y_{i1} < z_1^l \text{ and } z_2^{m-1} < y_{i2} < z_2^m)$ , where  $l, m = 1, \dots, k$ , and the  $z_j$  are the thresholds that separate the different consumption groups, with  $z_j^0 = -\infty$  and  $z_j^k = \infty$ , for

period  $j, j = 1, 2$ . Defining  $F^{l,m}$  as  $\Phi_2\left(\frac{z_1^l - \beta_1'x_{ij}}{\sigma_{\varepsilon_1}}, \frac{z_2^m - \beta_2'x_{ij}}{\sigma_{\varepsilon_2}}, \rho\right)$ , our point estimates are

distributed as

$$\sqrt{n} \left[ P^{lm} - (\hat{F}^{l,m} - \hat{F}^{l,(m-1)} - \hat{F}^{(l-1),m} + \hat{F}^{(l-1),(m-1)}) \right] \sim N(0, V) \quad (23)$$

Given  $k$  consumption groups in each period, there are  $k \times k$  transitions in total. The formulae for the standard errors for the general case can be much more than those for mobility for three periods or more (and are at least as complicated). Thus we suggest estimation of the standard errors by the bootstrap method.<sup>20</sup>

<sup>20</sup> Also note that our empirical estimates, discussed later, point to little, if any difference between the standard errors estimated using the analytical formulae offered in Proposition 3 and those using the bootstrap approach.

### III.5. Estimation Procedures

A practical concern not yet discussed is whether or not equations (1) and (2) should be estimated with household weights. There appear to be both advantages and disadvantages with both approaches. Weighted regressions are especially relevant when the provided household weights were constructed to account for non-response or attrition bias or specifically based on the dependent variables (informative sampling); on the other hand, unweighted regressions are most relevant when the proposed super-population (i.e., equations (1) and (2)) model is correct and can provide some causal interpretation. Estimation without weights in the former case results in biased estimates, while estimation with weights in the latter case yields inefficiency (i.e., larger standard errors).<sup>21</sup> Thus it seems advisable to estimate models both with and without weights and compare results, particularly where there is limited information on how the weights have been constructed.

Given the framework discussed above, we propose the following steps to obtain poverty mobility for two periods:

*Step 1:* Using the data in survey round 1, estimate equation (1) and obtain the predicted coefficients  $\hat{\beta}_1'$ , and the predicted standard error  $\hat{\sigma}_{\varepsilon_{i1}}$  for the error term  $\varepsilon_{i1}$ . Using the data in survey round 2, estimate equation (2) and obtain similar parameters  $\hat{\beta}_2'$  and  $\hat{\sigma}_{\varepsilon_{i2}}$ .

*Step 2:* Aggregate data in both survey rounds 1 and 2 by cohorts and obtain the estimated cohort-level simple correlation coefficient  $\hat{\rho}_{y_{i1}y_{i2}}$ . Calculate  $\hat{\rho}$  using Proposition 2, and check

that  $\hat{\rho}_{y_{i1}y_{i2}} \geq \hat{\rho}$  (and also  $\hat{\rho}_{y_{i1}y_{i2}} \geq \frac{\hat{\beta}_1' \text{var}(x_i) \hat{\beta}_2}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}}$ ).

*Step 3:* For each household in survey round  $j$ , calculate absolute quantities of poverty mobility as  $\Phi_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right)$ , where  $d_j$  is an indicator function that equals

---

<sup>21</sup> See also Deaton (1997), Lohr (2010), and Pfeffermann (2011) for further discussion on this topic.

1 if the household is poor and equals -1 if the household is non-poor in period  $j$ ,  $j= 1, 2$ , and  $\hat{\rho}_d = d_1 d_2 \hat{\rho}$ . Calculate the standard errors using Proposition 3. Make the appropriate adjustments to obtain population-level numbers.

*Step 4:* (If relevant) Calculate the population-level relative quantities of poverty mobility for

period  $j$  as  $\frac{\hat{\Phi}_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_1}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{ij}}{\hat{\sigma}_{\varepsilon_2}}, \hat{\rho}_d \right)}{\hat{P}_j}$ , where  $d_j$  is an indicator function that equals

1 if the household is poor and equals -1 if the household is non-poor in period  $j$ ,  $j= 1, 2$ , and  $\hat{\rho}_d = d_1 d_2 \hat{\rho}$ . Calculate the standard errors using Corollary 3.1.

The estimation procedures for three periods or more are similar, with poverty mobility rates and standard errors estimated in Steps 3 and 4 using Proposition 4 instead of Proposition 3.<sup>22</sup>

#### IV. Monte Carlo Simulation

We validate our method based on simulated data in this section. Assuming household  $i$ 's consumption in periods 1 and 2 can be generated using the following model

$$y_{i1} = \alpha_1 + \beta_{11}x_1 + \beta_{12}x_2 + \beta_{13}x_3 + \beta_{14}x_4 + \beta_{15}x_5 + \beta_{16}x_6 + \beta_{17}x_7 + \beta_{18}x_8 + v_{i1} \quad (24)$$

$$y_{i2} = \alpha_2 + \beta_{21}x_1 + \beta_{22}x_2 + \beta_{23}x_3 + \beta_{24}x_4 + \beta_{25}x_5 + \beta_{26}x_6 + \beta_{27}x_7 + \beta_{28}x_8 + v_{i2} \quad (25)$$

where the  $x$ 's are household head's time-invariant characteristics, and  $v$ 's the random error terms. Since we may not be able to fully observe all household time-invariant characteristics in a survey, this model allows us to simulate situations ranging from one extreme with minimal data collected on the household (say, we only observe  $x_1$ ) to the other extreme with fully observed household characteristics (i.e., we can observe all the  $x$ 's). When the  $x$ 's are unobserved, they would, together with the random error term  $v_j$ , be absorbed in the (total) random error term  $\varepsilon_j$ .

<sup>22</sup> We are working on a Stata program to calculate poverty mobility based on synthetic panels and will make it publicly available soon.

We then assume the following parameter values

$$\alpha_1 = 1, \alpha_2 = 1.5$$

$$\beta_{11} = \beta_{12} = \beta_{13} = \beta_{14} = \beta_{15} = \beta_{16} = \beta_{17} = \beta_{18} = 1$$

$$\beta_{21} = 1.2, \beta_{22} = 1.1, \beta_{23} = 1.05, \beta_{24} = 1.3, \beta_{25} = 0.9, \beta_{26} = 1.15, \beta_{27} = 1.4, \beta_{28} = 0.6$$

and

$$x_1 \sim N(0, 2.5), x_2 \sim N(0, 5), x_3 \sim N(0, 6), x_4 \sim N(0, 4), x_5 \sim N(0, 1), x_6 \sim N(0, 3),$$

$$x_7 \sim N(0, 2), x_8 \sim N(0, 1)$$

$$\begin{pmatrix} v_{i1} \\ v_{i2} \end{pmatrix} \sim BVN\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 6.5 & 1 \\ 1 & 6.5 \end{pmatrix}\right)$$

where  $N(0, c)$  stands for the normal distribution with mean 0 and variance  $c$ ;  $BVN(\cdot, \cdot)$  similarly represents the bivariate normal distribution with the vector of mean 0 and the given variance-covariance matrix.

Given these parameter values, we can calculate that  $\text{var}(y_{i1}) = 27$ ,  $\text{var}(y_{i2}) = 38.6$ ,  $\rho_{y_{i1}, y_{i2}} = 0.89$  as well as a range of values for the partial correlation coefficient  $\rho$  corresponding to the number of time-variant characteristics we can observe in the model. Table 1 provides the values for  $\rho$  for different data situations. These range from very limited information on the time-invariant variables where we only observe  $x_1$  (i.e.,  $\rho = 0.88$  and is almost identical to  $\rho_{y_{i1}, y_{i2}}$ ) to a typical setting with just a few such variables (i.e.,  $\rho = 0.58$ ), and to an unusual setting where we fully observe all the  $x$ 's (i.e.,  $\rho = 0.15$ ). Note that without loss of generality, we assume a certain degree of correlation over time for the error terms  $v$ 's ( $\rho = 0.15$ ) which may represent the correlation due to factors other than household time-invariant characteristics such as time-varying unexpected shocks. In data situations where we only observe some but not all of the time-invariant  $x$ 's, the unobserved  $x$ 's end up in the random error terms and contribute to the correlation  $\rho$  of these error terms over the two periods. In such situations, the contribution of the time-invariant component is

indistinguishable from that of the time-varying component, and it is the sum total of these two components that matters for our simulation purposes (as well as in most practical settings).

To save space, we provide simulation results only for Models 1, 5, and 8 at three different sample sizes  $N= 1000, 4000, \text{ and } 10000$ , with 1000 simulations for each model run. These sample sizes range from a small panel dataset to a decent-sized dataset and an unusually large one. We fix the poverty line in period 2 at the 30<sup>th</sup> percentile, and then graph in Figure 1 the true percentage of households that are poor in both periods (solid line) and its 95 percent confidence intervals (shaded bands) and the estimated percentage using simulated data (dashed line) against the whole spectrum of poverty rates in the first period.

Figure 1 shows that estimated poverty rates using simulated data closely track the true rates and fall within their 95 percent confidence intervals, but the more time-invariant variables that are available, the better prediction we have. In particular, the dashed lines are almost indistinguishable from the solid line for the graphs where  $\rho = 0.15$  or even where  $\rho = 0.58$ . When there is very limited information on these time-invariant variables ( $\rho = 0.88$ ) and the sample sizes for the actual panel data are mid-sized or unusually large ( $N= 4000$  or  $10000$ ), estimates (partially) fall outside the 95 percent confidence interval for the middle part of the distribution. But note that even in this case, if the sample size is not large enough ( $N= 1000$ ), the synthetic panel estimates still compares favorably well to true poverty rates (which are less accurately estimated given the small sample size). Varying the model parameters or the poverty lines gives us similar results (not shown).

Overall, these simulation results are very encouraging and indicate that our method performs rather well under the theoretical setting where our assumption about the bivariate

normality distribution of the error terms holds. We will examine in the next sections how our method performs in the real-life settings where this assumption may be violated.

## **V. Data**

To validate our method with real survey data, we analyze household panel survey data from Bosnia-Herzegovina (Bosnia-Herzegovina Living Standards Measurement Survey, BLSMS), Lao PDR (Expenditure and Consumption Survey, LECS), the United States (Panel Study of Income Dynamics, PSID), Peru (Peruvian National Household Survey, ENAHO), and Vietnam (Vietnam Household Living Standards Survey, VHLSS). We use two rounds from the first two surveys and three rounds from the last three surveys, with data from the BLSMS in 2001-2004,<sup>23</sup> the LECS in 2002/03-2007/08, the PSIDs in 2005, 2007, and 2009, the ENAHOs in 2004, 2005, and 2006, and the VHLSSs in 2004, 2006, and 2008. The number of households hovers around 2,376 households for Bosnia-Herzegovina, 6,500 households for the LECS, 9,189 households for each round of the VHLSSs, more than 5,000 households for the PSIDs,<sup>24</sup> and almost 20,000 households for the ENAHOs.

Except for the PSID that is implemented by the University of Michigan, all the other surveys are nationally representative surveys implemented by each country's statistical agencies, with previous or current technical assistance from international organizations (the World Bank with Peru and Vietnam), leading universities (University of Essex with Bosnia-Herzegovina) or statistical agencies in richer countries (Statistics Sweden with Lao PDR). Also except for the PSID, all the other surveys are similar to the LSMS-type (Living Standards Measurement Survey) surveys supported by the World Bank in a number of developing countries and provide detailed information on household consumption and

---

<sup>23</sup> We build this dataset up on the basis of data from Demirguc-Kunt, Klapper and Panos (2011).

<sup>24</sup> We only consider the sample persons in the PSID with non-zero longitudinal weights.

demographics, as well as schooling, health, employment, migration, and housing. The PSID has a more complex structure and provides similarly detailed, if not richer, information. All these surveys are widely used in academic studies (especially the PSID) as well as poverty assessment by the government and the donor community. We use the official poverty lines for Lao PDR, Peru, and Vietnam; for the USA, we use the poverty lines provided in the PSID data (which adjust for family size and demographics); for Bosnia-Herzegovina we use the 20<sup>th</sup> percentile of the consumption distribution in 2001 as the poverty line.

One particular feature the LECSs, VHLSSs and ENAHOs share is a rotating panel design, which collects panel data for a subset of each survey round between two adjacent years. Around one third and one half of the households in the first round are repeated in the next round for the LECs and VHLSSs respectively, and the corresponding repetition ratio for the ENAHOs is around one quarter. This combination of both cross-sectional data and panel data in one survey provides an appropriate setting for us to implement our procedures on the cross section components, and then validate our estimates against the true rates from the panel components for each country. For the BLSMSs and the PSIDs, there is no rotating panel design and so we use the panel halves, pretending that these are cross sectional data. To ensure comparability between estimates based on the panel and cross section components, we use household weights with our estimates for the ENAHOs and population weights for the remaining surveys.<sup>25</sup> We use income and household consumption as a household welfare measure respectively for the US and all the other countries.<sup>26</sup>

---

<sup>25</sup> There can be both pros and cons with using true panel data versus a survey with a rotating panel design for validation purposes. On one hand, a rotating panel design may be more suitable since actual panel surveys usually have a smaller sample size than those of cross sections, and a reasonably large sample size is required to obtain accurate estimates for the cohort-level simple correlation coefficient as well as for our asymptotic results. On the other hand, an important requirement for using rotating panel surveys is that the data from the cross section component be similar to those from the panel component. For Peru, the household-weighted headcount poverty rates based on the actual panel component are around 5 percent lower than those based on

Appendix 3 provides a more detailed description of these surveys and other data quality checks.

## VI. Estimates on Poverty Dynamics

### VI.1. Approximation for $\rho$

Consistent with the literature on pseudo-panel data, we restrict household heads' age range to 25-55 for the first survey round and adjust this appropriately for later survey rounds (e.g., looking at the age cohort 27-57 if the next survey round is two years later). While this age range can be extended to include older people, it may be ill-advised to include those who are younger, at least since most household heads tend to be older than 25 in all the countries we look at.

After obtaining an estimate for  $\rho_{y_{i1}y_{i2}}$  from the synthetic panels (which are all highly statistically significant with p-values less than 0.01) based on age cohorts, we calculate  $\rho$  in two different ways using Proposition 2 and Corollary 2.1. Estimation results are very similar with the differences being at most 0.01; thus, we only show the estimates based on Proposition 2 in Table 2. Comparing estimates using the actual panels and the synthetic panels, the absolute difference for  $\hat{\rho}$  ranges from 0.01 (Bosnia-Herzegovina) to 0.18 (the US during 2005-2007), which corresponds to a range of relative differences of 2 to 28 percent.

Nevertheless, it is reassuring to see that estimates for  $\rho$  are always less than those for  $\rho_{y_{i1}y_{i2}}$ , which is consistent with the hypothesis in Corollary 2.2; similarly, estimates for  $\rho_{y_{i1}y_{i2}}$  are also larger than its lower value estimates as can be estimated in Corollary 2.3(ii)

---

the cross section component, and the population-weighted estimates are even more different. Thus while the Peruvian data are not perfect for validation purposes, we believe it is still useful to show estimates for this country using household weights. But note that this difference will introduce more noise into our estimates.

<sup>26</sup> The PSID also has some information on household consumption but this measure is not commonly used to measure poverty and is much less comprehensive than those for other surveys.



(not shown). These results are very encouraging and suggest that our framework could well be applied to these cross sections to obtain an approximation for  $\rho$  in the absence of actual panel data.

## **VI.2. Overall Poverty Mobility**

To save space, we show estimates for poverty dynamics using the latest survey rounds available, that is for Bosnia-Herzegovina during 2001-2004, Lao PDR during 2002/03-2007/08, Peru during 2005- 2006, the US during 2007-2009, and Vietnam during 2006-2008 in Table 3.<sup>27</sup> Estimates for earlier survey rounds for the last three countries provided in Appendix 2, Table 2.2 offer qualitatively similar results.

We provide two ways to evaluate the goodness-of-fit for estimation results. First, we consider the precision of the synthetic panel point estimates by enumerating the number of times they fall within the 95 percent confidence interval around those from the actual panels. A more demanding test is to consider a similar statistics on a narrower band of one standard error around the actual panel estimates. Second, we consider the efficiency of the synthetic panel estimates by looking at the share of the overlap between the 95 percent confidence intervals of the synthetic panel estimates and the true estimates over the former's 95 percent confidence interval. This test focuses on the proportional contribution of its variance in the total mean squared error and scores the accuracy of each estimate on a 100 percent scale.<sup>28</sup>

Our estimation results using actual panel data and synthetic panel data are respectively displayed under the columns “Actual Panel” and “Synthetic Panel” in Table 3, and the

---

<sup>27</sup> Unless otherwise noted, we use data in the second survey round ( $x_{i2}$ )—as the base year—for predictions. Estimates corresponding to those in Table 2 but use data in the first survey round are provided in Appendix 2, Table 2.3. Weighted regressions provide qualitatively similar but somewhat less accurate results, thus we use unweighted regressions.

<sup>28</sup> See also Dorfman (2011) for a related test for the bias of the point estimate in the context of small area estimation.

goodness-of-fit results are shown in Table 4.<sup>29</sup> Results appear very encouraging with the synthetic panel point estimates being close to the true point estimates and lying within the 95 percent confidence intervals around the true estimates for almost all the cases (i.e., 19 out of 20). In fact, more than half of the synthetic panel point estimates fall within one standard error of the actual panel estimates (i.e., 11 out of 20; not shown). Furthermore, the difference between the true rates and our estimates seems negligible for certain mobility categories; for example, the percentage of those remaining poor in both periods for Peru during 2005-2006 were 29.9 percent and 30.9 percent respectively for the true rates and our estimates.

For the second test, Table 4 shows that all the synthetic panel estimates have half or more of their 95 percent confidence intervals being contained inside those around the true estimates. If we only consider the cases where the former should be completely contained inside the latter, more than four fifths (i.e., 17 out of 20) of the synthetic panel estimates pass this higher bar of 100 percent inclusion.

As discussed above, the standard errors for the synthetic panel estimates consist of two components, the model errors and the sampling errors, with the latter's variance expected to be larger than the former's variance when the regressions have good fits.<sup>30</sup> This is indeed the case where (results not shown) the variances of the sampling errors are significantly larger than those for the model errors. Thus, since the sampling errors account for most of the errors with the synthetic estimates and the cross sections used for the synthetic estimates have larger sample sizes than panel data, the synthetic estimates unsurprisingly have smaller

---

<sup>29</sup> Estimated parameters for equations (1) and (2) are provided in Appendix 2, Table 2.1.

<sup>30</sup> We also calculate the bootstrap standard errors by bootstrapping  $(y_{ij}, x_{ij})$  from its empirical distribution function (1,000 times) and applying the estimated parameters for equations (1) and (2) from the original samples. Estimation results are very similar to the analytical standard errors.

standard errors than those based on actual panel data. Table 3 shows that, for the US and Bosnia-Herzegovina where the sample size is the same for both actual panel and synthetic panel estimates, the standard errors for the latter are smaller than those for the former.<sup>31</sup>

Table 5 provides the proportions of the population that exit or fall into poverty in the second period given their poverty status in the first period, using the results derived in Corollary 3.1. The goodness-fit-of tests are now shown at the bottom of this table to save space. Estimation results are, not surprisingly, slightly less accurate than those in Table 3 since both the numerators and denominators in the ratios in Corollary 3.1 are estimated. Due to one additional layer of estimates, now just more than four fifths (i.e., 18 out of 20) and one half (i.e., 10 out of 20) of the synthetic estimates are respectively within the 95 percent confidence intervals and one standard error of the true rates. Interestingly, all the estimates for Lao PDR fall within one standard error and the 95 percent confidence intervals of the true rates.

Similarly with the efficiency test, now almost all (i.e., 19 out of 20) and four fifths (i.e., 16 out of 20) of the synthetic panel estimates pass the 50 percent mark and 100 percent mark respectively.

### **VI.3. Poverty Mobility for Population Sub-Groups**

It is important to investigate poverty dynamics for population sub-groups for at least two reasons. First, policy makers are usually interested in focusing on smaller population groups rather than the whole population in designing social safety net programs; and second, synthetic panels usually have larger sample sizes than actual panel data, and thus the larger sample sizes the former has, the more accurate estimates it can bring.

---

<sup>31</sup> More generally, this is consistent with the well-known result in survey sampling that the model-based variances (synthetic panel estimates in our case) are usually smaller than the design-based variances (weighted estimates based on true panel data). See, e.g, Binder and Roberts (2009) for a recent review on this topic.

We estimate and plot the estimated rates with their 95 percent confidence intervals for the absolute and relative measures of poverty dynamics against the true rates for the population categorized by ethnicity (i.e., ethnic minority groups), gender of household heads (i.e., female-headed households), education achievement (i.e., primary education or higher, lower secondary education or higher), and residence areas (i.e., urban households or regions the household live in) respectively for Peru in Figures 2 to 5 and Vietnam in Figures 6 to 9.<sup>32</sup> Clearly, these categorizations can overlap but they can provide a first cut at profiling poverty mobility for different groups, and we would expect an overlap between the synthetic panel estimates and the true rates for the groups whose heterogeneity mimics that of the whole country.

Not surprisingly, the 95 percent confidence intervals for synthetic panels estimates for both Peru and Vietnam are much smaller than those for the true rates with the gaps between the standard errors amplified roughly twice (i.e., multiplied by 1.96). Our estimates appear to be reasonably good, especially for those who are poor in either period or both periods. Except for a few cases (e.g., households where heads only achieve primary education or living in urban Selva in Figure 2, or households living in urban areas in Figure 7), there is much overlap between the true rates and our estimated rates.

#### **VI.4. Poverty Mobility in Three Periods**

We turn next in Table 6 to examining our estimates on poverty mobility for three periods using data from all three survey rounds for the US in 2005-2007, and 2009, Vietnam in 2004, 2006, and 2008, and Peru in 2004, 2005, and 2006, where there are 8 possible poverty categories that each household can fall in in these three periods (for the

---

<sup>32</sup> We do not show all eight graphs for one country or for other countries to save space, but results are qualitatively similar.

unconditional probabilities).<sup>33</sup> As discussed above, we should expect estimates to be less accurate than those for two periods; however, our proposed method turns out to work quite well with more than two thirds (i.e., 17 out of 24) of all the point estimates being contained in the 95 percent confidence intervals around the true rates; the corresponding figure for the stricter test of one standard error is just one half (i.e., 12 out of 24). The efficiency test points to somewhat better results with more than two thirds (i.e., 17 out of 24) and almost two thirds (i.e., 15 out of 24) of the point estimates respectively pass the 50 percent mark and 100 percent mark.

#### **VI.5. Consumption Dynamics between Different Groups**

A common way to study consumption dynamics is to divide the consumption distribution into five groups (quintiles), and look at the 5x5 transition matrix for mobility between two periods. We provide such estimates using data from Vietnam in 2006-2008 in Table 7, where the synthetic panel estimates are shown in panel B and the actual panel estimates shown in panel A. While estimates are off with some of the row and column totals, it is perhaps more useful to focus on the inner transitions since the former do not offer as much insight into mobility as the latter.<sup>34</sup> Estimation results are, again, rather encouraging with the majority of the different consumption categories (i.e., four-fifths of the inner transitions) falling within the 95 percent CIs of the true estimates, which are presented in

---

<sup>33</sup> Estimation results for the trivariate normal probabilities in this table are calculated using the Stata algorithm by Cappellari and Jenkins (2006) with 100 Halton draws.

<sup>34</sup> The row or column totals should sum up to 20 percent by definition and serve mostly as an indicator of prediction accuracy for these totals only. In addition, it may be useful to highlight the fact that our validation is predicated on the assumption that the true panel data for Vietnam have good quality. If the mobility in the true panel data is partly caused by spurious changes due to measurement errors (or attrition bias) in household consumption, our estimates based on the synthetic panel data would be more accurate since cross sections are free of such data issues. See also Dang and Lanjouw (2014) for a related study of mobility between three consumption groups: poor, vulnerable, and middle class.

bold. These estimates also pass the 100 percent mark of the coverage test.<sup>35</sup> Furthermore, some of the remaining estimates that fall just outside these 95 percent CIs around the true estimates appear practically close to the latter (e.g., the transition from quintile 3 to quintile 4 or from the richest quintile to quintile 2).

## **VII. Conclusion**

For effective poverty reduction policies to take place, we need to understand well poverty mobility over time. In the absence of panel data, it has historically been difficult to study poverty mobility in developing countries. However, there are by now at least two rounds of cross-sectional household survey data available in a large majority of developing countries. Our proposed method, which generalizes that initiated by Dang et al. (2014), offers a means to convert these cross-sectional survey data into synthetic panel data.

In particular, by moving away from the bound estimates in Dang et al. (2014) to point estimates, we show that our estimates are quite accurate and do not depend on additional information from ancillary data. Our method would thus seem applicable in most settings where two cross sections are available. We find that estimation results are good not only for the general population but for smaller population groups as well, and are associated with much tighter confidence intervals than even direct, panel-data based estimates in those settings where the sample sizes for the cross sections are large enough. Our estimates are validated against both simulated data and actual panel data spanning different income levels and geographical regions. In addition, our method can be readily extended to settings with three survey rounds or more, although predictions are more accurate for shorter periods.

---

<sup>35</sup> Other useful statistics that can be calculated from Table 7, panel B, include the percentages of the population that have seen either an improvement or a decline or remained in the same quintile over time, which are respectively 24.7 percent, 27.3 percent, and 48 percent. These estimates are within the 95 CIs around those based on the actual panels.

It should be noted that our method needs not be restricted only to the analysis of poverty mobility. It can in principle also be employed to study dynamics more generally. Possible application can involve dynamics in the areas of labor (e.g., what is the percentage of those who move from the rural farm sector to the non-farm sector?), health (e.g., what is the percentage of those who have medical insurance in the first period but do not in the second period?) or finance (e.g., what is the percentage of borrowers who do not default on loans in the first period but do in the second period?).

We come away with a growing sense that this basic methodology offers significant potential towards a better understanding of poverty dynamics in settings where panel data are absent, and can serve as a rather promising avenue for further research. For example, since measurement errors may be larger with richer households' consumption, future research may investigate the heterogeneity of the error terms for different types of households. Another potential direction is to help better target social transfers by identifying a "vulnerability line" that is higher than the poverty line but households with a consumption level lower than this line are still subject to a considerable probability of falling back into poverty in the next period (Dang and Lanjouw, 2014). Another is to examine the extent that the estimates offered by our synthetic panels can improve on and correct for "bad" panel data estimates resulting from serious attrition problems.

## References

- Agostini, Claudio A. and Philip H. Brown. (2010). "Local Distributional Effects of Government Cash Transfers in Chile". *Review of Income and Wealth*, 56(2): 366- 388.
- Alkire, Sabina and James E. Foster. (2011). "Counting and Multidimensional Poverty Measurement." *Journal of Public Economics*, 95: 476-487.
- Anderson, Theodore W. (2003). "*An Introduction to Multivariate Statistical Analysis*". USA: John Wiley & Sons.
- Banks, James, Richard Blundell, and Agar Brugiavini. (2001). "Risk Pooling, Precautionary Saving and Consumption Growth". *Review of Economic Studies*, 68(4): 757-779.
- Bierbaum, Mira and Franziska Gassmann. (2012). "Chronic and Transitory Poverty in the Kyrgyz Republic: What Can Synthetic Panels Tell Us?" *UNU-MERIT Working Paper #2012-064*.
- Binder, David A. and Georgia Roberts. (2009). "Design- and Model-Based Inference for Model Parameters". In D. Pfeiffermann and C.R. Rao. *Handbook of Statistics, Vol. 29B-Sample Surveys: Inference and Analysis*. North-Holland: Elsevier.
- Blundell, Richard, Alan Duncan, and Costas Meghir. (1988). "Estimating Labor Supply Responses Using Tax Reforms". *Econometrica*, 66(4): 827- 861.
- Botman, Steven L. and Susan S. Jack. (1995). "Combining National Health Interview Survey Datasets: Issues and Approaches". *Statistics in Medicine*, 14: 669-677.
- Bourguignon, Francois, Chor-Ching Goh, and Dae Il Kim. (2004). "Estimating Individual Vulnerability to Poverty with Pseudo-Panel Data", *World Bank Policy Research Working Paper No. 3375*. Washington DC: The World Bank.
- Calvo, César and Stefan Dercon. (2009). "Chronic Poverty and All That: The Measurement of Poverty Over Time". In Tony Addison, David Hulme, and Ravi Kanbur. (Eds.) *Poverty Dynamics: Interdisciplinary Perspectives*. Oxford University Press: New York.
- Cappellari, Lorenzo, and Stephen P. Jenkins. (2006). "Calculation of Multivariate Normal Probabilities by Simulation, with Applications to Maximum Simulated Likelihood Estimation". *Stata Journal*, 6(2): 156- 189.
- Casella, George and Roger L. Berger. (2002). *Statistical Inference*, 2nd Edition. California: Duxbury Press.
- Cross, Philip J. and Charles F. Manski. (2002). "Regressions, Short and Long". *Econometrica*, 70(1): 357-368.



- Cruces, Guillermo, Peter Lanjouw, Leonardo Lucchetti, Elizaveta Perova, Renos Vakis, and Mariana Viollaz. (forthcoming). “Estimating Poverty Transitions Repeated Cross-Sections: A Three-country Validation Exercise”. *Journal of Economic Inequality*.
- Dang, Hai-Anh, Peter Lanjouw, Jill Luoto, and David McKenzie. (2014). “Using Repeated Cross-Sections to Explore Movements in and out of Poverty”. *Journal of Development Economics*, 107: 112-128.
- Dang, Hai-Anh and Peter Lanjouw. (2014). “Welfare Dynamics Measurement: Two Definitions of a Vulnerability Line”. World Bank Policy Research Paper # 6944. Washington DC: The World Bank.
- Deaton, Angus. (1985). “Panel Data from Time Series of Cross-Sections”. *Journal of Econometrics*, 30: 109- 126.
- Deaton, Angus. (1997). “*The Analysis of Household Surveys: A Microeconomic Approach to Development Policy*.” MD: The Johns Hopkins University Press.
- Deaton, Angus and Christina Paxson. (1994). “Intertemporal Choice and Inequality”. *Journal of Political Economy*, 102(3): 437- 467.
- Demirguc-Kunt, Asli, Leora F. Klapper, and Georgios A. Panos. (2011). “Entrepreneurship in Post-Conflict Transition: The Role of Informality and Access to Finance”. *Economics of Transition*, 19(1): 27-78.
- Demombynes, Gabriel and Berk Özler. (2005). “Crime and Local Inequality in South Africa,” *Journal of Development Economics*, 76(2): 265–292.
- Devereux, Paul J. (2007). “Small-Sample Bias in Synthetic Cohort Models of Labor Supply”. *Journal of Applied Econometrics*, 22: 839-848.
- Dorfman, Alan H. (2011). “A Coverage Approach to Evaluating Mean Square Error”. *Pakistan Journal of Statistics*, 27(4): 493-506.
- Elbers, Chris, Jean O. Lanjouw, and Peter Lanjouw. (2002) “Micro-Level Estimation of Welfare”. *World Bank Policy Research Working Paper # 2911*.
- . (2003). “Micro-level Estimation of Poverty and Inequality”. *Econometrica*, 71(1): 355-364.
- Elbers, Chris, Tomoki Fujii, Peter Lanjouw, Berk Özler, and Wesley Yin. (2007). “Poverty Alleviation Through Geographic Targeting: How Much Does Disaggregation Help?” *Journal of Development Economics*, 83: 198–213.
- Foster, James E. (2009). “A Class of Chronic Poverty Measures”. In Tony Addison, David Hulme, and Ravi Kanbur. (Eds.) *Poverty Dynamics: Interdisciplinary Perspectives*. Oxford University Press: New York.

- Friedman, Lynn and Melanie Wall. (2005). "Graphical Views of Suppression and Multicollinearity in Multiple Linear Regression". *American Statistician*, 59(2): 127-136.
- Glewwe, Paul and Hanan Jacoby. (2000). "Recommendations for Collecting Panel Data". In Margaret Grosh and Paul Glewwe. (Eds). *Designing Household Survey Questionnaires for Developing Countries: Lessons from 15 Years of the Living Standards Measurement Study*. Washington DC: The World Bank.
- Güell, Maia and Luojia Hu. (2006). "Estimating the Probability of Leaving Unemployment Using Uncompleted Spells from Repeated Cross-Section Data". *Journal of Econometrics*, 133: 307-341.
- Inoue, Atsushi. (2008). "Efficient Estimation and Inference in Linear Pseudo-Panel Data Models". *Journal of Econometrics*, 142: 449- 466.
- Kalton, Graham. (2009). "Designs for Surveys over Time". In D. Pfeffermann and C.R. Rao. *Handbook of Statistics, Vol. 29A- Sample Surveys: Design, Methods and Applications*. North-Holland: Elsevier.
- Kish, Leslie. (1999). "Cumulating/ Combining Population Surveys". *Survey Methodology*, 25(2): 129- 138.
- . (2002). "New Paradigms (Models) for Probability Sampling". *Survey Methodology*, 28(1): 31- 34.
- Little, Roderick J. A. and Donald B. Rubin. (2002). *Statistical Analysis with Missing Data*. 2nd Edition. New Jersey: Wiley.
- Lorh, Sharon L. (2010). *Sampling, Design and Analysis*. Massachusetts: Duxbury Press.
- Martinez, Arturo Jr., Mark Western, Michele Haynes, Wojtek Tomaszewski. (2013). "Measuring Income Mobility Using Pseudo-Panel Data". *Philippine Statistician*, 62(2): 71-99.
- Matloff, Norman S. (1981). "Use of Regression Functions for Improved Estimation of Means". *Biometrika*, 68(3): 685-689.
- Mcintosh, Steven. (2006). "Further Analysis of the Returns to Academic and Vocational Qualifications". *Oxford Bulletin of Economics and Statistics*, 68(2): 225- 251.
- McKenzie, David. (2004). "Asymptotic Theory for Heterogeneous Dynamic Pseudo-Panels". *Journal of Econometrics*, 120, 235-262.
- Moffitt, Robert. (1993). "Identification and Estimation of Dynamic Models with a Time Series of Repeated Cross- Sections". *Journal of Econometrics*, 59: 99-123.

- Montgomery, Douglas C. (2012). *Introduction to Statistical Quality Control*. USA: Wiley.
- Mullahy, John. (2011). “Marginal Effects in Multivariate Probit and Kindred Discrete and Count Outcome Modes, with Applications in Health Economics”. *NBER Working paper* 17588.
- Pfeffermann, Danny. (2011). “Modelling of Complex Survey Data: Why model? Why Is It a Problem? How Can We Approach It?” *Survey Methodology*, 37(2): 115- 136.
- Pencavel, John. (2007). “A Life Cycle Perspective on Changes in Earnings Inequality among Married Men and Women”. *Review of Economics and Statistics*, 88(2): 232-242.
- Peruvian Statistics Bureau (INEI). [http://www.inei.gob.pe/srienaho/Consulta\\_por\\_Encuesta.asp](http://www.inei.gob.pe/srienaho/Consulta_por_Encuesta.asp) Accessed October 2012.
- Pham-Gia, T, N. Turkkan, and E. Marchand (2006). “Density of the Ratio of Two Normal Random Variables and Applications”. *Communications in Statistics- Theory and Method*, 35(9): 1569-1591.
- Plackett, R. L. (1954). “A Reduction Formula for Normal Multivariate Integrals”. *Biometrika*, 41:351-360.
- Prekopa, Andras. (1970). “On Probabilistic Constrained Programming”. *In Proceedings of the Princeton Symposium on Mathematical Programming*. New Jersey: Princeton Press.
- Propper, Carol, Hedley Rees, and Katherine Green. (2001). “The Demand for Private Medical Insurance in the UK: A Cohort Analysis”. *Economic Journal*, 111: C180-C200.
- PSID Main Interview User Manual: Release 2012.1. Institute for Social Research, University of Michigan, January 23, 2012. Available on the Internet at <http://psidonline.isr.umich.edu/data/Documentation/UserGuide2009.pdf> Accessed October 2012.
- Rao, J. N. K. (2003). *Small Area Estimation*. New Jersey: Wiley.
- Ridder, Geert and Robert Moffitt. (2007). “The Econometrics of Data Combination”. In Heckman and Leamer. (Eds). *Handbook of Econometrics*, Volume 6B. Elsevier: the Netherlands.
- Stock, James H. and Motohiro Yogo. (2005). “ Testing for Weak Instruments in Linear IV Regression.” In D. W. K. Andrews and J. H. Stock. (Eds.) *Identification and Inference for Econometric Models: Essays in Honor of Thomas Rothenberg*. Cambridge: Cambridge University Press.
- Sungur, Engin A. (1990). “Dependence Information in Parameterized Copulas”. *Communications in Statistics- Simulation and Computation*, 19(4): 1339-1360.

Tung, Phung Duc and Nguyen Phong. (undated). “*Vietnam Household Living Standards Surveys (VHLSSs) in 2002 and 2004- Basic Information*”. Available on the World Bank’s LSMS website at [http://siteresources.worldbank.org/INTLSMS/Resources/3358986-1181743055198/3877319-1207074161131/BINFO\\_VHLSS\\_02\\_04.pdf](http://siteresources.worldbank.org/INTLSMS/Resources/3358986-1181743055198/3877319-1207074161131/BINFO_VHLSS_02_04.pdf) Accessed October 2012.

Verbeek, Marno (2008) “Synthetic panels and repeated cross-sections”, pp.369-383 in L. Matyas and P. Sevestre (eds.) *The Econometrics of Panel Data*. Berlin: Springer-Verlag.

Verbeek, M. and T. Nijman. (1992). “Can Cohort Data Be Treated as Genuine Panel Data?” *Empirical Economics*, 17: 9- 23.

**Table 1: Model Parameters for Simulation**

Model	X	$V(\varepsilon_1)$	$V(\varepsilon_2)$	$\rho$
1	1	24.5	35.0	0.88
2	1, 2	19.5	28.9	0.85
3	1, 2, 3	17.5	22.3	0.71
4	1, 2, 3, 4	13.5	15.6	0.60
5	1, 2, 3, 4, 5	12.5	14.7	0.58
6	1, 2, 3, 4, 5, 6	9.5	10.8	0.43
7	1, 2, 3, 4, 5, 6, 7	7.5	6.9	0.22
8	1, 2, 3, 4, 5, 6, 7, 8	6.5	6.5	0.15

**Note:** The x's variables are added sequentially and cumulatively to each simulation model. For example, Model 1 includes only the intercept and x1, Model 2 includes the intercept, x1 and x2, and so on. The error terms v1 and v2 are assumed to follow a standard bivariate normal distribution where  $\text{Var}(v_j) = 6.5$ , for  $j = 1, 2$ , and  $\text{cov}(v_1, v_2) = 1$ . The vectors of coefficients are  $b_1 = (1, 1, 1, 1, 1, 1, 1, 1)$  and  $b_2 = (1.5, 1.2, 1.1, 1.05, 1.3, 0.9, 1.15, 1.4, 0.6)$ . The x's are assumed to follow a standard normal distribution, where their variances are respectively  $(2.5, 5, 6, 4, 1, 3, 2, 1)$  for  $x_k$ ,  $k = 1, \dots, 8$ .

**Table 2: Estimated  $\rho$  from Actual Panels and Synthetic Panels for Different Countries**

Country	Survey Year	Actual panels		Synthetic panels		Relative difference (%)	
		$P_{y_1y_2}$	$\rho$	$P_{y_1y_2}$	$\rho$	$P_{y_1y_2}$	$\rho$
Bosnia-Herzegovina	2001	0.48	0.45	0.43	0.40	-10.4	-11.1
	2004						
Lao PDR	2002-03	0.51	0.43	0.56	0.46	9.8	7.0
	2007-08						
Peru	2004	0.82	0.64	0.82	0.69	0.0	7.8
	2005						
	2005	0.82	0.66	0.80	0.63	-2.4	-4.5
	2006						
	2004	0.79	0.63	0.73	0.51	-7.6	-19.0
	2006						
Vietnam	2004	0.81	0.66	0.85	0.73	4.9	10.6
	2006						
	2006	0.78	0.63	0.85	0.76	9.0	21.6
	2008						
	2004	0.75	0.58	0.84	0.74	12.0	27.6
	2008						
United States	2005	0.76	0.66	0.89	0.84	17.1	27.3
	2007						
	2007	0.82	0.70	0.86	0.79	4.9	12.9
	2009						
	2005	0.72	0.57	0.71	0.59	-1.4	3.5
	2009						

**Note:** The synthetic panel estimates are based on cross sectional data except for Bosnia-Herzegovina and the US, where these estimates are based on two rounds of actual panel data.  $P_{y_1y_2}$  is the simple correlation across two survey rounds for household consumption for all countries except for the US, where it is the correlation for household income.  $\rho$  is the partial correlation between the residuals of the regression of household consumption on household head's gender, years of schooling, ethnicity, and residence areas. All estimates for  $P_{y_1y_2}$  are significant at the 0.01 level. Household heads' ages are restricted to between 25 and 55 in the first survey round and adjusted accordingly for the second survey round.

**Table 3: Poverty Dynamics Based on Synthetic Panel Data for Two Periods, Joint Probabilities (Percentage)**

Poverty Status  First Period & Second Period	Bosnia- Herzegovina		Lao PDR		Peru		United States		Vietnam	
	2001- 2004		2002/03- 2007/08		2005-06		2007-09		2006-08	
	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel
Poor, Poor	10.3 (1.7)	8.2 (0.2)	13.8 (1.2)	13.2 (0.4)	29.9 (1.3)	30.9 (0.4)	6.0 (0.4)	6.2 (0.2)	9.9 (0.8)	9.6 (0.3)
Poor, Nonpoor	12.6 (1.2)	12.6 (0.3)	14.3 (1.1)	13.2 (0.1)	11.6 (0.9)	12.3 (0.1)	3.8 (0.3)	3.2 (0.1)	5.9 (0.5)	4.9 (0.1)
Nonpoor, Poor	10.5 (1.4)	12.1 (0.2)	10.9 (1.0)	11.4 (0.2)	8.9 (0.8)	10.0 (0.1)	4.6 (0.4)	4.0 (0.1)	4.9 (0.5)	5.0 (0.1)
Nonpoor, Nonpoor	66.5 (2.2)	67.2 (0.6)	61.0 (1.6)	62.2 (0.6)	49.7 (1.6)	46.8 (0.4)	85.7 (0.6)	86.6 (0.3)	79.3 (1.0)	80.4 (0.4)
N	1342	1342	1989	3215	2250	9084	3368	3368	2722	3701

**Note:** Synthetic panels are constructed from cross sections for Lao PDR, Peru, and Vietnam and from panel halves for Bosnia-Herzegovina and the US. Predictions are obtained using the estimated parameters from the first and second survey rounds on data in the second survey round. Standard errors are obtained adjusting for complex survey design for all countries, except for the US PSID. All numbers are weighted using household weights for Peru, and population weights for other countries. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round.

**Table 4: Coverage Test for Poverty Dynamics Based on Synthetic Panel Data for Two Periods (Percentage)**

Poverty status	Bosnia- Herzegovina		Lao PDR		Peru		United States		Vietnam	
	2001- 2004		2002/03- 2007/08		2005-06		2007-09		2006-08	
	95% CI	Coverage (%)	95% CI	Coverage (%)	95% CI	Coverage (%)	95% CI	Coverage (%)	95% CI	Coverage (%)
Poor, Poor	Yes	100	Yes	100	Yes	100	Yes	100	Yes	100
Poor, Nonpoor	Yes	100	Yes	100	Yes	100	No	60.6	Yes	100
Nonpoor, Poor	Yes	100	Yes	100	Yes	100	Yes	100	Yes	100
Nonpoor, Nonpoor	Yes	100	Yes	100	Yes	66.3	Yes	72.5	Yes	100

**Note:** Coverage tests are calculated based on estimates from Table 2. For each country, column 1 (95% CI) shows whether estimates based on the synthetic panels fall within the 95% confidence interval (CI) of estimates based on the actual panels, and column 2 (Coverage) shows the proportion of the 95% confidence intervals around estimates based on the synthetic panels that fall within those of the actual panels.



**Table 5: Poverty Dynamics Based on Synthetic Panel Data for Two Periods, Conditional Probabilities (Percentage)**

Poverty Status  First Period--> Second Period	Bosnia- Herzegovina		Lao PDR		Peru		United States		Vietnam	
	2001- 2004		2002/03- 2007/08		2005-06		2007-09		2006-08	
	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual panel	Synthetic panel
Poor--> Poor	45.0 (4.6)	39.4 (1.2)	49.0 (3.0)	50.0 (1.6)	72.0 (1.9)	71.5 (1.0)	61.2 (2.2)	66.0 (1.8)	62.8 (2.8)	66.0 (1.5)
Poor--> Nonpoor	55.0 (4.6)	60.6 (1.7)	51.0 (3.0)	50.0 (1.1)	28.0 (1.9)	28.5 (0.3)	38.8 (2.2)	34.0 (0.8)	37.2 (2.8)	34.0 (0.6)
Nonpoor--> Poor	13.6 (1.8)	15.3 (0.2)	15.2 (1.3)	15.5 (0.3)	15.1 (1.3)	17.6 (0.2)	5.0 (0.4)	4.4 (0.1)	5.9 (0.6)	5.9 (0.1)
Nonpoor--> Nonpoor	86.4 (1.8)	84.7 (0.7)	84.8 (1.3)	84.5 (0.8)	84.9 (1.3)	82.4 (0.7)	95.0 (0.4)	95.6 (0.3)	94.1 (0.6)	94.1 (0.3)
<i>Goodness-of-fit Tests</i>										
Within 95% CI	4/4		4/4		4/4		2/4		4/4	
Within 1 standard error	2/4		4/4		2/4		0/4		2/4	
Coverage of 50% or more	4/4		4/4		4/4		3/4		4/4	
Coverage of 100%	4/4		4/4		2/4		2/4		4/4	
N	1342	1342	1989	3215	2250	9084	3368	3368	2723	3701
<p><b>Note:</b> Synthetic panels are constructed from cross sections for Lao PDR, Peru, and Vietnam and from panel halves for Bosnia-Herzegovina and the US. Predictions are obtained using the estimated parameters from the first and second survey rounds on data in the second survey round. Standard errors are obtained adjusting for complex survey design for all countries, except for the US PSID. All numbers are weighted using household weights for Peru, and population weights for other countries. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. The "Within 95% CI" row shows the number of times that the estimates based on the synthetic panels fall within the 95% confidence interval (CI) of the estimates based on the actual panels; the "Within 1 standard error" row shows a similar figure but using one standard error around the estimates based on the actual panels. The "Coverage of 50% or more" row shows the number of times that half or more of the 95% CI around the synthetic panel estimates overlap with those based on the actual panels overlap with those based on the actual panels; the "Coverage of 100%" row shows a similar figure for the number of times that the former fall completely inside the latter.</p>										

**Table 6: Poverty Dynamics Based on Synthetic Panel Data for Three Periods (Percentage)**

Poverty Status	Peru		United States		Vietnam	
	2004-05-06		2005-07-09		2004-06-08	
First, Second & Third Period	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual panel	Synthetic panel
Poor, Poor, Poor	26.6 (1.4)	24.0 (0.4)	4.0 (0.4)	4.0 (0.2)	8.1 (1.0)	7.6 (0.3)
Poor, Poor, Nonpoor	6.9 (0.7)	7.1 (0.1)	1.4 (0.2)	2.0 (0.0)	3.1 (0.6)	2.8 (0.0)
Poor, Nonpoor, Poor	4.4 (0.6)	3.3 (0.0)	1.0 (0.2)	0.5 (0.0)	2.3 (0.5)	2.9 (0.1)
Poor, Nonpoor, Nonpoor	7.2 (0.7)	6.4 (0.0)	2.7 (0.3)	2.8 (0.0)	6.6 (0.8)	5.0 (0.1)
Nonpoor, Poor, Poor	3.9 (0.5)	6.1 (0.1)	1.8 (0.2)	1.7 (0.1)	0.8 (0.2)	1.1 (0.0)
Nonpoor, Poor, Nonpoor	5.4 (0.6)	5.0 (0.0)	2.0 (0.3)	1.1 (0.0)	1.7 (0.4)	2.6 (0.0)
Nonpoor, Nonpoor, Poor	4.6 (0.6)	6.4 (0.0)	3.1 (0.3)	3.2 (0.1)	2.9 (0.5)	2.7 (0.0)
Nonpoor, Nonpoor, Nonpoor	41.0 (1.7)	41.6 (0.4)	84.0 (0.7)	84.6 (0.4)	74.5 (1.5)	75.3 (0.4)
<i>Goodness-of-fit Tests</i>						
Within 95% CI	6/8		5/8		6/8	
Within 1 standard error	3/8		5/8		4/8	
Coverage of 50% or more	6/8		5/8		6/8	
Coverage of 100%	4/8		5/8		6/8	
N	1987	8608	3036	3036	1282	3808

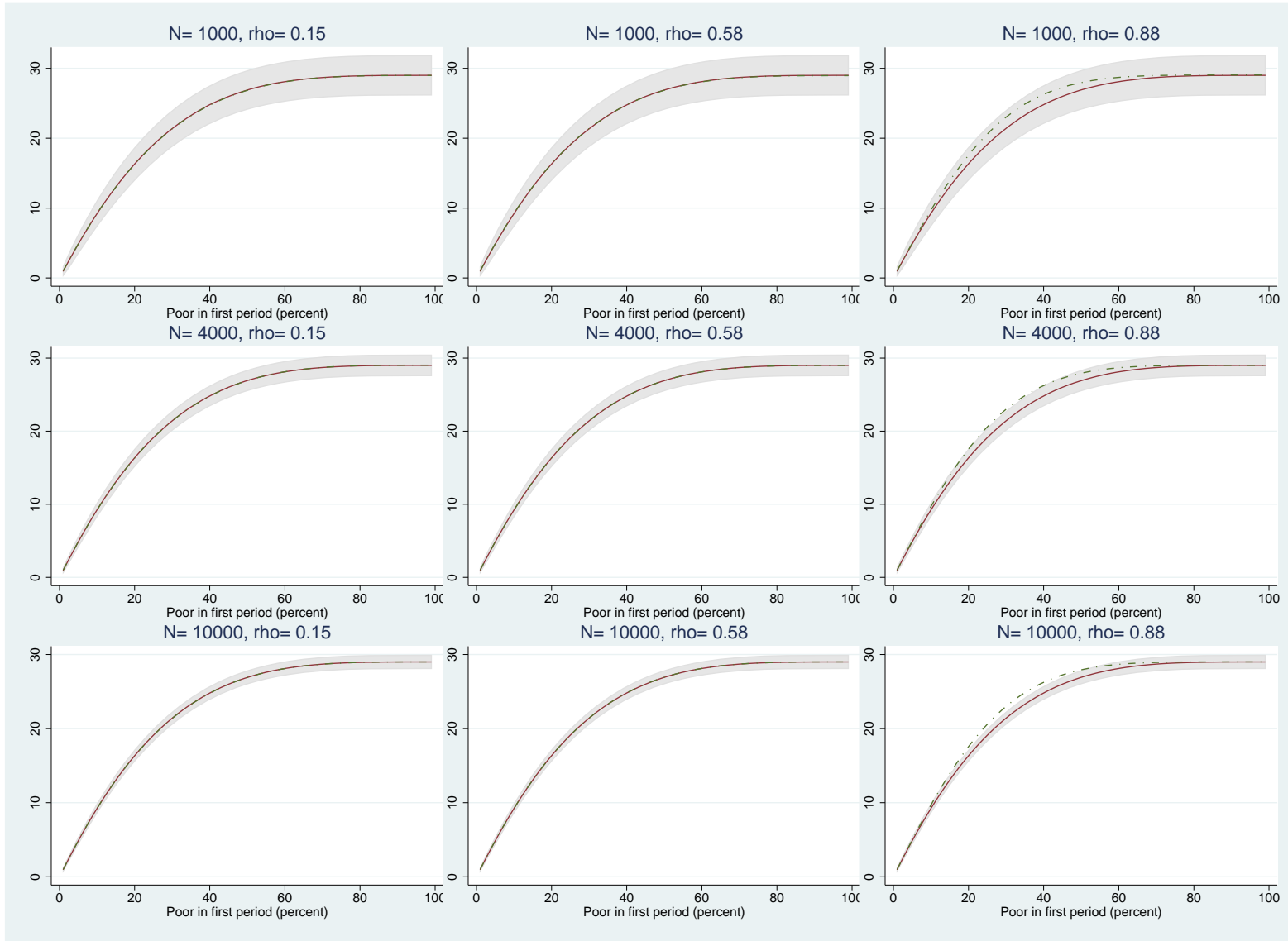
**Note:** Synthetic panels are constructed from cross sections for Peru and Vietnam and from panel halves for the US. Predictions are obtained using the estimated parameters from the first, second, and third survey rounds on data in the third survey round. Standard errors are obtained adjusting for complex survey design for all countries, except for the US PSID. All numbers are weighted using household weights for Peru, and population weights for other countries. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. The "Within 95% CI" row shows the number of times that the estimates based on the synthetic panels fall within the 95% confidence interval (CI) of the estimates based on the actual panels; the "Within 1 standard error" row shows a similar figure but using one standard error around the estimates based on the actual panels. The "Coverage of 50% or more" row shows the number of times that half or more of the 95% CI around the synthetic panel estimates overlap with those based on the actual panels; the "Coverage of 100%" row shows a similar figure for the number of times that the former fall completely inside the latter.

**Table 7: Consumption Dynamics for Two Periods, Vietnam 2006-2008 (Percentage)**

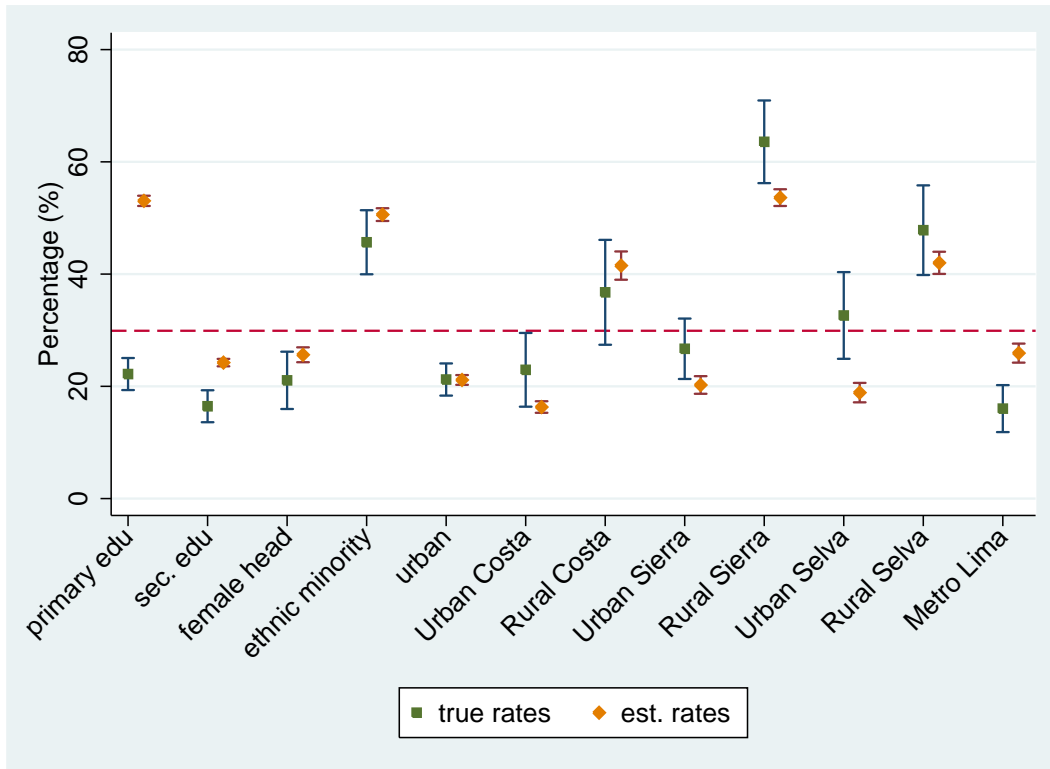
		2008						
		Poorest	Quintile 2	Quintile 3	Quintile 4	Richest	Total	
Panel A: True Panels	2006	Poorest	12.7 (0.8)	4.7 (0.4)	1.7 (0.3)	0.6 (0.2)	0.2 (0.1)	19.7 (0.9)
		Quintile 2	4.8 (0.4)	7.5 (0.6)	4.6 (0.5)	2.0 (0.3)	0.6 (0.1)	19.6 (0.9)
		Quintile 3	1.8 (0.3)	5.2 (0.5)	6.9 (0.5)	4.6 (0.5)	1.5 (0.2)	20.0 (0.9)
		Quintile 4	0.6 (0.2)	2.0 (0.3)	5.0 (0.5)	7.8 (0.6)	4.8 (0.5)	20.2 (0.9)
		Richest	0.1 (0.1)	0.6 (0.2)	1.8 (0.3)	4.9 (0.5)	12.9 (0.7)	20.5 (0.8)
		Total	20.0 (1.0)	20.0 (0.9)	20.0 (0.9)	20.0 (0.9)	20.0 (0.9)	100
Panel B: Synthetic Panels	2006	Poorest	<b>13.7</b> (0.3)	3.6 (0.0)	<b>1.6</b> (0.0)	<b>0.4</b> (0.0)	<b>0.0</b> (0.0)	19.2 (0.3)
		Quintile 2	<b>5.6</b> (0.1)	5.4 (0.0)	<b>4.5</b> (0.0)	<b>2.2</b> (0.0)	0.3 (0.0)	17.8 (0.1)
		Quintile 3	<b>2.3</b> (0.0)	<b>4.5</b> (0.0)	<b>6.4</b> (0.0)	5.6 (0.0)	<b>1.5</b> (0.0)	20.4 (0.1)
		Quintile 4	<b>0.6</b> (0.0)	<b>2.1</b> (0.0)	<b>5.1</b> (0.0)	<b>8.5</b> (0.1)	<b>5.2</b> (0.0)	21.4 (0.1)
		Richest	<b>0.0</b> (0.0)	0.3 (0.0)	<b>1.4</b> (0.0)	<b>5.4</b> (0.0)	<b>14.0</b> (0.2)	21.1 (0.2)
		Total	22.2 (0.3)	15.8 (0.1)	18.9 (0.1)	22.2 (0.1)	20.9 (0.2)	100

**Note:** Synthetic panels are constructed from cross sections for Vietnam. Predictions are obtained using the estimated parameters from the first and second survey rounds on data in the second survey round. Standard errors are obtained adjusting for complex survey design. All numbers are weighted using population weights. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. Joint probabilities are shown. Estimates based on the synthetic panels that fall within the 95% CI of those based on the actual panels are shown in bold.

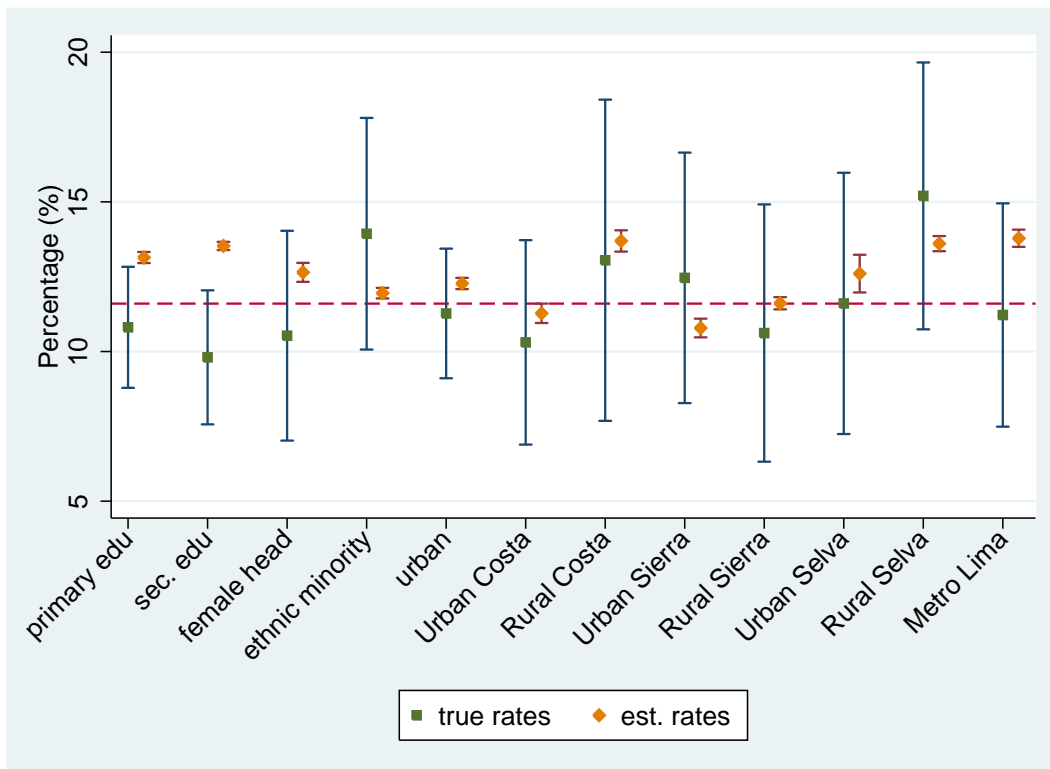
**Figure 1: Predicted Poverty Rates vs. True Poverty Rates for Two Periods Based on Simulated Data**



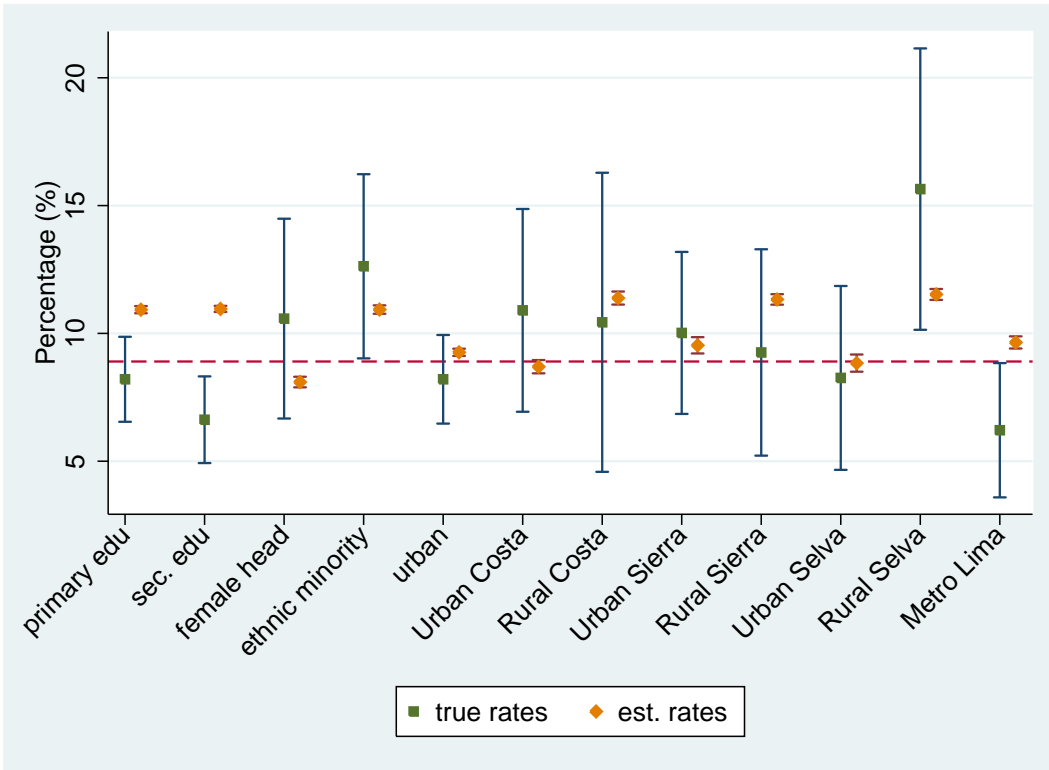
**Figure 2: Profiles for Those Who Remained Poor in Both Periods, Peru 2005- 2006**



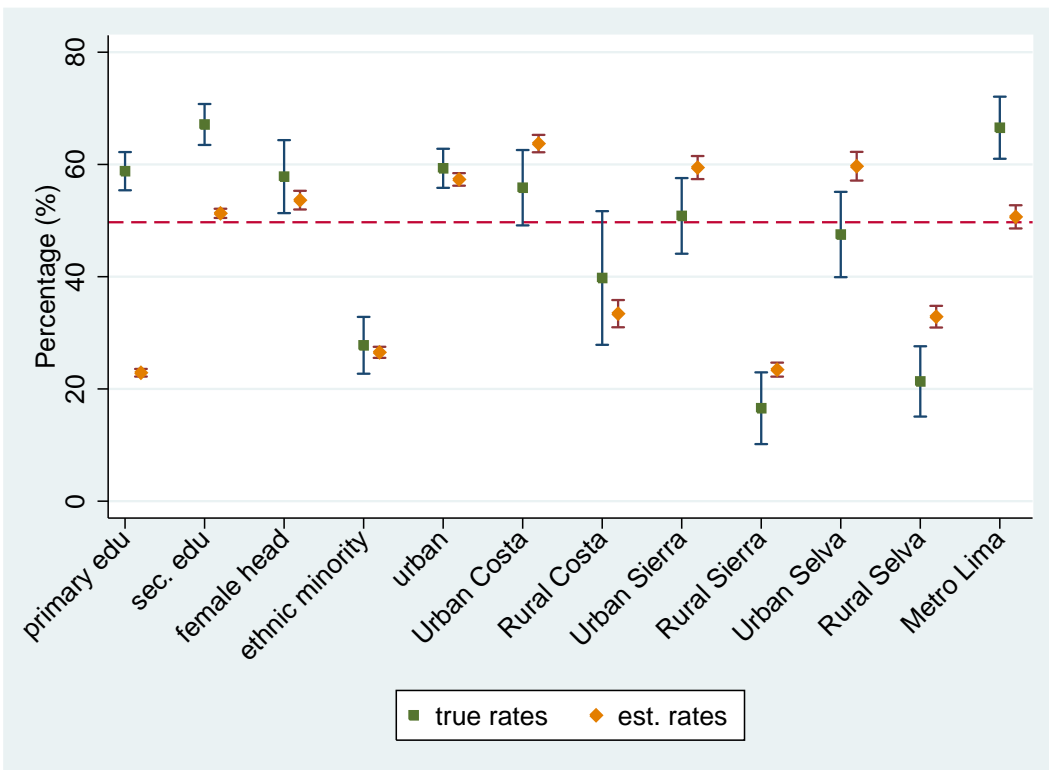
**Figure 3: Profiles for Those Who Were Poor in the First Period but Non-poor in the Second Period, Peru 2005- 2006**



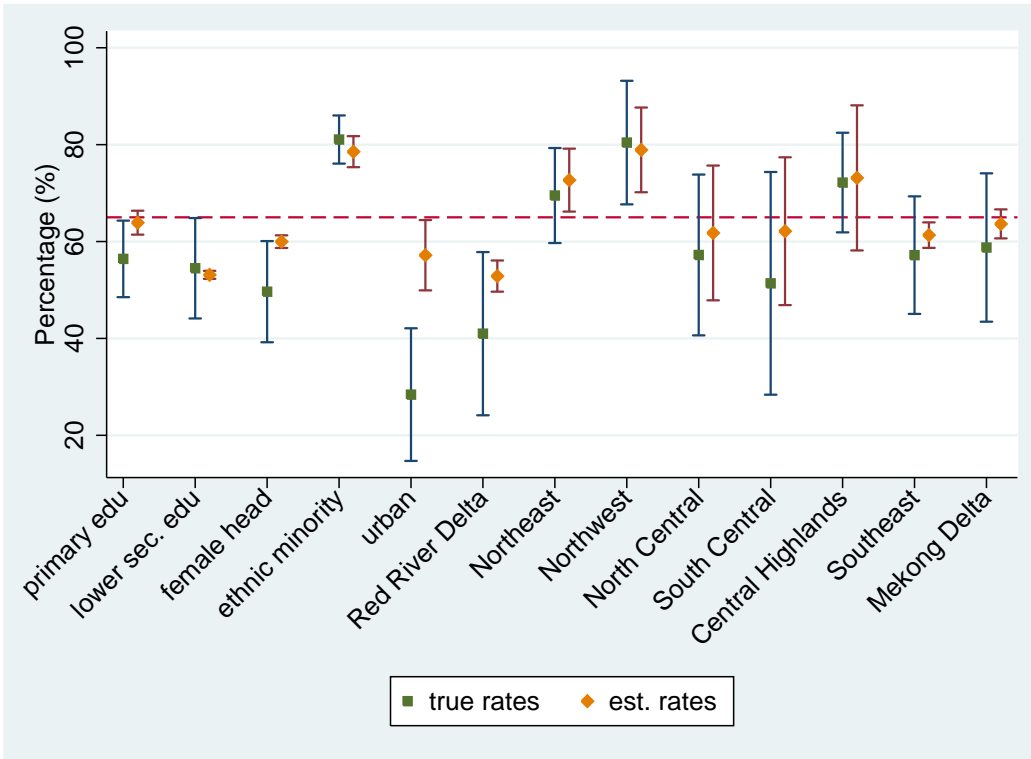
**Figure 4: Profiles for Those Who Were Non-poor in the First Period but Poor in the Second Period, Peru 2005- 2006**



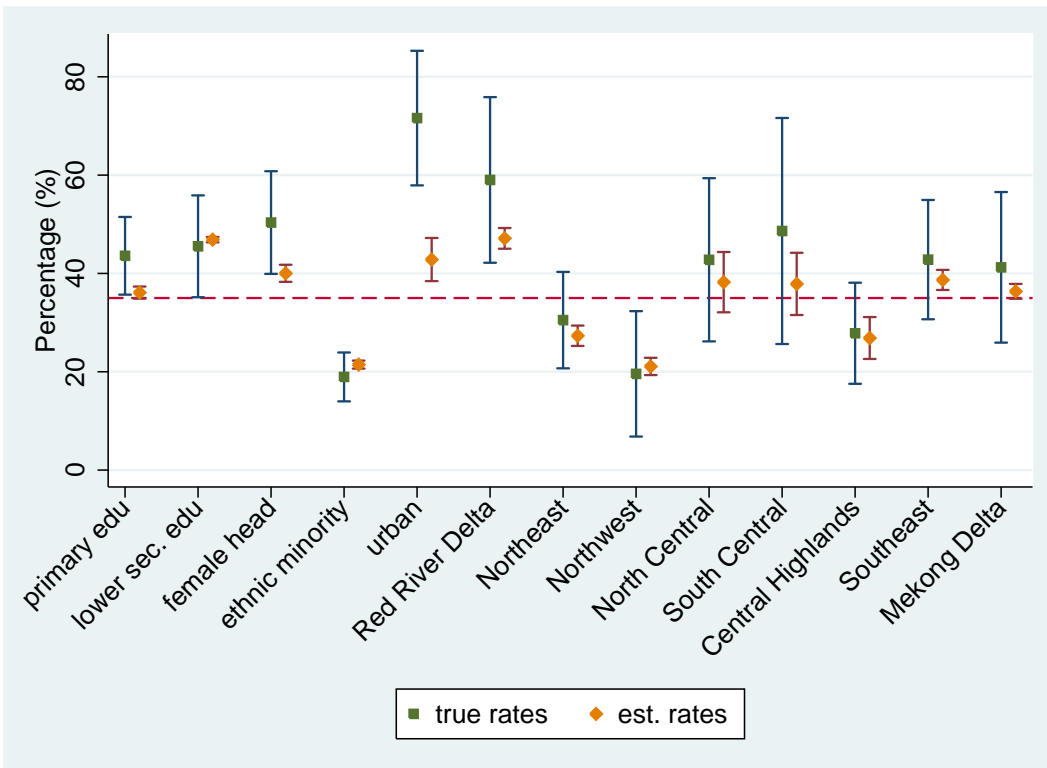
**Figure 5: Profiles for Those Who Remained Non-poor in Both Periods, Peru 2005- 2006**



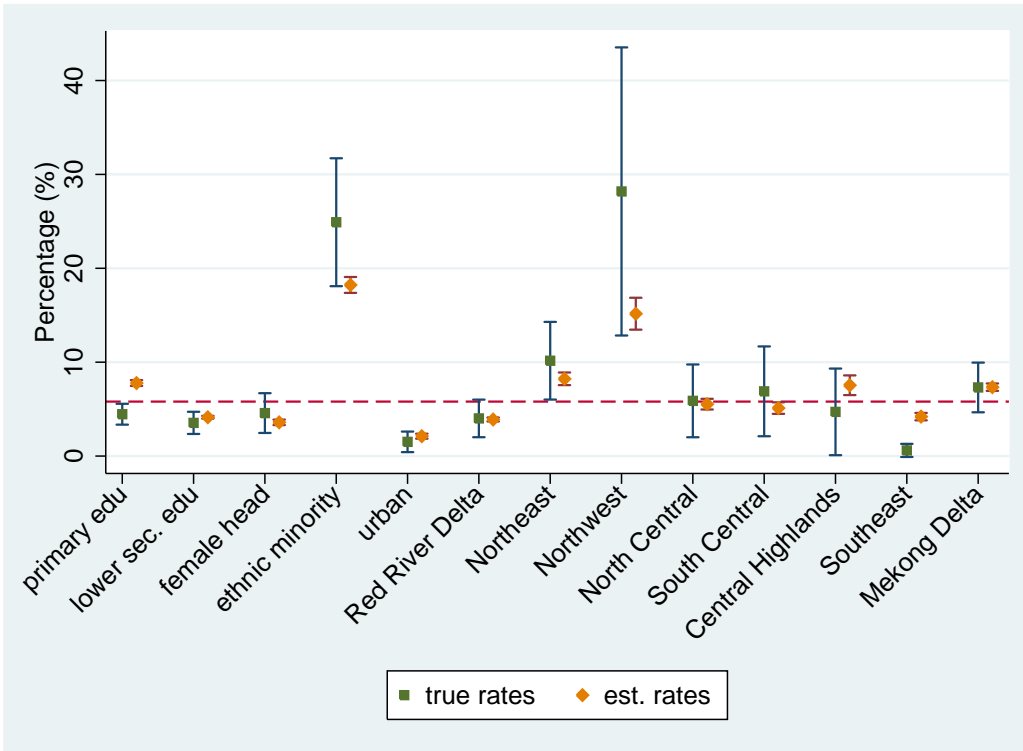
**Figure 6: Profiles for the Proportion of the Population Who Were Poor in the Second Period Given that They Were Poor in the First Period, Vietnam 2006- 2008**



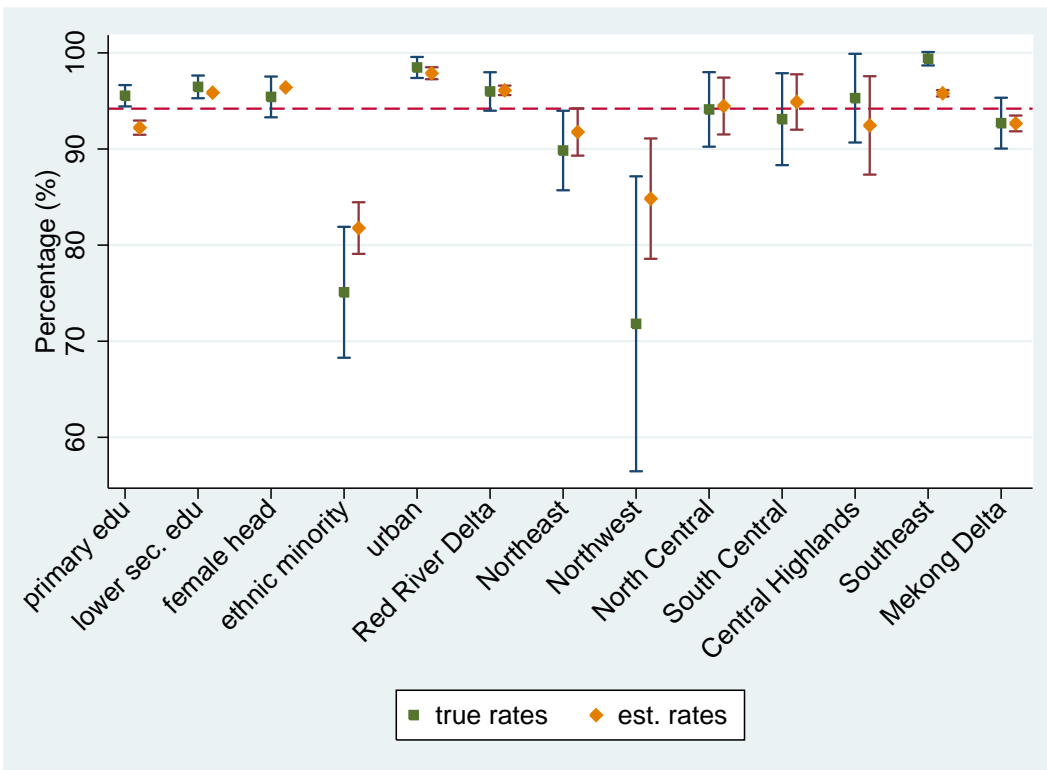
**Figure 7: Profiles for the Proportion of the Population Who Were Non-poor in the Second Period Given that They Were Poor in the First Period, Vietnam 2006- 2008**



**Figure 8: Profiles for the Proportion of the Population Who Were Poor in the Second Period Given that They Were Non-poor in the First Period, Vietnam 2006- 2008**



**Figure 9: Profiles for the Proportion of the Population Who Were Non-poor in the Second Period Given that They Were Non-poor in the First Period, Vietnam 2006- 2008**





## Appendix 1: Proofs

### Proof for Proposition 1

Consider a simple linear dynamic data-generating process for household consumption given by

$$y_{i2} = \alpha + \delta' y_{i1} + \eta_{i2} \quad (1.1)$$

where  $y_{it}$  is household  $i$ 's consumption in period  $t$ ,  $t = 1, 2$ , and  $\eta_{i2}$  is a random error term. Note that in the absence of actual panel data we do not observe  $y_{i1}$  for the same household, and we only have two repeated cross sections. Our objective of obtaining the simple correlation coefficient  $\rho_{y_{i1}y_{i2}}$  in this case is closely related to getting a consistent estimate for  $\delta$ , since by

definition  $\rho_{y_{i1}y_{i2}} = \frac{\text{cov}(y_{i1}, y_{i2})}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} = \sqrt{\frac{\text{var}(y_{i1})}{\text{var}(y_{i2})}} \delta$ . A consistent estimate for  $\delta$  can be

obtained by instrumenting for it with the cohort dummy variables, as long as these instrumental variables are relevant and exogenous. Thus estimation of (1.1) this way is identical to applying OLS to the same model where all variables are aggregated to the cohort level (Verbeek, 2008)

$$\bar{y}_{c2} = \delta' \bar{y}_{c1} + \bar{\eta}_{c2} \quad (1.2)$$

Thus from (1.2) we can consistently estimate  $\delta$ , and  $\rho_{y_{i1}y_{i2}}$  as  $\rho_{y_{c1}y_{c2}} = \frac{\text{cov}(\bar{y}_{c1}, \bar{y}_{c2})}{\sqrt{\text{var}(\bar{y}_{c1}) \text{var}(\bar{y}_{c2})}}$ .

### Proof for Proposition 2

If actual panel data were available, the simple correlation coefficient for household consumption between the two survey rounds would be

$$\begin{aligned} \rho_{y_{i1}y_{i2}} &= \frac{\text{cov}(y_{i1}, y_{i2})}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} = \frac{\text{cov}(\beta_1' x_{i1} + \varepsilon_{i1}, \beta_2' x_{i2} + \varepsilon_{i2})}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} \\ &= \frac{\beta_1' \text{var}(x_i) \beta_2 + \rho \sqrt{\sigma_{\varepsilon_1}^2 \sigma_{\varepsilon_2}^2}}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} \end{aligned}$$

where the second line follows from Assumption 2 in Dang et al.. The third line follows from Assumption 1 that the underlying population being sampled in survey rounds 1 and 2 are the same, thus the time-invariant household characteristics  $x_{i1}$  and  $x_{i2}$  are replaced with  $x_i$ . Solving

for  $\rho$  from the above equality, we have  $\rho = \frac{\rho_{y_{i1}y_{i2}} \sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})} - \beta_1' \text{var}(x_i) \beta_2}{\sigma_{\varepsilon_1} \sigma_{\varepsilon_2}}$ .

### Proof for Corollary 2.1

If  $\beta_1 \approx \beta_2$ , we have

$$\begin{aligned}
\rho &= \frac{\rho_{y_{i1}y_{i2}} \sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})} - \beta_1' \text{var}(x_i) \beta_2}{\sigma_{\varepsilon_1} \sigma_{\varepsilon_2}} \\
&= \frac{\rho_{y_{i1}y_{i2}} - \frac{\beta_1' \text{var}(x_i) \beta_2}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}}}{\sqrt{1 - R_1^2} \sqrt{1 - R_2^2}} \\
&= \frac{\rho_{y_{i1}y_{i2}} - \sqrt{\frac{\beta_1' \text{var}(x_i) \beta_2 \beta_1' \text{var}(x_i) \beta_2}{\text{var}(y_{i1}) \text{var}(y_{i2})}}}{\sqrt{1 - R_1^2} \sqrt{1 - R_2^2}} \\
&\approx \frac{\rho_{y_{i1}y_{i2}} - \sqrt{\frac{\beta_1' \text{var}(x_i) \beta_1 \beta_2' \text{var}(x_i) \beta_2}{\text{var}(y_{i1}) \text{var}(y_{i2})}}}{\sqrt{1 - R_1^2} \sqrt{1 - R_2^2}} \\
&= \frac{\rho_{y_{i1}y_{i2}} - \sqrt{R_1^2 R_2^2}}{\sqrt{1 - R_1^2} \sqrt{1 - R_2^2}}
\end{aligned}$$

where the denominator in the second row follows from the definition for  $R^2$ , and the last equality follows from the definition of  $R^2$ .

In fact, given that  $\beta_1 \approx \beta_2$ , another perhaps more intuitive way to prove Corollary 2.1 is for us to think about the systemic part of predicted household consumption as a single variable (i.e.,  $\hat{\beta}_1' x_{i1} \approx \hat{\beta}_2' x_{i2}$ ). By definition, the correlation between household consumption and this predicted variable is the multiple correlation coefficient  $R_j^2$ . Thus using the familiar expression that links the simple and partial correlation coefficients for bivariate normal variables<sup>36</sup>

$$\rho_{12.3} = \frac{\rho_{12} - \rho_{13}\rho_{23}}{\sqrt{1 - \rho_{13}^2} \sqrt{1 - \rho_{23}^2}}$$

and replacing  $\rho_{12}$  with  $\rho_{y_{i1}y_{i2}}$  and the simple correlation coefficients  $\rho_{13}$  and  $\rho_{23}$  respectively with  $R_1^2$  and  $R_2^2$ , we have the expression in (6).

### Proof for Corollary 2.2

The partial correlation coefficient between the error terms is then given as usual

$$\rho = \frac{\text{cov}(\varepsilon_{i1}, \varepsilon_{i2})}{\sqrt{\text{var}(\varepsilon_{i1}) \text{var}(\varepsilon_{i2})}} = \frac{\sigma_u^2}{\sigma_u^2 + \sigma_v^2},$$

Assumption 2. If we hold  $\sigma_v^2$  fixed, it is straightforward to show that  $\rho$  is an increasing and concave function of  $\sigma_u^2$ . Thus  $\rho$  reaches its maximum when  $\sigma_u^2$  reaches its maximum value.

When  $\sigma_u^2$  reaches its maximum value  $\sigma_u^{2+}$  (or  $R_1^2 = R_2^2 = 0$ ), the estimation equations (1) and (2) have zero predictive power or all the terms with the estimated coefficients  $\beta$ 's are equal to 0 and will drop out. Thus we have

<sup>36</sup> See, for example, equation (36) in Ridder and Moffitt (2007) or equation (20) in Anderson (2003, p. 39).

$$\rho \leq \rho(\sigma_u^{2+}) = \frac{\sigma_u^{2+}}{\sigma_u^{2+} + \sigma_v^2} = \frac{\text{cov}(y_{i1}, y_{i2})}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} = \rho_{y_{i1}y_{i2}}$$

with equality occurring only when the model has zero predictive power. Then the simple correlation coefficient  $\rho_{y_{i1}y_{i2}}$  would be identical to the partial correlation coefficient  $\rho$ .

If  $\beta_1 \approx \beta_2$ , we can use Corollary 2.1 for an alternative proof. Holding  $R_2^2$  fixed in the expression in (6), taking the first derivative of  $\rho$  with regards to  $R_1^2$ , we have

$$\frac{\partial \rho}{R_1^2} = \frac{1}{\sqrt{1-R_2^2}} \left( \frac{-\frac{R_2^2}{2\sqrt{R_1^2 R_2^2 (1-R_1^2)}} - \frac{\rho_{y_{i1}y_{i2}} - \sqrt{R_1^2 R_2^2}}{2\sqrt{(1-R_1^2)^3}}}{1-R_1^2} \right)$$

By definition, both  $R_1^2$  and  $R_2^2$  are bounded by  $[0, 1]$ , and both two terms in the numerators are non-positive,<sup>37</sup> we have  $\frac{\partial \rho}{R_1^2} \leq 0$ . Thus given any value for  $R_2^2$ ,  $\rho$  is a decreasing function of  $R_1^2$  and  $\rho$  reaches its maximum value when  $R_1^2$  equals 0. By a similar argument, given any value for  $R_1^2$ , we can show that  $\rho$  is a decreasing function of  $R_2^2$ , and  $\rho$  reaches its maximum value when  $R_2^2$  also equals 0. Thus  $\rho$  is less than or equal to  $\rho_{y_{i1}y_{i2}}$ , with equality occurring when both  $R_1^2$  and  $R_2^2$  equal 0.

### Proof for Corollary 2.3

Corollary 2.3 directly follows from our assumption that  $\rho$  is non-negative. From the derivations in the proofs for Proposition 2 and Corollary 2.1 above, we see that

$$\text{i) } \rho_{y_{i1}y_{i2}} = \frac{\beta_1' \text{var}(x_i) \beta_2}{\sqrt{\text{var}(y_{i1}) \text{var}(y_{i2})}} \text{ when } \rho \sqrt{\sigma_{\varepsilon_1}^2 \sigma_{\varepsilon_2}^2} = 0 \text{ or}$$

ii)  $\rho_{y_{i1}y_{i2}} - \sqrt{R_1^2 R_2^2}$  when the estimation model fully captures all the variations in the dependent variable (i.e., all the error terms are zero).

### Proof for Propositions 3 and 4

To save space, we only show the proof for Proposition 4 since the two survey rounds scenario in Proposition 3 is a special case of the k survey rounds case in Proposition 4.

Given that household consumption can be explained by household characteristics in equations (1) and (2) and the standard regularity conditions are satisfied, our estimator  $\hat{\Phi}_2(\cdot)$  is a continuous and differentiable function of  $\hat{\beta}_m, \hat{\sigma}_{\varepsilon_m}, \hat{\rho}_{y_{im}y_{in},d}$ , for  $m=1, \dots, k-1, n=m+1, \dots, k$ ,

---

<sup>37</sup> Note that by Assumption 2,  $\rho$  is non-negative thus  $\rho_{y_{i1}y_{i2}} - \sqrt{R_1^2 R_2^2}$  is non-negative.

and  $j \neq m, n$ , which are consistent estimators of the parameters. Thus  $\hat{\Phi}_2(\cdot)$  is a consistent estimator of  $\Phi_2(\cdot)$ .

We can then decompose the variance for  $P - \hat{\Phi}_k(\cdot)$  into two parts, one due to sampling errors and the other due to model errors

$$\begin{aligned} \text{Var}(P - \hat{\Phi}_k(\cdot)) &= \text{Var}\left((P - \Phi_k(\cdot)) + (\Phi_k(\cdot) - \hat{\Phi}_k(\cdot))\right) \\ &= \Sigma_s + \Sigma_m \end{aligned}$$

assuming that these two errors are uncorrelated with each other.

The variance for the sampling errors  $\Sigma_s$  can be estimated using the bootstrap method.

Using the delta method, the variance for the model errors  $\Sigma_m$  can be written as

$$\sum_{m=1}^k \nabla'_{\hat{\beta}_m} V(\hat{\beta}_m) \nabla_{\hat{\beta}_m} + \sum_{m=1}^k \nabla'_{\hat{\sigma}_{\varepsilon_m}} V(\hat{\sigma}_{\varepsilon_m}) \nabla_{\hat{\sigma}_{\varepsilon_m}} + \sum_{m=1}^{k-1} \sum_{n=m+1}^k \nabla'_{\hat{\rho}_{y_m y_n, d}} V(\hat{\rho}_{y_m y_n, d}) \nabla_{\hat{\rho}_{y_m y_n, d}}$$

where applying the chain rule and taking the first partial derivative with regards to  $\hat{\beta}_m$  and  $\hat{\sigma}_{\varepsilon_m}$  (see, for example, Prekopa (1970)) and  $\hat{\rho}_{y_m y_n, d}$  (see, for example, Plackett (1954)) we have the

stated results.<sup>38</sup> Note that the approximation formula for  $V(\hat{\sigma}_{\varepsilon_m}) = \frac{(8N-7)\hat{\sigma}_{\varepsilon_m}^2}{(4N-3)^2}$  is based on

Montgomery (2012, pp. 720) where  $N > 25$ .

### Proof for Corollary 3.1

Since  $\hat{\Phi}(\cdot) = \frac{1}{N} \sum_{i=1}^N \hat{\Phi} \left( d_j \frac{z_j - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_{ij}}} \right)$  is a consistent estimator of  $P_{ij}$  and  $\hat{\Phi}_2(\cdot)$  is a consistent

estimator of  $P_{i,12}$  as discussed in the proof for Proposition 3 above, it follows that  $\frac{\hat{\Phi}_2(\cdot)}{\hat{\Phi}(\cdot)}$  is a

consistent estimator of  $\frac{P_{i,12}}{P_{ij}}$ . Then note that, since  $\frac{\partial(\hat{\Phi}_2(\cdot)/\hat{\Phi}(\cdot))}{\partial \hat{\Phi}_2(\cdot)} = \frac{1}{\hat{\Phi}(\cdot)}$  and

$\frac{\partial(\hat{\Phi}_2(\cdot)/\hat{\Phi}(\cdot))}{\partial \hat{\Phi}(\cdot)} = \frac{-\hat{\Phi}_2(\cdot)}{(\hat{\Phi}(\cdot))^2}$ , using the delta method<sup>39</sup> we have

$$\sqrt{n} \left[ \frac{\frac{P_{i,12}}{P_{ij}} - \frac{\hat{\Phi}_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{i1}}{\hat{\sigma}_{\varepsilon_{i1}}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{i2}}{\hat{\sigma}_{\varepsilon_{i2}}}, \hat{\rho}_d \right)}{\hat{\Phi} \left( d_j \frac{z_j - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_{ij}}} \right)}}{\frac{\hat{\Phi}_2 \left( d_1 \frac{z_1 - \hat{\beta}_1' x_{i1}}{\hat{\sigma}_{\varepsilon_{i1}}}, d_2 \frac{z_2 - \hat{\beta}_2' x_{i2}}{\hat{\sigma}_{\varepsilon_{i2}}}, \hat{\rho}_d \right)}{\hat{\Phi} \left( d_j \frac{z_j - \hat{\beta}_1' x_{ij}}{\hat{\sigma}_{\varepsilon_{ij}}} \right)}} \right] \sim N(0, V_r) \text{ where the covariance-variance}$$

<sup>38</sup> See also Mullahy (2011) for a related derivation.

<sup>39</sup> See, for example, theorem 5.5.28 in Casella and Berger (2002).

matrix  $V_r$  can be estimated as

$$\begin{aligned} V_r &= \frac{1}{\left(\hat{\Phi}(\cdot)\right)^2} \text{Var}(\hat{\Phi}_2(\cdot)) + \frac{\left(\hat{\Phi}_2(\cdot)\right)^2}{\left(\hat{\Phi}(\cdot)\right)^4} \text{Var}(\hat{\Phi}(\cdot)) - 2 \frac{\hat{\Phi}_2(\cdot)}{\left(\hat{\Phi}(\cdot)\right)^3} \text{Cov}(\hat{\Phi}_2(\cdot), \hat{\Phi}(\cdot)) \\ &= \left(\frac{\hat{\Phi}_2(\cdot)}{\hat{\Phi}(\cdot)}\right)^2 \left[ \frac{\text{Var}(\hat{\Phi}_2(\cdot))}{\left(\hat{\Phi}_2(\cdot)\right)^2} + \frac{\text{Var}(\hat{\Phi}(\cdot))}{\left(\hat{\Phi}(\cdot)\right)^2} - 2 \frac{\text{Cov}(\hat{\Phi}_2(\cdot), \hat{\Phi}(\cdot))}{\hat{\Phi}_2(\cdot)\hat{\Phi}(\cdot)} \right] \end{aligned}$$

where similar to  $\text{Var}(\hat{\Phi}_2(\cdot))$ ,  $\text{Var}(\hat{\Phi}(\cdot))$  can be decomposed into a model error  $\Sigma_{jm}$  and a sampling error  $\Sigma_{js}$  assuming these two errors are uncorrelated.<sup>40</sup> The model error can be

$$\begin{aligned} \text{estimated as } \Sigma_{jm} &= \nabla'_{\hat{\beta}_j} V(\hat{\beta}_j) \nabla_{\hat{\beta}_j} + \nabla'_{\hat{\sigma}_{\varepsilon_j}} V(\hat{\sigma}_{\varepsilon_j}) \nabla_{\hat{\sigma}_{\varepsilon_j}} \text{ with } \nabla_{\hat{\beta}_j} = d_j \begin{pmatrix} -x_{ij} \\ \hat{\sigma}_{\varepsilon_j} \end{pmatrix} \phi \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right) \text{ and} \\ \nabla_{\hat{\sigma}_{\varepsilon_j}} &= -d_j \begin{pmatrix} z_j - \hat{\beta}_j' x_{ij} \\ \hat{\sigma}_{\varepsilon_j}^2 \end{pmatrix} \phi \left( d_j \frac{z_j - \hat{\beta}_j' x_{ij}}{\hat{\sigma}_{\varepsilon_j}} \right). \end{aligned}$$

### Proof for Proposition 5

The consistency part of this proof is similar as that for Proposition 4 above.

For an illustration with the 5x5 transition, we can write out the formulae for the transition from the 2<sup>nd</sup> quintile in the first year as follows

$$\begin{aligned} P^{21} &= P(z_1^1 < y_{i1} < z_1^2 \text{ and } y_{i2} < z_2^1) = F^{21} - F^{11} \\ P^{22} &= P(z_1^1 < y_{i1} < z_1^2 \text{ and } z_2^1 < y_{i2} < z_2^2) = F^{22} - F^{21} - F^{12} + F^{11} \\ P^{23} &= P(z_1^1 < y_{i1} < z_1^2 \text{ and } z_2^2 < y_{i2} < z_2^3) = F^{23} - F^{22} - F^{13} + F^{12} \\ P^{24} &= P(z_1^1 < y_{i1} < z_1^2 \text{ and } z_2^3 < y_{i2} < z_2^4) = F^{24} - F^{23} - F^{14} + F^{13} \\ P^{25} &= P(z_1^1 < y_{i1} < z_1^2 \text{ and } z_2^4 < y_{i2} < +\infty) = F^{25} - F^{24} - F^{15} + F^{14} \end{aligned}$$

where, for example,  $F^{11} = \Phi_2 \left( \frac{z_1^1 - \beta_1' x_{i2}}{\sigma_{\varepsilon_1}}, \frac{z_2^1 - \beta_2' x_{i2}}{\sigma_{\varepsilon_2}}, \rho \right)$ . For the case where  $k=0$  or 5, note

that we have, for example,  $F^{15} = \Phi \left( \frac{z_1^1 - \beta_1' x_{i2}}{\sigma_{\varepsilon_1}} \right)$  and  $F^{01} = 0$  by the definition of the bivariate

normal probability function. See also Dang and Lanjouw (2014) for a discussion of the 3x3 transition.

<sup>40</sup> Pham-Gia, Turkkan and Marchand (2006) offer an alternative expression of the density of a ratio of two normal random variables in terms of Hermite and confluent hypergeometric functions.

## Appendix 2: Additional Tables

**Table 2.1: Estimated Parameters of Household Consumption for Each Year**

	Bosnia-Herzegovina		Lao PDR		Peru				United States				Vietnam			
	2001	2004	2002/03	2007/08	2004-05		2005-06		2005-07		2007-09		2004-06		2006-08	
Age	0.006*** (0.002)	0.012*** (0.002)	0.004*** (0.001)	0.006*** (0.001)	0.010*** (0.001)	0.012*** (0.001)	0.012*** (0.001)	0.013*** (0.001)	0.008*** (0.001)	0.006*** (0.001)	0.011*** (0.001)	0.008*** (0.001)	0.009*** (0.001)	0.010*** (0.001)	0.011*** (0.001)	0.009*** (0.001)
Female	0.190*** (0.041)	0.277*** (0.043)	0.086* (0.048)	0.137*** (0.041)	0.166*** (0.022)	0.153*** (0.016)	0.144*** (0.016)	0.192*** (0.016)	-0.306*** (0.016)	-0.463*** (0.020)	-0.433*** (0.020)	-0.516*** (0.024)	0.133*** (0.023)	0.094*** (0.021)	0.084*** (0.022)	0.113*** (0.022)
Years of schooling	0.035*** (0.005)	0.038*** (0.005)	0.032*** (0.003)	0.046*** (0.003)	0.064*** (0.002)	0.068*** (0.002)	0.068*** (0.002)	0.067*** (0.002)	0.419*** (0.022)	0.579*** (0.028)	0.573*** (0.028)	0.794*** (0.036)	0.051*** (0.003)	0.053*** (0.002)	0.053*** (0.003)	0.056*** (0.003)
Bosnian	-0.227*** (0.051)	-0.042 (0.053)														
Serb	-0.128** (0.051)	-0.068 (0.053)														
Ethnic majority group			0.239*** (0.021)	0.261*** (0.022)	0.209*** (0.025)	0.197*** (0.018)	0.188*** (0.018)	0.205*** (0.017)	0.150*** (0.016)	0.182*** (0.020)	0.200*** (0.020)	0.253*** (0.024)	0.393*** (0.027)	0.389*** (0.025)	0.361*** (0.026)	0.383*** (0.026)
Urban	-0.151*** (0.030)	-0.020 (0.031)	0.132*** (0.026)	0.133*** (0.024)	0.352*** (0.027)	0.430*** (0.020)	0.439*** (0.020)	0.446*** (0.019)	0.004*** (0.001)	0.008*** (0.002)	0.008*** (0.002)	0.006*** (0.002)	0.529*** (0.026)	0.447*** (0.024)	0.433*** (0.024)	0.310*** (0.023)
Constant	7.525*** (0.119)	7.022*** (0.131)	11.264*** (0.051)	11.658*** (0.055)	4.091*** (0.057)	3.928*** (0.040)	3.937*** (0.040)	3.946*** (0.040)	10.822*** (0.040)	10.538*** (0.053)	10.360*** (0.049)	10.086*** (0.065)	6.901*** (0.053)	7.192*** (0.048)	7.166*** (0.051)	7.492*** (0.050)
$\sigma_v$	0.522	0.543	0.518	0.537	0.547	0.556	0.553	0.546	0.407	0.519	0.511	0.628	0.473	0.482	0.485	0.489
Adjusted R <sup>2</sup>	0.077	0.077	0.157	0.218	0.407	0.443	0.443	0.463	0.293	0.328	0.335	0.340	0.454	0.421	0.407	0.370
N	1342	1342	3032	3215	4493	9169	8593	9084	3275	3275	3368	3368	3527	3674	3596	3701

**Note:** \*p<0.1, \*\*p<0.05, \*\*\*p<0.01. Standard errors are in parentheses. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. For the US, dummy variables for college degree and being white are used instead of years of schooling and ethnic majority group respectively. Other control variables used for the US include dummy variables indicating high school education and dummy variables indicating religion. Estimation is provided using the cross sections for Lao PDR, Peru, and Vietnam and from panel halves for Bosnia-Herzegovina and the US.

**Table 2.2: Poverty Dynamics Based on Synthetic Data for Two Periods Using Earlier Survey Rounds (Percentage)**

Poverty Status	Peru		United States		Vietnam	
	2004-05		2005-07		2004-06	
	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel
Poor, Poor	32.4 (1.5)	32.7 (0.4)	5.6 (0.4)	7.2 (0.2)	11.3 (0.8)	11.0 (0.3)
Poor, Nonpoor	9.8 (0.8)	9.7 (0.1)	3.9 (0.3)	3.8 (0.1)	9.2 (0.6)	7.8 (0.1)
Nonpoor, Poor	9.7 (0.9)	11.2 (0.1)	4.0 (0.3)	3.1 (0.1)	4.0 (0.5)	3.9 (0.0)
Nonpoor, Nonpoor	48.1 (1.8)	46.4 (0.4)	86.5 (0.6)	85.8 (0.3)	75.5 (1.0)	77.3 (0.4)
<i>Goodness-of-fit Tests</i>						
Within 95% CI	4/4		2/4		3/4	
Within 1 standard error	3/4		1/4		2/4	
Coverage of 50% or more	4/4		2/4		3/4	
Coverage of 100%	4/4		1/4		2/4	
N	2087	9169	3275	3275	2703	3674

**Note:** Synthetic panels are constructed from cross sections for Peru and Vietnam and from panel halves for the US. Predictions are obtained using the estimated parameters from the first and second survey rounds on data in the second survey round. Standard errors are obtained adjusting for complex survey design for all countries, except for the US PSID. All numbers are weighted using household weights for Peru, and population weights for other countries. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. The "Within 95% CI" row shows the number of times that the estimates based on the synthetic panels fall within the 95% confidence interval (CI) of the estimates based on the actual panels; the "Within 1 standard error" row shows a similar figure but using one standard error around the estimates based on the actual panels. The "Coverage of 50% or more" row shows the number of times that half or more of the 95% CI around the synthetic panel estimates overlap with those based on the actual panels; the "Coverage of 100%" row shows a similar figure for the number of times that the former fall completely inside the latter.

**Table 2.3: Poverty Dynamics Based on Synthetic Data for Two Periods, Using Data in the First Survey Round as the Base (Percentage)**

Poverty Status  First Period & Second Period	Bosnia- Herzegovina		Lao PDR		Peru		United States		Vietnam	
	2001- 2004		2002/03- 2007/08		2005-06		2007-09		2006-08	
	Actual panel	Synthetic panel	Actual panel	Synthetic panel	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel	Actual Panel	Synthetic Panel
Poor, Poor	10.8 (2.3)	9.8 (0.3)	13.5 (1.2)	15.1 (0.4)	29.4 (1.4)	32.2 (0.4)	6.0 (0.4)	6.2 (0.2)	9.6 (0.7)	10.2 (0.3)
Poor, Nonpoor	13.4 (1.5)	12.3 (0.3)	16.0 (1.1)	13.6 (0.1)	11.7 (0.9)	11.9 (0.1)	3.8 (0.3)	3.2 (0.1)	6.2 (0.6)	5.2 (0.1)
Nonpoor, Poor	11.5 (1.7)	14.4 (0.2)	8.9 (0.9)	12.4 (0.2)	8.8 (0.7)	9.7 (0.1)	4.6 (0.4)	4.0 (0.1)	4.5 (0.5)	5.2 (0.1)
Nonpoor, Nonpoor	64.3 (2.7)	63.5 (0.7)	61.7 (1.6)	59.0 (0.6)	50.1 (1.6)	46.2 (0.4)	85.7 (0.6)	86.6 (0.3)	79.7 (1.0)	79.4 (0.4)
<i>Goodness-of-fit Tests</i>										
Within 95% CI	4/4		2/4		2/4		3/4		4/4	
Within 1 standard error	3/4		0/4		1/4		2/4		2/4	
Coverage of 50% or more	4/4		2/4		2/4		3/4		4/4	
Coverage of 100%	4/4		1/4		2/4		2/4		4/4	
N	1342	1342	1989	3032	2250	8593	3368	3368	2723	3596

**Note:** Synthetic panels are constructed from cross sections for Lao PDR, Peru, and Vietnam and from panel halves for Bosnia-Herzegovina and the US. Predictions are obtained using the estimated parameters from the first and second survey rounds on data in the first survey round. Standard errors are obtained adjusting for complex survey design for all countries, except for the US PSID. All numbers are weighted using household weights for Peru, and population weights for other countries. Poverty rates are in percent. Household heads' ages are restricted to between 25 and 55 for the first survey round and adjusted accordingly with the year difference for the second survey round. The "Within 95% CI" row shows the number of times that the estimates based on the synthetic panels fall within the 95% confidence interval (CI) of the estimates based on the actual panels; the "Within 1 standard error" row shows a similar figure but using one standard error around the estimates based on the actual panels. The "Coverage of 50% or more" row shows the number of times that half or more of the 95% CI around the synthetic panel estimates overlap with those based on the actual panels; the "Coverage of 100%" row shows a similar figure for the number of times that the former fall completely inside the latter.



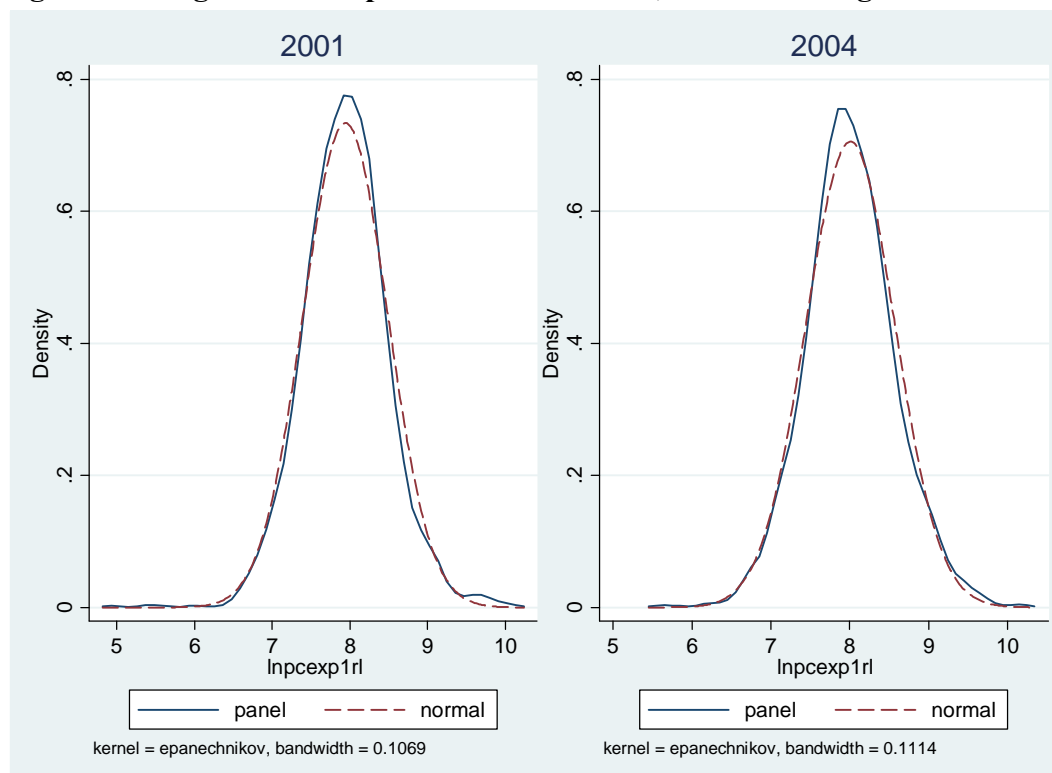
### Appendix 3: Data Appendix

#### Bosnia-Herzegovina

We use two rounds of the panel data for the Bosnia-Herzegovina Living Standards Measurement Survey (also known as Living in BiH Survey) in 2001 and 2004, which is publicly available on the World Bank LSMS website. We build our data based on the files made available by Demirguc-Kunt, Klapper and Panos (2011).

There are 2,376 panel households between 2001 and 2004. After restricting household heads' age to between 25 and 55 for the first survey round and adjust accordingly for the second round, and keeping cases with non-missing observation for the modelling variables, we are left with 1,342 panel households for analysis. We implement our analysis on the two halves of these panel data pretending that they are two cross sections. Figure 3.1 below provides the density graphs for household consumptions for Bosnia-Herzegovina in 2001 and 2004.

**Figure 3.1. Log of Consumptions for Panel Data, Bosnia-Herzegovina 2001-2004**



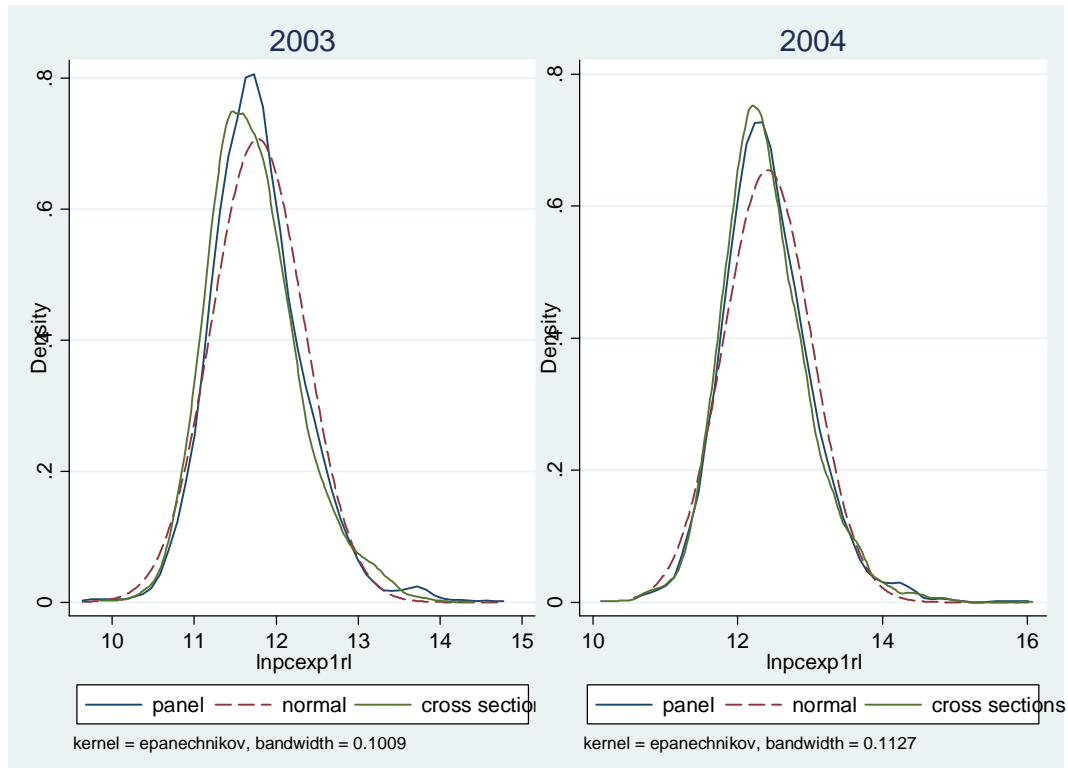
#### Lao PDR

We use two rounds of the Lao Expenditure and Consumption Survey (LECS) in 2002/03 and 2007/08, which is provided to us by the World Bank office in Lao PDR.

There are 2,357 panel households between 2002/03 and 2007/08. After restricting household heads' age to between 25 and 55 for the first survey round and adjust accordingly for the second round (i.e., increasing this age range to 30-60), we are left with 1,989 panel households for analysis. The corresponding numbers of cross sectional households with non-missing observation for the modelling variables we analyze are 3,032 and 3,215 respectively for 2002/03 and 2007/08.

Two sample t-tests with unequal variances show that household consumptions in the panel component are statistically but negligibly higher than those in cross section component for the 2002/03 round (e.g., with a difference of 0.05 between the two means of 11.78 and 11.73 respectively for the panel and cross section); however, these consumptions are not statistically different at the 5% level for the 2007/08 round. Figure 3.2 below provides as an example the density graphs for household consumptions in the panel and cross section components for Lao PDR in 2002/03 and 2007/08.

**Figure 3.2. Log of Consumptions for Panel and Cross Section Components, Lao PDR 2002/03- 2007/08**



## Peru

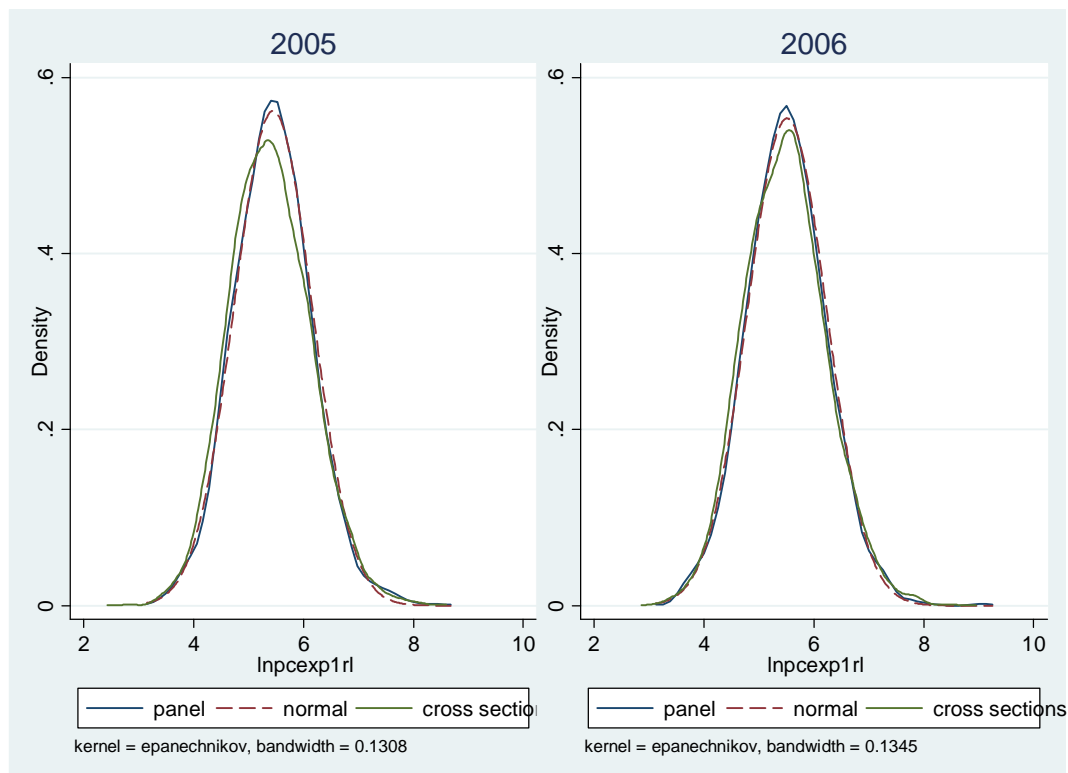
We use three rounds of the Peruvian National Household Survey (ENAHO) in 2004, 2005, and 2006, which is publicly available on the Peruvian Statistics Bureau (INEI)'s website. Both the panel and cross sectional households constructed from the ENAHOs are graciously provided to us by Renos Vakis and Leonardo Lucchetti based on their paper (Cruces et al., forthcoming).

There are 3,247 and 3,559 panel households respectively between 2004-2005 and 2005-2006. After restricting household heads' age to between 25 and 55 for the first survey round and adjust accordingly for the second round (i.e., increasing this age range to 26-56), we are left with 2,087 and 2,250 panel households for analysis. The corresponding numbers of cross sectional households with non-missing observation for the modelling variables we analyze are 4,493 and 9,169 respectively for 2004 and 2005 for the survey pair 2004-2005, and 8,593 and 9,084 respectively for 2005 and 2006 for the survey pair 2005-2006.

Two sample t-tests with unequal variances show that household consumptions in the panel component are statistically but negligibly higher than those in cross section component for the survey pair 2005-2006 (e.g., with a difference of 0.05 between the two means of 5.39 and 5.44 respectively for the panel and cross section); however, these consumptions are not statistically different for the survey pair 2004-2005. Figure 3.3 below provides as an example the density graphs for household consumptions in the panel and cross section components for Peru in 2006.

For all three years 2004-2005-2006, we have 2,668 panel households which is reduced to 1,987 panel households after a similar restriction on heads' age ranges. The corresponding number of cross sectional households we analyze is 8,608 for 2006.

**Figure 3.3. Log of Consumptions for Panel and Cross Section Components, Peru 2005-2006**



### United States

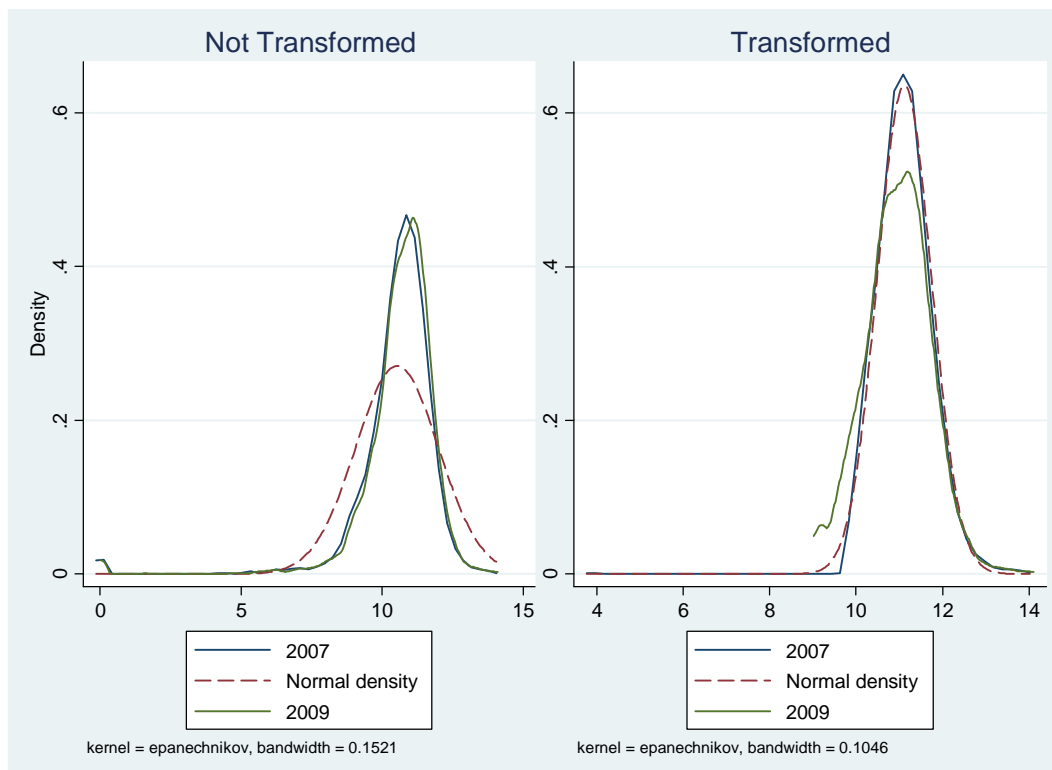
We use the three most recent rounds of the Panel Study of Income Dynamics (PSIDs) in 2005, 2007, and 2009, which is publicly available on the University of Michigan Institute for Social Research's website. The PSID started in 1968 and is the longest-running panel household survey implemented in the United States. The PSID was implemented annually between 1968 and 1997, and biannually after 1997. A useful documentation is provided in the PSID Main Interviewer User Manual Release 2012.1.

We use the sample persons in the PSID (i.e., those with a positive longitudinal weight), and after restricting household heads' age to between 25 and 55 for the first survey round and adjust accordingly for the second round (i.e., increasing this age range to 27-57) and keeping

the age difference across the two survey rounds between one and three,<sup>41</sup> we are left with 3,275 and 3,368 panel households respectively in 2005-2007 and 2007-2009 for analysis. For all three years 2005-2007-2009, we have 3,036 panel households after the restriction on heads' age range.

Since no comprehensive consumption aggregates are available in the PSID, we use income as a measure of households' welfare. Different from consumption measures, there are two potential issues with income measures: one is that the latter can be zero or negative (even though there are generally less than one percent of households at this welfare level in the PSID), which will translate into missing observations when we take the logarithm; the other is that income measures can have a lop-sided distribution that not as close to the normal distribution as consumption measures. We thus deal with both of these issues by implementing a Box-Cox transformation on the income variables using the *lnskew0* command in Stata, which effectively adds a positive constant  $k$  to the income before taking logarithm to minimize the skewness of the income variables. This constant  $k$  is 39,077, 19,727, and 8,279 respectively for incomes in 2005, 2007, and 2009.<sup>42</sup> Figure 3.4 below provides as an example the density graphs for log of incomes in 2007 and 2009 before and after the Box-Cox transformation.

**Figure 3.4. Log of Incomes Before and After Box-Cox Transformation, USA 2007-2009**



## Vietnam

We use three rounds of the Vietnam Household Living Standards Surveys (VHLSSs) in 2004, 2006, and 2008, which is provided to us by Vietnam's General Statistical Office (GSO)

<sup>41</sup> This helps ensure the household heads remain the same across the two surveys.

<sup>42</sup> Since the number of panel observations change slightly between the two pairs of survey years, this constant  $k$  also changes slightly for 2007.

and the World Bank office in Vietnam. The VHLSSs have a rotating panel design with around one half of the data for the 2004 round of the VHLSS is repeated in the 2006 round, and one half of the 2006 round consisting of one half of the 2004-2006 panel and one half of the 2006 cross section is repeated in the 2008 round. An introduction to the VHLSSs (with a focus on the years 2002 and 2004) is provided by Tung and Phong (undated).

We construct panel data for the VHLSSs using household identification codes. Where we suspect mismatching between panel households due to incorrect identification codes, we double check and correct these cases with household heads' names. As a result, we could match 4,276 panel households between 2004 and 2006 out of 9,189 households for each year. After restricting heads' age to between 25 and 55 for 2004 and 27 and 57 for 2006 and keeping the age difference across the two survey rounds between one and three, we are left with 2,723 panel households for analysis.

Following a similar procedure, we could match 4,120 panel households between 2004 and 2006 out of 9,189 households for each year. After placing the restriction on heads' age range and difference, we are left with 2,703 panel households for analysis. The corresponding numbers of cross sectional households we analyze are 3,527 and 3,674 respectively for 2004 and 2006 for the survey pair 2004-2006, and 3,596 and 3,701 respectively for 2006 and 2008 for the survey pair 2006-2008.

Two sample t-tests with unequal variances show that household consumptions in the panel component are statistically but negligibly less than those in cross section component in the first round for each of the survey pairs 2004-2006 and 2006-2008 (e.g., with a difference of 0.04 between the two means of 8.44 and 8.48 respectively for the panel and cross section); however, these consumptions are not statistically different for the second rounds. Figure 3.5 below provides as an example the density graphs for household consumptions in the panel and cross section components for Vietnam in 2008.

The numbers for the panel households we could match between 2004, 2006, and 2008 are 1,848 (before restriction) and 1,282 (after restriction). The corresponding number of cross sectional households we analyze is 3,808 for 2008.

**Figure 3.5. Log of Consumptions for Panel and Cross Section Components, Vietnam 2006-2008**

