

Estimating Poverty and Inequality from Grouped Data: How Well Do Parametric Methods Perform?¹

Camelia Minoiu² and Sanjay G. Reddy³

Abstract. Poverty and inequality are often estimated from grouped data as complete household surveys are neither always available to researchers nor easy to analyze. In this study we assess the performance of functional forms proposed by Kakwani (1980a) and Villasenor and Arnold (1989) to estimate the Lorenz curve from grouped data. The methods are implemented using the computational tools POVCAL and SimSIP, developed and distributed by the World Bank. To identify biases associated with these methods, we use unit data from several household surveys and theoretical distributions. We find that poverty and inequality are better estimated when the true distribution is unimodal than multimodal. For unimodal distributions, biases associated with poverty measures are rarely larger than one percentage point. For data from multi-peaked or heavily skewed distributions, the biases are likely to be higher and of unknown sign.

Keywords: grouped data, Lorenz curve, poverty, inequality, income distribution, POVCAL, SimSIP

JEL Classifications: C13, C14, C15, C16, D31, D63, I32

Wordcount: 7,001

¹ We are grateful for the support of the Bureau of Development Policy and the International Poverty Centre of the United Nations Development Program. We would like to thank Terry McKinley for facilitating that support. We would also like to thank Shaohua Chen, Gaurav Datt, and Martin Ravallion for helpful comments, Sergey Kivalov for assisting us with software development, and Catherine Choi and Prabhjot Kaur for research assistance. Income data from the 2004 Brazilian National Household Sample Survey was generously provided by Marcelo Medeiros of the Instituto de Pesquisa Econômica Aplicada (IPEA). We are grateful to the editor and two anonymous referees for useful comments.

² African Department, International Monetary Fund. Email: CMinoiu@imf.org. Tel. 202-623-9731. Fax 202-623-6947.

³ Dept. of Economics, Barnard College and School of International and Public Affairs, Columbia University. Email: sr793@columbia.edu. Tel. 212-854-3790, Fax. 212-854-8947.

1 Introduction

Poverty and inequality are often estimated from grouped data (i.e., mean incomes of population quantiles such as quintiles or deciles)⁴ for two reasons: first, complete household surveys are not always available to researchers. Second, the analysis of unit data is often labour and time-intensive. Consequently, estimates of regional and global poverty and inequality are often derived from grouped data (Yotopoulos, 1989; Bourguignon and Morrison, 2002; Chen and Ravallion, 2001, 2004; Sala-i-Martin, 2006; Ackland, Dowrick, and Freyens, 2007). In particular, data in summary form has been the sole source of information on income distributions of large countries such as China which greatly influence the extent and trend of global poverty (Reddy and Minoiu, 2007).⁵ Data on the size distribution of firms is also often published by governments only in grouped form (Golan, Judge, and Perloff, 1996) as is much historical data (Shorrocks and Wan, 2008).

In applied work, poverty and inequality are estimated from grouped data almost exclusively by using POVCAL and SimSIP, two software programs developed and distributed by the World Bank.⁶ The computational tools are essentially identical except that SimSIP provides an Excel-based user-friendly interface, while POVCAL operates in MS-DOS. Both programs have been used extensively in poverty analysis (Belkacem and Limam, 2004; Bhalla 2002; Chen and Ravallion 2001, 2004, 2006; Chen and Wang, 2001; Figini and Santarelli, 2006; Karshenas, 2004; Pritchett 2006; Son and Kakwani, 2006; Woo et al, 2004). POVCAL-based Gini coefficients have been included in the UNU-WIDER World Income Inequality Database (2005) and the World Bank's Measuring Income Inequality Database (Deininger and Squire, 1996). A large number of cross-country econometric analyses have subsequently been undertaken using these databases (e.g., Lundberg and Squire, 2003; Banerjee and Duflo, 2003; Milanovic, 2002; Forbes, 2000; Easterly, 1999; Deininger and Squire, 1998). Finally, POVCAL and SimSIP have been widely used in the preparation of national poverty assessments (for instance Ali and Elbadawi, 2008; Asra, 2002; Belkacem, 2001; Eele et al, 2002; Acharya, 2004; Joeques et al, 2002) and have been widely recommended to practitioners by development agencies (USAID, 2004; World Bank, 2003).

Despite their widespread use, few studies have provided a systematic analysis of the grouped data methods used by POVCAL and SimSIP. This paper attempts to fill this gap. We analyze the accuracy of estimates based on two Lorenz curve functional forms used in the programs. These functional forms were developed respectively by Villasenor and Arnold (1989) and Kakwani (1980a), and are also known as the Generalized Quadratic (GQ) and Beta parameterizations. We compare estimates of Lorenz curves, poverty, and inequality derived from grouped data using this approach with those estimated directly from unit data, for a wide range of income

⁴ In what follows, we bear in mind that poverty and inequality analysis can be applied to distributions of consumption, income, wealth, or other dimensions of personal advantage. However, without loss of generality we refer to *income* distributions throughout the paper.

⁵ China's State Statistics Bureau has not made full household survey data available to outside researchers.

⁶ Some of the functions of POVCAL and SimSIP are built in to the PovcalNet website of the World Bank (<http://iresearch.worldbank.org/PovcalNet/jsp/index.jsp>) which allows users to obtain poverty estimates for any country in the period covered for a specified poverty line. The two programs are briefly described in the Appendix.

distributions. We undertake both Monte Carlo simulations as well as deterministic comparisons. Our data sources comprise: household surveys from five countries (Brazil, China, Nicaragua, Tanzania, and Vietnam) and data generated from two distributions (the theoretical Dagum distribution and an empirical global income distribution).

We find that the two functional forms perform relatively well in estimating poverty from unimodal distributions. Larger biases were identified in the case of multi-peaked distributions. We also find that the biases vary (albeit not systematically) with the sample size, functional forms, distributions, poverty lines, and poverty indicators. Inequality (measured by the Gini coefficient) is accurately estimated in most cases considered. While we have attempted to analyze a wide range of possible income distributions, we caution that our results should be regarded as conditional on the data used and may therefore not hold for data derived from a different underlying distribution (or data generating process).

The remainder of the paper is structured as follows. Section 2 briefly discusses previous assessments of various techniques for parametrically estimating Lorenz curves from grouped data. In section 3 we present the biases identified based on Monte Carlo simulations. Section 4 presents our findings in the case of deterministic comparisons using household survey data. Conclusions are drawn in Section 5. Plots of all distributions used in the paper are presented in Figures 1 and 2 in the Appendix.

2 Previous studies

A limited literature examines the theoretical validity and empirical performance of alternative Lorenz curve functional forms. Villasenor and Arnold (1989) used the 1967–68 Australian Survey of Consumer Expenditure to find that GQ is superior to three alternative parameterizations (Kakwani and Podder, 1976; Pakes, 1981; and the classical Pareto distribution) based on the sum of squared errors over the entire support. Larger samples were associated with a better fit. Although the best fit of the Lorenz curve was achieved for unimodal distributions, the authors judged their parameterization to be satisfactory for bimodal income distributions as well. Kakwani (1980a) assessed the goodness-of-fit of the Beta functional form using the 1974 Australian Household Expenditure Survey. The R-squared statistics from the estimating regressions (not reported in the paper) were close to 0.99, while the fitted Lorenz curve values were within two decimal places of the survey values.

Ravallion and Huppi (1989) used household consumption data for 50,000 Indonesian households and compared Lorenz curve estimates obtained from three functional forms: Villasenor and Arnold (1989), Kakwani and Podder (1976), and Kakwani (1980a). They found that the worst fit was provided by the two-parameter specification of Kakwani and Podder (1976), while the other two specifications—also the subject of our study—gave broadly similar results.⁷ The GQ

⁷ Kakwani and Podder (1976) discuss the goodness of fit of their two-parameter specification. An empirical exercise which they undertake (using data from the 1967–68 Australian Survey of Consumer Expenditures and Finance) reveals underestimation of the mean income of the poorest 5 percent of the population, and overestimation of the mean income of the poorest 10 percent. Anand (1983), and Anand and Kanbur (1993a, 1993b) note that there are reasons that Kakwani and Podder's (1976) proposed functional form for parametric estimation violates the

parameterization provided a better fit towards the high end of the income distribution, while the Beta form did better in the left tail.

In relation to inequality, Cheong (2002) undertook a comparison of four Lorenz curve functional forms (Kakwani and Podder, 1976; Rasche et al, 1980; Kakwani, 1980a; and Ortega et al, 1991) in estimating the U.S. income Gini (for one hundred income classes). The author concluded that Kakwani's Beta form, although theoretically invalid (on which point see Anand 1983; Anand and Kanbur, 1993a, 1993b; Rasche et al, 1980; and Ortega et al, 1991), provided as good a fit to the data as did the theoretically valid parameterization proposed by Rasche et al (1980). In an earlier study, Schader and Schmid (1994) used household survey data to compare Gini coefficients obtained through parametric Lorenz curve estimation with nonparametric Gastwirth bounds for the true Gini coefficient (Gastwirth, 1972). Estimates of inequality based on the Kakwani (1980a) parameterization of the Lorenz curve were found to lie between the nonparametric Gastwirth bounds for all datasets considered.

While many Lorenz curve functional forms have been proposed in the literature (other examples include Gastwirth and Glaubergerman, 1976; Gupta, 1984; Mazzarino, 1986; Basmann et al, 1990; Ogwang and Rao, 1996), an exhaustive assessment of all the alternatives is not the object of this study. The main reason for restricting our attention to the GQ and Beta forms is their almost exclusive use in the estimation of poverty and inequality from grouped data due to their use in the World Bank's POVCAL and SimSIP programs.

3 Findings from Monte Carlo simulations

In this section, we describe our Monte Carlo analysis of the performance of the GQ and Beta Lorenz curve parameterizations for grouped data. Our data are drawn first from a theoretical distribution—the Dagum—and second from a notional multimodal distribution.

Bandourian, McDonald and Turley (2003) show that the Dagum distribution is the best-fitting among three parameter distributions to survey income data. Jenkins and Cox (1999) note that the Dagum provides a good fit to income distributions largely due to its ability to model skewed distributions. We parameterize it with median values from the reported best-fitting Dagum parameters for recent income distributions from 27 countries (Bandourian, McDonald and Turley, 2003). These median parameter values happen to be closest to those fit by the authors for Russia's 1992 distribution.⁸

The multimodal distribution is the 2004 population-weighted world distribution of income, in which every individual is assigned the per capita PPP-adjusted GDP of her country (as reported in the World Development Indicators, 2006). The two higher peaks of the distribution

theoretical requirements for a valid Lorenz Curve. Dhongde (2004) derived the small sample bias of Lorenz curve estimates associated with the earlier parameterization of Kakwani and Podder (1973).

⁸ The Dagum distribution has the following parameter values: $a = 2.742$, $b = 100,000$ and $c = 0.337$.

correspond to population mass concentrated at China's and India's per capita GDP, while the lower peak corresponds to the population mass of the rich nations (Figure 1).

From each hypothesized density, we draw 100 random samples of 1,000 observations.⁹ We restrict the number of draws to 100 due to the high volume of manual work involved in running the software on these samples. Each sample is then collapsed into grouped data (quintile, decile and ventile means) and entered into POVCAL. We do not consider cases beyond ventiles since in practice at most twenty datapoints are typically available to researchers. Furthermore, we use multiple poverty lines in order to identify biases at multiple points along the support of the curve, including money-metric international poverty lines (\$1/day and \$2/day), nutritionally-anchored poverty lines (constructed by Reddy, Visaria, and Asali 2008), and thresholds representing the population or survey median multiplied by various constant proportions.

3.1 The Lorenz curve

We first compare the Lorenz curve estimates from grouped data with the true curve. The goodness-of-fit is assessed both along the entire curve and in the left tail (up to a poverty headcount ratio of 20 percent) using the sum of squared errors (SSE) and the sum of absolute errors (SAE). We find that a higher number of quantile means implies a lower SSE and SAE (for the GQ parameterization) but (surprisingly) the reverse is the case for the Beta parameterization (Table 1). If the SSE and SAE are computed up to a headcount ratio of 20 percent, the Beta parameterization underperforms the GQ method, except in the case of quintile data.

Average Lorenz curve estimates for population proportions up to 10 percent (Table 2) give a fine-grained indication of the biases in the left tail of the distribution. A series of interesting patterns arise: for example, the GQ parametrization overestimates the income share accruing to each population proportion toward the left side of the support. In contrast, the Beta parameterization yields negative, and hence invalid, average Lorenz curve estimates for the bottom population centiles. POVCAL alerts the user to the invalidity of the Lorenz curve each time negative income shares are fitted.

The Lorenz curve is consistently overestimated across the support (Table 3). As a result, the true Lorenz curve is dominated by the estimated one (Figure 3), which implies that distortions in the Lorenz curve arise along the entire support and any Lorenz-consistent measure of inequality derived from these estimates would under-estimate inequality. The magnitude of the biases is very similar across sample sizes and estimation methods.

We repeated the exercise for the notional multimodal distribution. Using the SAE criterion, the GQ parametrization provides a worse fit. Interestingly, in larger samples, the goodness-of-fit does not vary monotonically with the number of datapoints for the Beta parameterization (Table 4). Furthermore, both functional forms occasionally give rise to invalid estimated Lorenz curves

⁹ In the case of the Dagum distribution, we draw 100 random samples from a universe of 1,000,000 observations. For the multimodal distribution, we draw 100 random samples from a universe of almost 600,000 observations each representing 1/10,000 of the world's population.

(Table 5). Finally, both methods underestimate the Lorenz curve (unlike in the case of the Dagum distribution) (Table 6).

Figure 4 offers a visual representation of these findings by superimposing the true and fitted Lorenz curves. Since both positive and negative biases occur along the support, the estimated and true Lorenz curves cross. It follows that whether a Lorenz-consistent inequality index will over- or underestimate the true level of inequality from these estimates depends on specific features of the index.

3.2 Poverty and inequality

We find that poverty and inequality are better estimated for the Dagum distribution than for the multimodal distribution (Tables 7–9). For the Dagum distribution, the two parameterizations perform exceptionally well. Specifically, the average bias associated with various poverty indicators is rarely higher than one percentage point (Table 7). Furthermore, the GQ parameterization slightly outperforms the Beta parameterization. As before, a larger sample is not always associated with a lower bias.

In contrast, the biases are often sizable for the multimodal distribution and larger samples are frequently associated with larger biases (Table 8).¹⁰ In addition, the sign of the biases changes from positive (for quintile means) to negative (for decile and ventile means) and as the poverty line rises. The Gini coefficient is more accurately estimated by the Beta parameterization, regardless of the distribution (Table 9). Biases are extremely small for the Dagum distribution, but as large as 5 percentage points (Dagum) and 3.7 percentage points (multimodal) in samples of quintiles.

From these results, it is difficult to identify any regularities other than that the biases are larger when the true distribution is multimodal. Indeed, Villasenor and Arnold (1989) themselves note that the GQ parameterization provides a good fit to data from the unimodal family of densities, but less so to data with bimodal histograms.

4 Findings from household surveys

In this section we report on the performance of the two functional forms using unit data from five household surveys (Tables 10–12).¹¹ We use disposable income and consumption

¹⁰ We often cannot report poverty biases for the \$2/day poverty line due to the frequency with which POVCAL shut down, failed to write to the output files, provided meaningless output (e.g., higher than one poverty headcount ratios) or generated infeasible bounds for the poverty lines. The biases are thus computed across successful program runs. For more examples on technical problems encountered when running POVCAL, see Minoiu and Reddy (2007a).

¹¹ The findings presented in this section are consistent with those from analyzing theoretical distributions such as Weibull, Log-normal, Pareto, and Generalized Beta II. (The results are reported in Minoiu and Reddy, 2007a.)

distributions from nationally representative surveys of Brazil, China, Tanzania, Nicaragua, and Vietnam. A description of the data sources and variables is provided in the Appendix.

For many surveys, the biases in poverty estimates are very small: they are seldom larger than one percentage point (in either direction). The poverty headcount ratio in particular is generally estimated well. As shown in Table 11, a marked difference is however observed when comparing the cases of China (a well-behaved distribution) and Brazil (a multi-peaked, heavily skewed distribution)¹². For lower poverty lines, the poverty headcount ratio is substantially more underestimated for Brazil, with biases as high as 5 to 14 percentage points. The two parameterizations give rise to concerning results for the poverty gap ratio and the squared poverty gap as well. Overall, the magnitude of the biases and the manner in which they change sign depending on the sample size, poverty line, and poverty indicator gives rise to concern regarding the use of grouped data from underlying multi-peaked, highly unequal distributions.

As before, inequality is better estimated than poverty (Table 12). While our chosen household surveys have Gini coefficients ranging between 35 (Vietnam) and 71 (Brazil), we find that in all cases considered, the biases are negligible. For Brazil, the Beta functional form underestimates the Gini index by at most 2 percentage points. Generally, inequality is well estimated by both parameterizations.

5 Discussion and conclusions

In this paper we have analyzed the biases associated with Lorenz curve, poverty, and inequality estimates obtained from grouped data using two parameterizations: Villasenor and Arnold (1989) and Kakwani (1980a). These estimation methods have been widely used through the World Bank's POVCAL and SimSIP computational tools. They have also been the basis for many regional and global assessments of poverty and inequality. We have used data drawn from a wide range of income distributions with single and multiple modes, and have undertaken both Monte Carlo simulations as well as deterministic comparisons.

We found that the two parameterizations perform relatively well in estimating poverty and inequality for unimodal distributions. Larger biases were identified, however, in the case of the multimodal distribution considered. The extent of misestimation of poverty does not vary predictably with the level of poverty line, the poverty indicator, or even the sample size. Inequality (measured by the Gini index) is usually well estimated.

We often found that the two parameterizations yielded negative fitted income shares, giving rise to invalid estimated Lorenz curves. Should one use or discard the poverty and inequality output associated with them? We found no straightforward correspondence between the validity of the Lorenz curve and the quality of the poverty and inequality estimates. Notably, Kakwani (1980b) defended a parameterization which had been shown to give rise to theoretically invalid Lorenz

¹² We employ for Brazil a survey that includes (rather than discards) zero values for income, recognizing that this is only one plausible treatment of the underlying data.

curve estimates by arguing that its *overall* superior performance in fitting income distributions was a sufficient basis for its use.

The results presented in this study offer qualified support for the use of existing software such as POVCAL and SimSIP in poverty and inequality analysis based on grouped data. It should be noted, however, that other techniques have also been recently proposed to analyze such data. For example, Minoiu and Reddy (2007b) assessed kernel density estimators and concluded that Lorenz curves functional forms often outperformed kernel methods in poverty analysis. Wu and Perloff (2003) evaluated the maximum entropy density estimator for grouped data and found that it performed well. Similarly, Chotikapanich et al (2007) have shown the the Generalized Beta 2 distribution is a good parametric choice. Since grouped data will continue to be an important source of information—and often, the only one—in poverty and inequality analysis, future research should aim to more fully establish the relative empirical performance of these alternatives and to develop new methods.

References

- Acharya, S. (2004) “Measuring and Analyzing Poverty (with particular reference to the case of Nepal)”, *European Journal of Comparative Economics*, Vol. 1(2), pp. 195–215.
- Ackland, R., Dowrick, S. and Freyens, B. (2007) “Measuring Global Poverty: Why PPP Methods Matter”, forthcoming, *Economic Journal*.
- Ali, A.A.G. and Elbadawi, I.A. (2008) “Explaining Sudan’s Economic Performance”, in *The Political Economy of Economic Growth in Africa, 1960–2000*, Volume 2, Country Case Studies (eds. Benno J. Ndullu, Stephen A. O’Connell, Jean-Paul Azam, Robert H. Bates, Augustin K. Fosu, Jan Willem Gunning, and Dominique Nijinkeu), Cambridge University Press, Cambridge, U.K.
- Anand, S. (1983) *Inequality and Poverty in Malaysia. Measurement and Decomposition*, Volume 1, International Bank for Reconstruction and Development / World Bank, Oxford University Press, Washington, D.C.
- Anand, S. and Kanbur, S.M.R. (1993a) “Inequality and Development: A Critique”, *Journal of Development Economics*, Vol. 41, pp. 19–43.
- Anand, S. and Kanbur, S.M.R. (1993b) “The Kuznets Process and the Inequality-Development Relationship”, *Journal of Development Economics*, Vol. 40, pp. 25–52.
- Asra, A. (2000) “Poverty and Inequality in Indonesia”, *Journal of the Asia Pacific Economy*, Vol. 5, pp. 91–111.
- Basman, R.L., Hayes, K.J., Slottje, D.J., and J.D. Johnson (1990) “A General Functional Form for Approximating the Lorenz Curve”, *Journal of Econometrics*, Vol. 43, pp. 77–90.
- Bandourian, R., McDonald, J.B. and Turley, R.S. (2003) “A comparison of parametric models of income distributions across countries and over time”, *Revista Estadística*, Vol. 55, pp. 164–165.
- Banerjee, A.V. and Duflo, E. (2003) “Inequality and growth: what can the data say?”, *Journal of Economic Growth*, Vol. 8, pp. 267–299.
- Belkacem, L. (2001) “Poverty Dynamics in Algeria”, API Working Paper No. 0103 (Kuwait: Arab Planning Institute).
- Belkacem, L. and Limam, I. (2004) “Impact of Public Policies on Poverty, Income Distribution and Growth”, API Working Paper No. 0401 (Kuwait: Arab Planning Institute).
- Bhalla, S. (2002) *Imagine There’s No Country: Poverty, Inequality, and Growth in the Era of Globalization*, IIE (Washington: Institute for International Economics).

- Bourguignon, F. and C. Morrisson (2002) “Inequality among World Citizens: 1820–1992”, *American Economic Review*, Vol. 92 No. 4.
- Chen, S., Datt, G. and Ravallion, M. (2001) POVCAL, A program for calculating poverty measures for grouped data, World Bank Poverty and Human Resource Division (Washington: The World Bank).
- Chen, S. and Ravallion, M. (2001) “How Did the World’s Poorest Fare in the 1990s?”, *Review of Income and Wealth*, Vol. 47(3), pp. 283–300.
- Chen, S. and Ravallion, M. (2004) “How Have the World’s Poorest Fared Since the Early 1980s?”, World Bank Development Research Group Working Paper No. 3341 (Washington: The World Bank) and *World Bank Research Observer*, Vol. 19(2), pp. 141–169.
- Chen, S. and Ravallion, M. (2006) “China’s (Uneven) Progress against Poverty”, *Journal of Development Economics*, Vol. 82(1), pp. 1–42.
- Chen, S. and Wang, Y., (2001) “China’s growth and poverty reduction: recent trends between 1990 and 1999”, World Bank Policy Research Working Paper No. 2651 (Washington: The World Bank).
- Cheong, K.S. (2002) “A Comparison of Alternative Functional Forms for Parametric Estimation of the Lorenz Curve”, *Applied Economics Letters*, Vol. 9, Issue 3, pp. 171–176.
- Chotikapanich, D., Griffiths, W.E., and Rao, D.S. P. (2007) “Estimating and Combining National Income Distributions using Limited Data”, *Journal of Business and Economic Statistics*, Vol. 25, pp. 97–109.
- Datt, G. (1998) “Computational Tools for Poverty Measurement and Analysis”, IFPRI Food Consumption and Nutrition Division Discussion Paper No. 50 (Washington: International Food Policy Research Institute).
- Deininger, K. and Squire, L. (1998) “New Ways of Looking at Old Issues: Inequality and Growth”, *Journal of Development Economics*, Vol. 57(2), pp. 259–287.
- Deininger, K. and Squire, L. (1996) “A new data set measuring income inequality”, *World Bank Economic Review*, Vol. 10(3), pp. 565–591.
- Dhongde, S. (2004) “On the Bias of Estimating Lorenz Curve Parameters”, unpublished manuscript, Department of Economics, University of California, Riverside.
- Easterly, W. (1999) “Life during growth”, *Journal of Economic Growth*, Vol. 4(3), pp. 239–276.
- Eele, G., Semboja, J., Likwelile, S., and Ackroyd, S. (2000) “Meeting International Poverty Targets in Tanzania”, *Development Policy Review*, Vol. 18(1), pp. 63–83.

Figini, P. and Santarelli, E. (2006) “Openness, Economic Reforms, and Poverty: Globalization in the Developing Countries”, *Journal of Developing Areas*, Vol. 39(2), pp. 129–151.

Forbes, K.J. (2000) “A reassessment of the relationship between inequality and growth”, *American Economic Review*, Vol. 90(4), pp. 869–887.

Gastwirth, J.L. (1972) “The Estimation of the Lorenz Curve and Gini Index”, *Review of Economics and Statistics*, Vol. 54(3), pp. 306–316.

Gastwirth, J.L. and Glauber, M. (1976) “The Interpolation of the Lorenz Curve and Gini Index from Grouped Data”, *Econometrica*, Vol. 40, pp. 479–483.

Golan, A., Judge, G., and Perloff, J.M. (1996) “Estimating the Size Distribution of Firms Using Government Summary Statistics”, *Journal of Industrial Economics*, Vol. 44, pp. 69–80.

Gupta, M.R. (1984) “Functional Form for Estimating the Lorenz Curve”, *Econometrica*, Vol. 52(5), pp. 1313–1314.

Jenkins, S.P and Cox, N.J. (1999) “DAGUMFIT: Stata module to fit a Dagum distribution to unit record data”, Statistical Software Components Paper No. S366101 (Boston: Boston College, Department of Economics).

Joekes, S., Ahmed, N., Ercelawn, A. and Zaidi, S. A. (2000) “Poverty Reduction without Human Development in Pakistan: Money Doesn’t Buy You Everything”, *Development Policy Review*, Vol. 18(1), pp. 37–62.

Kakwani, N. C. (1980a) “On A Class of Poverty Measures”, *Econometrica*, Vol. 48, Issue 2, pp 437–446.

Kakwani, N. C. (1980b) “Functional Forms for Estimating the Lorenz Curve: A Reply”, *Econometrica*, Vol. 48, Issue 4, pp 1063–1064.

Kakwani, N. C. and Podder, N. (1976) “Efficient Estimation of the Lorenz Curve and Associated Inequality Measures from Grouped Observations”, *Econometrica*, Vol. 44(1), pp 137–148.

Kakwani, N. C. and Podder, N. (1973) “On the Estimation of Lorenz Curves from Grouped Observations”, *International Economic Review*, Vol. 14, pp. 278–292.

Karshenas, M. (2004) “Global Poverty Estimates and the Millennium Goals: Towards a Unified Framework”, ILO Employment Strategy Paper No. 2004/5 (Geneva: International Labor Organization).

Lundberg, M. and Squire, L. (2003) “The simultaneous evolution of growth and inequality”, *The Economic Journal*, Vol. 113(4), pp. 326–344.

- Mazzarino, G. (1986) "Fitting the Distribution Curves to Grouped Data", *Oxford Bulletin of Economics and Statistics*, Vol. 48, pp. 189–200.
- Milanovic, B. (2002) "True World Income Distribution, 1988 and 1993: First Calculation Based on Household Surveys Alone", *Economic Journal*, Vol. 112, pp 51–92.
- Minoiu, C. and Reddy, S. (2007a) "The Assessment of Poverty and Inequality through Parametric Estimation of Lorenz Curves: An Evaluation", ISERP Working Paper No. 2007–02 (New York: Columbia University, Institute for Social and Economic Research and Policy).
- Minoiu, C. and Reddy, S. (2007b) "Kernel Density Estimation Based on Grouped Data: The Case of Poverty Assessment". Available at SSRN: <http://ssrn.com/abstract=991503>
- Ogwang, T. and U.L.G. Rao (1996) "A New Functional Form for Approximating the Lorenz Curve", *Economics Letters*, Vol. 52, pp. 21–29.
- Ortega, P., Martin, G., Fernandez, A., Ladoux, M. and A. Garcia (1991) "A New Functional Form for Estimating Lorenz Curves", *Review of Income and Wealth*, Vol. 37, pp. 447–452.
- Pakes, A.G. (1981) "On Income Distributions and their Lorenz Curves", Department of Mathematics Technical Report (Nedlands: University of Western Australia).
- Pritchett, L. (2006) "Who Is Not Poor? Dreaming of a World Truly Free of Poverty", *The World Bank Research Observer*, Vol. 21(1), pp. 1–23.
- Ramadas, K., van der Mensbrugge, D., and Wodon, Q. (2002) "SimSIP Poverty: Poverty and Inequality Comparisons using Grouped Data" (Washington: The World Bank).
- Rasche, R.H., Gaffney, J., Koo, A.Y.C. and N. Obst (1980) "Functional Forms for Estimating the Lorenz Curve", *Econometrica*, Vol. 48(4), pp. 1061–1062.
- Ravallion, M. and Huppi, M. (1989) "Poverty and Under-nutrition in Indonesia during the 1980s", World Bank Agriculture and Rural Development Department, Policy, Planning and Research Working Paper No. 286 (Washington: The World Bank).
- Reddy, S. and Minoiu, C. (2007) "Has World Poverty Really Fallen?", *Review of Income and Wealth*, Vol 53(3), pp. 484–502.
- Reddy, S. and Minoiu, C. (2006) "Chinese Poverty: Assessing the Impact of Alternative Assumptions", ISERP Working Paper No. 06–04 (New York: Columbia University Institute for Social and Economic Research and Policy).
- Reddy, S., Visaria, S. and Asali, M. (2008) "Inter-Country Comparisons of Poverty Based on a Capability Approach: An Empirical Exercise", in *Arguments for a Better World: Essays in Honor of Amartya Sen* (eds. Kaushik Basu and Ravi Kanbur), Oxford University Press.

Sala-i-Martin, X. (2006) “The World Distribution of Income: Falling Poverty and Convergence, Period”, *Quarterly Journal of Economics*, Vol. 121(2), pp. 351–397.

Schader, M. and Schmid, F. (1994) “Fitting Parametric Lorenz Curves to Grouped Income Distributions—A Critical Note”, *Empirical Economics*, Vol. 19(3), pp. 361–370.

Shorrocks, A. and Wan, G., (2008) “Ungrouping Income Distributions: Synthesizing Samples for Inequality and Poverty Analysis”, in *Arguments for a Better World: Essays in Honor of Amartya Sen* (eds. Kaushik Basu and Ravi Kanbur), Oxford University Press.

Son, H.H. and Kakwani, N. (2006) “Global Estimates of Pro-Poor Growth”, UNDP–IPC Working Paper No. 31 (Brasilia: United Nations Development Programme, International Poverty Center).

USAID (2004) “Review of Poverty Assessment Tools—Implementation Training Materials”, USAID Accelerated Microenterprise Advancement Project and Developing Poverty Assessment Tools Project (Washington: United States Agency for International Development).

UNU-WIDER (2005) World Income Inequality Database V. 2.0a, June 2005.

Villasenor, J.A. and Arnold, B.C. (1989) “Elliptical Lorenz Curves”, *Journal of Econometrics*, Vol. 40, pp. 327–338.

Yotopoulos, P.A. (1989) “Distributions of Real Income: Within Countries and by World Income Classes”, *Review of Income and Wealth*, Vol. 35(4), pp. 357–376.

Woo, W.T., Shi, Li., Ximing, Y., Xiaoying, H.W., Xingpeng, X. (2004) “The Poverty Challenge for China in the New Millennium”, Report to the Poverty Reduction Taskforce of the Millennium Development Goals of the United Nations.

World Bank (2003) “A User’s Guide to Poverty and Social Impact Analysis”, Poverty Reduction Group (PRMPR) and Social Development Department (SDV) (Washington: The World Bank).

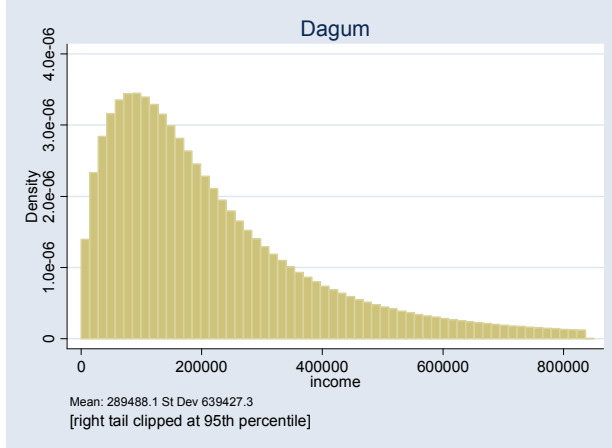
World Development Indicators online database (2006) (Washington: The World Bank).

Wu, X. and Perloff, J. (2003) “Calculation of Maximum Entropy Densities with Application to Income Distributions”, *Journal of Econometrics*, Vol. 115, pp. 347–354.

APPENDIX

Figure 1. Theoretical distributions used in Monte Carlo analysis

Dagum distribution



Multimodal (notional) distribution

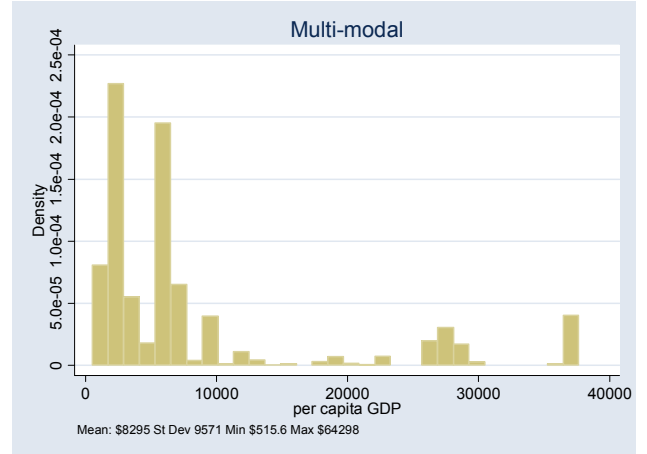
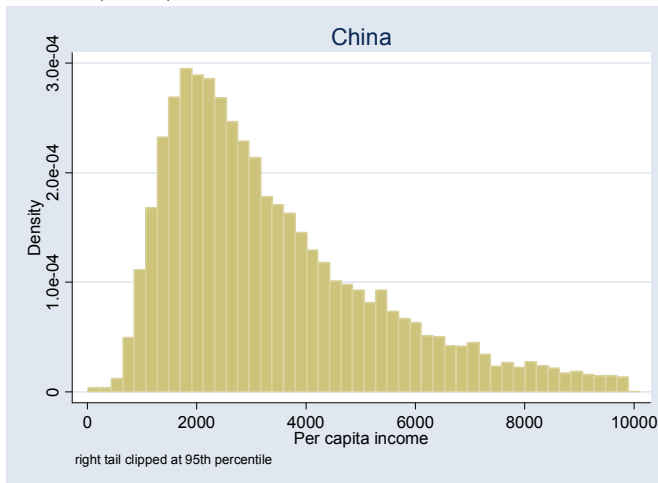
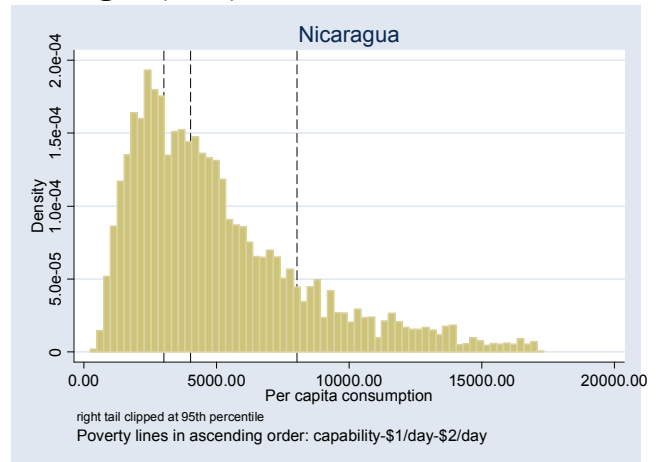


Figure 2. Household survey distributions used for deterministic comparisons

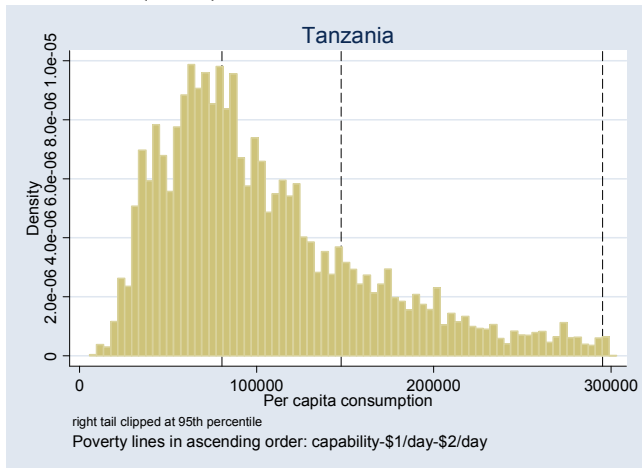
China (1995)



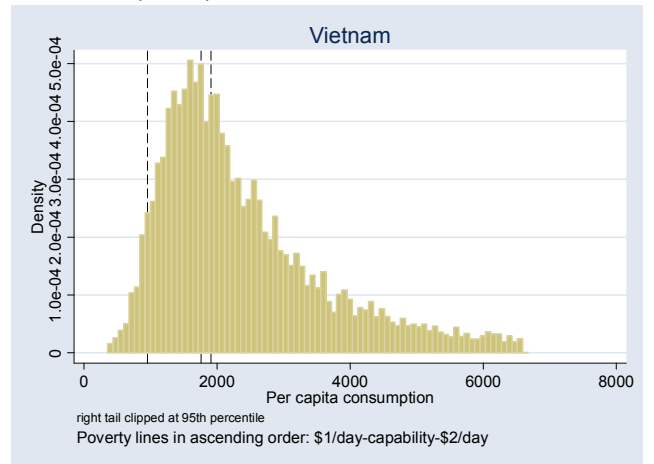
Nicaragua (1998)



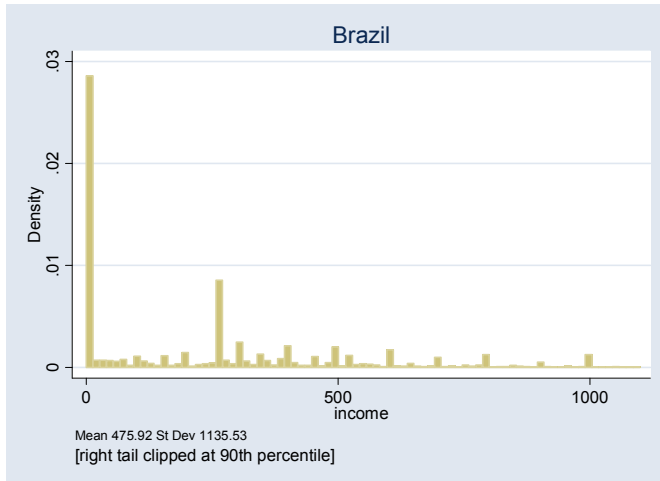
Tanzania (2001)



Vietnam (1998)



Brazil (2004)



Description of POVCAL and SimSIP

POVCAL (www.worldbank.org/lsms/tools/povcal) and SimSIP-Simulations for Social Indicators and Poverty (www.worldbank.org/simsip) are poverty assessment tools produced and distributed by the World bank. POVCAL functions in MS-DOS, whereas SimSIP is Excel-based. SimSIP has the additional features that it enables sector-level and decomposition analyses of poverty.

The Lorenz curve is estimated from grouped data by regression analysis based on two functional forms—the GQ parameterization of Villasenor and Arnold (1989) and the Beta parameterization of Kakwani (1980a)—each involving three parameters (Chen, Datt, and Ravallion, 2001; Datt, 1998). The grouped data read by the programs may take different forms (for example, income shares or mean incomes of population quantiles, share of the population in given income intervals, etc.). Datt (1998) and Ramadas, van der Mensbrugghe, and Wodon (2002) provide detailed formulas for poverty and inequality indices as functions of the parameters of the two functional forms.

For a user-specified poverty line, the output includes: the poverty headcount ratio (FGT0), the poverty gap index (FGT1), the squared poverty gap (FGT2), and the elasticity of these poverty measures with respect to the mean income (assuming a constant distribution). Also reported are the Gini coefficient of inequality, the Lorenz curve, and the parameterization which provides a better fit to the data.

The programs also report on the consistency of the estimated Lorenz curve with the requirements for a valid Lorenz curve. It can easily be shown algebraically that the Beta functional form *always* violates conditions required for the validity of the Lorenz curve (in particular by implying a negative slope at the origin). The GQ parameterization gives rise to valid Lorenz curves only under certain conditions on its parameters. Villasenor and Arnold (1989) find that the estimated GQ Lorenz curve is sometimes negative. We also find this to be the case, as the range of the regressors, and consequently the domain of the fitted values, are unrestricted in the estimation.

Description of income (consumption) variables in the household surveys

The 2004 Brazilian National Household Sample Survey (PNAD) is a representative household survey of the entire population of the country with the exception of remote areas in the Amazon. The data are weighted. The income variable represents *individual total earned income* (labor income and pension transfers). The two modes represent a high mass at zero (reported income by groups such as children and housewives) and the minimum wage (which is also the minimum value for pensions).

The 1995 Chinese Household Income Project is publicly available through the Inter-University Consortium for Political and Social Research, 2000. We pooled the rural and urban surveys to obtain the variable, which represents *per capita consumption* (with no adjustment for household composition). The data are not weighted. For a detailed description of this variable, see Reddy and Minoiu (2006).

The 1998 Vietnam Living Standards Survey (VLSS) contains information on *per capita expenditure* of households at current prices for 22,510 individuals. The data are weighted. Source: World Bank Living Standards Measurement Study (LSMS), Development Economics Research Group (DECRG), Washington, D.C.

The 1997-98 Nicaragua Living Standards and Measurement Survey (LSMS) contains information on *per capita consumption* for 18,383 individuals. The data are weighted. Source: World Bank Living Standards Measurement Study, Development Economics Research Group (DECRG), Washington, D.C.

The 2000-01 Tanzania Household Budget Survey contains information on *per capita consumption* for 22,176 households. The data are weighted. Source: National Bureau of Statistics, Tanzania, 2002.

**Monte Carlo simulation results for the Lorenz curve
(Dagum distribution)**

Table 1. Sum of Squared Errors (SSE) and Sum of Absolute Errors (SAE) of the Lorenz curve estimate from grouped data

	GQ			Beta		
	Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
<i>Along the entire support</i>						
SSE	0.3751	0.3706	0.3689	0.3637	0.3638	0.3620
SAE	5.3502	5.3165	5.3007	5.2458	5.2697	5.2704
<i>Up to the 20th percentile of the population</i>						
SSE	0.0021	0.0020	0.0020	0.0020	0.0024	0.0027
SAE	0.1622	0.1564	0.1526	0.1537	0.1707	0.1818

Table 2. Lorenz curve estimates from grouped data in the left tail of the distribution

Cum. pop. prop.	<i>True</i>	GQ			Beta		
		Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
0.01	0.0001	0.0001	0.0000	0.0000	-0.0016	-0.0010	-0.0007
0.02	0.0002	0.0005	0.0004	0.0003	-0.0017	-0.0008	-0.0004
0.03	0.0005	0.0012	0.0010	0.0008	-0.0011	-0.0001	0.0005
0.04	0.0009	0.0021	0.0018	0.0017	-0.0002	0.0009	0.0016
0.05	0.0014	0.0032	0.0029	0.0028	0.0011	0.0023	0.0030
0.06	0.0020	0.0046	0.0043	0.0041	0.0026	0.0039	0.0047
0.07	0.0028	0.0062	0.0059	0.0056	0.0044	0.0058	0.0066
0.08	0.0037	0.0080	0.0077	0.0074	0.0064	0.0078	0.0086
0.09	0.0047	0.0101	0.0097	0.0094	0.0087	0.0101	0.0109

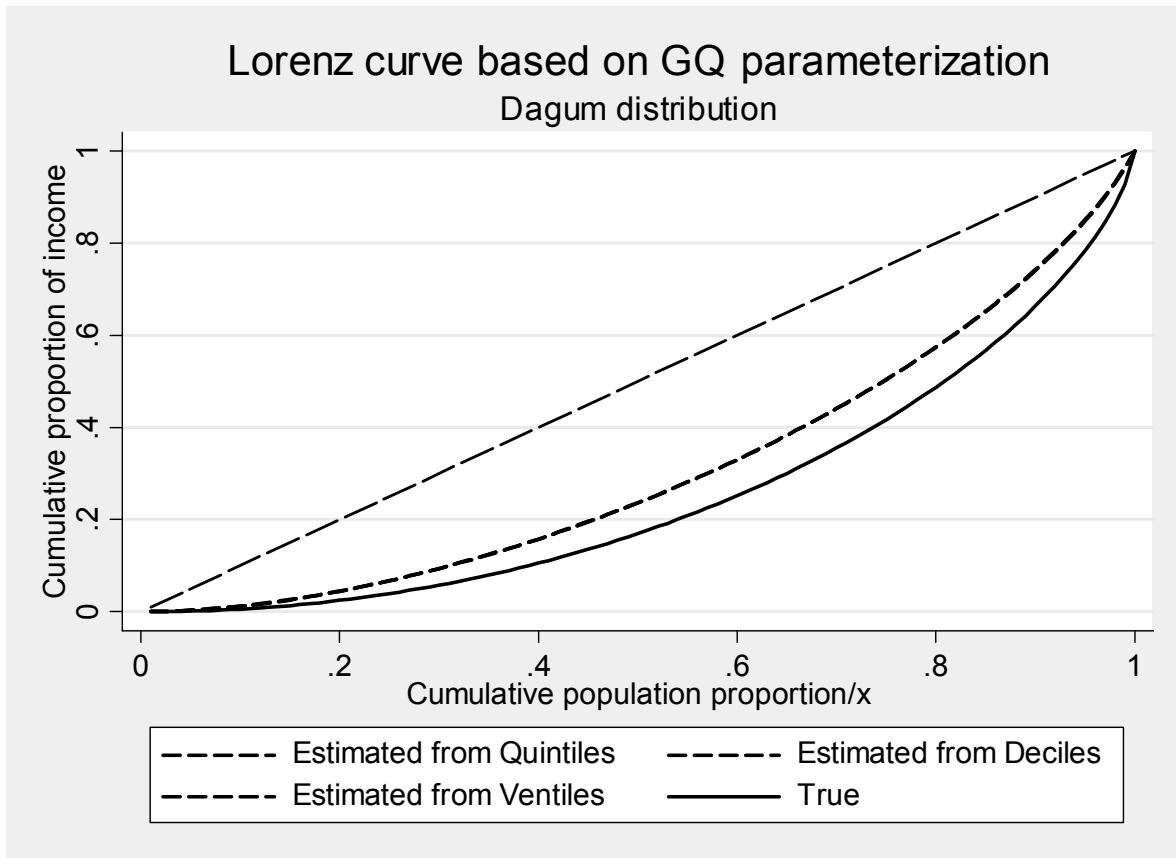
Note: Invalid Lorenz curve estimates are printed in boldface.

Table 3. Extent of Lorenz curve misestimation at along the entire support

Note: positive values represent overestimate and negative values represents underestimate

Cum. pop. proportion	GQ			Beta		
	Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
0.1	111%	104%	100%	90%	115%	129%
0.2	81%	80%	80%	82%	86%	89%
0.3	63%	63%	63%	63%	64%	64%
0.4	49%	50%	50%	49%	49%	49%
0.5	39%	39%	40%	39%	39%	38%
0.6	31%	31%	31%	31%	31%	30%
0.7	24%	24%	24%	24%	24%	24%
0.8	18%	18%	18%	18%	18%	18%
0.9	12%	12%	12%	11%	12%	12%

Figure 3. True and fitted Lorenz curve, Dagum distribution



Note: the fitted Lorenz curves from quintiles, deciles, and ventiles are observationally equivalent. The graph is identical for the Beta parameterization and is not shown.

**Monte Carlo simulation results for the Lorenz curve
(Multimodal distribution)**

Table 4. Sum of Squared Errors (SSE) and Sum of Absolute Errors (SAE) of the Lorenz curve estimate from grouped data

	GQ			Beta		
	Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
<i>Along the entire support</i>						
SSE	0.5951	0.4558	0.5048	0.2506	0.0663	0.0620
SAE	2.8596	3.4789	3.6097	1.8910	1.5521	1.7965
<i>Up to the 20th percentile of the population</i>						
SSE	0.0000	0.0149	0.0141	0.0000	0.0001	0.0002
SAE	0.0194	0.5435	0.5300	0.0157	0.0389	0.0536

Table 5. Lorenz curve estimates from grouped data in the left tail of the distribution

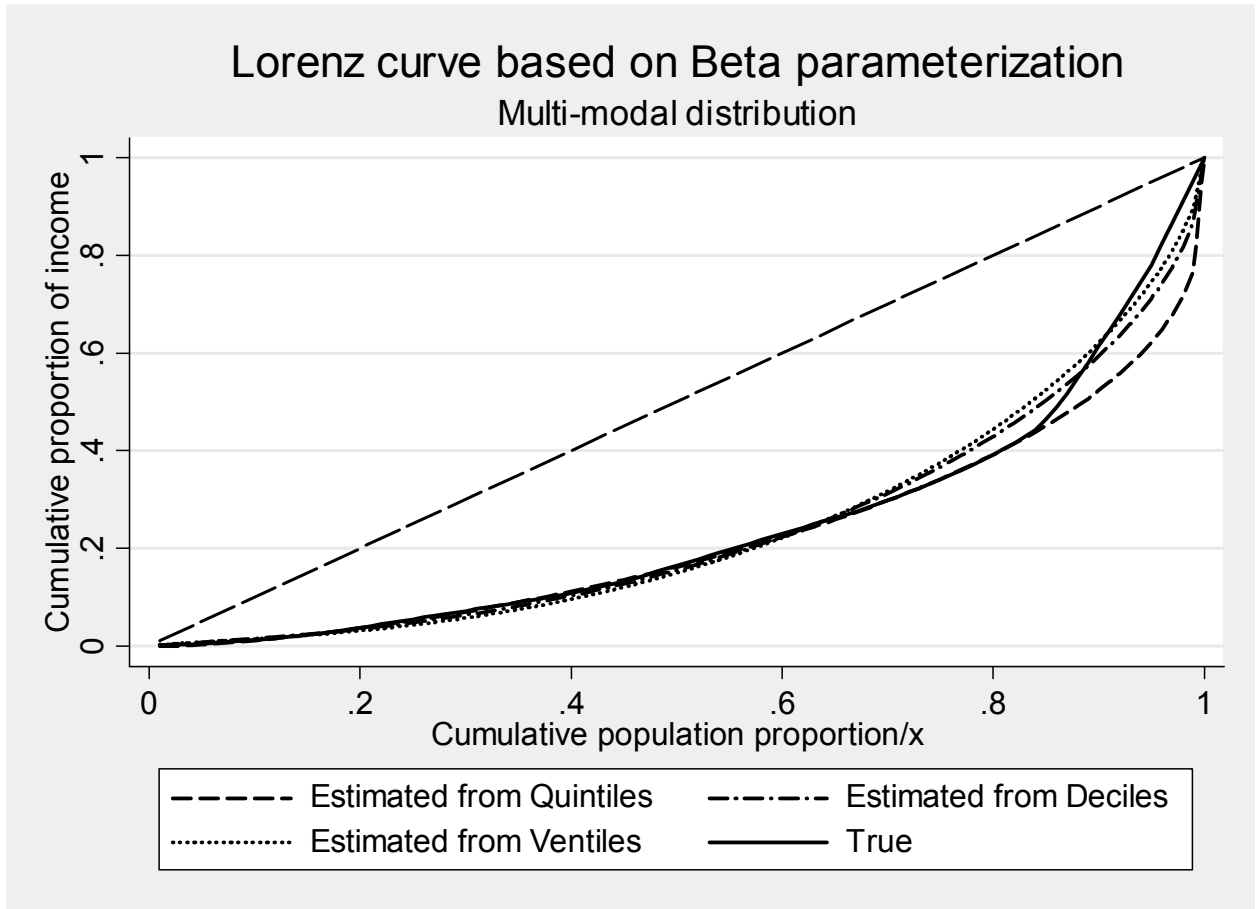
Cum. pop. prop.	<i>True</i>	GQ			Beta		
		Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
0.01	0.0007	0.0005	-0.0291	-0.0260	-0.0007	0.0017	0.0023
0.02	0.0015	0.0012	-0.0288	-0.0261	-0.0003	0.0031	0.0040
0.03	0.0023	0.0020	-0.0281	-0.0257	0.0005	0.0045	0.0054
0.04	0.0033	0.0029	-0.0270	-0.0250	0.0015	0.0059	0.0067
0.05	0.0044	0.0040	-0.0257	-0.0240	0.0028	0.0072	0.0080
0.06	0.0057	0.0052	-0.0241	-0.0228	0.0042	0.0086	0.0092
0.07	0.0070	0.0065	-0.0224	-0.0213	0.0058	0.0101	0.0105
0.08	0.0083	0.0080	-0.0204	-0.0196	0.0075	0.0116	0.0117
0.09	0.0099	0.0095	-0.0183	-0.0177	0.0094	0.0131	0.0130

Table 6. Extent of Lorenz curve misestimation along the entire support

Note: positive values represent overestimation and negative values represent underestimation

Cum. pop. proportion	GQ			Beta		
	Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
0.1	-5%	-235%	-231%	-4%	24%	21%
0.2	-7%	-61%	-63%	-2%	-7%	-16%
0.3	-6%	-22%	-23%	-3%	-12%	-20%
0.4	2%	-1%	-1%	3%	-5%	-12%
0.5	-1%	2%	2%	-1%	-5%	-9%
0.6	-1%	3%	3%	-3%	-3%	-4%
0.7	1%	6%	7%	0%	5%	6%
0.8	0%	5%	5%	0%	9%	13%
0.9	-18%	-14%	-14%	-15%	-3%	1%

Figure 4. True and fitted Lorenz curve, multimodal distribution



Note: The graph corresponding to the GQ gives rise to the same qualitative conclusions and is not shown.

Monte Carlo simulation results for poverty and inequality

Table 7. Poverty biases (in percentage points), Dagum distribution

Poverty indicator	Poverty line Median x:	<i>True</i>	GQ			Beta		
			Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
Poverty headcount ratio	1.33	66.9	0.18	0.34	0.53	-0.04	-0.21	-0.21
	0.5	21.7	-0.49	-0.67	-0.76	-0.83	-0.48	-0.21
	0.33	13.0	-0.10	-0.25	-0.27	-1.07	-1.07	-0.98
	0.25	9.1	0.01	0.03	0.05	-1.07	-1.13	-1.18
Poverty gap ratio	1.33	31.7	0.09	0.05	0.10	0.08	0.09	0.19
	0.5	9.6	0.03	0.06	0.45	-0.03	-0.28	-0.30
	0.33	5.8	0.27	0.36	0.62	0.52	0.00	-0.22
	0.25	4.0	0.35	0.51	0.67	1.03	0.41	0.04
Squared poverty gap	1.33	19.9	0.07	0.04	0.07	0.19	0.07	0.09
	0.5	5.9	0.23	0.35	0.48	1.07	0.37	0.01
	0.33	3.5	0.39	0.55	0.68	2.20	1.16	0.54
	0.25	2.4	0.47	0.68	0.84	3.36	1.80	1.02

Table 8. Poverty biases (in percentage points), Multimodal distribution

Poverty indicator	Poverty line Median x:	<i>True</i>	GQ			Beta		
			Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
Poverty headcount ratio	1.33	71.3	1.95	-1.72	-1.79	0.38	-7.04	-9.74
	0.75	43.5	-2.55	-6.90	-7.83	-0.43	-1.05	-1.38
	\$2/day	4.2	0.26	-	-	-0.60	-	-
	\$1/day	0.0	-	-	-	-	-	-
Poverty gap ratio	1.33	34.90	0.34	-1.20	-1.40	0.65	-0.51	-0.44
	0.75	18.17	-0.92	0.26	0.32	-1.00	0.81	2.30
	\$2/day	0.97	0.38	-	-	1.65	-	-
	\$1/day	0.0	-	-	-	-	-	-
Squared poverty gap	1.33	21.33	0.41	1.47	1.15	0.61	1.00	1.97
	0.75	9.80	0.16	4.78	4.07	0.00	0.85	2.33
	\$2/day	0.25	0.35	-	-	3.78	-	-
	\$1/day	0.0	-	-	-	-	-	-

Note: The international \$1.08/day poverty line (at 1993 international US\$) produces zero poverty estimates since the lowest per capita GDP in 2004 was \$515 (Sierra Leone) while the yearly equivalent of the \$1.08/day poverty line is \$448. Interestingly, this suggests that, for the \$1/day international poverty line, the vast majority of global poverty can be associated (against the counterfactual of even national income distributions) with intra-national inequalities.

Table 9. Inequality biases (in percentage points)

Distribution:	<i>True</i>	GQ			Beta		
		Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
Dagum	38.17	-0.22	-0.15	-0.12	0.03	-0.04	-0.04
Multimodal	52.43	5.10	3.59	4.53	3.73	0.64	-0.12

Deterministic comparison results for poverty and inequality

Table 10. Poverty headcount ratio biases (in percentage points)

Survey:	Poverty line:	<i>True</i>	GQ			Beta		
			Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
Vietnam	\$1/day	5.2	-1.00	-1.10	-0.19	-0.77	-0.51	n/a
	Capability	35.7	0.17	0.30	1.71	0.11	0.04	n/a
	\$2/day	41.9	-0.40	-0.35	1.52	-0.35	-0.64	n/a
Tanzania	Capability	40.4	-0.30	-0.15	-0.13	-0.25	-0.23	-0.20
	\$1/day	75.4	1.00	0.64	0.57	1.04	0.33	-0.23
	\$2/day	97.8	0.34	0.19	0.16	n/a	0.96	0.94
Nicaragua	Capability	30.6	-0.45	-0.30	-0.24	-0.62	-0.42	-0.27
	\$1/day	44.6	0.09	0.14	0.18	0.35	0.31	0.29
	\$2/day	79.0	0.58	0.26	0.17	0.40	0.10	-0.12

Table 11. Poverty biases (in percentage points). GQ estimation method.

Poverty headcount ratio:						
	Quintiles		Deciles		Ventiles	
Poverty line:	China	Brazil	China	Brazil	China	Brazil
Median x						
2.00	-0.13	-0.2	0.05	-0.3	0.10	-0.4
1.50	0.17	0.0	0.30	-0.3	0.33	0.3
1.00	-0.55	6.5	-0.64	6.1	-0.67	7.7
0.50	0.94	-5.2	0.83	-5.2	0.79	-3.1
0.25	-2.43	-12.3	-2.04	-11.7	-1.94	-9.9
0.20	-6.30	-14.1	-5.71	-13.4	-5.56	-11.7

Poverty gap ratio:						
	Quintiles		Deciles		Ventiles	
Poverty line:	China	Brazil	China	Brazil	China	Brazil
Median x						
2.00	0.00	-2.5	0.04	-2.2	0.05	-1.0
1.50	-0.08	-3.2	-0.08	-2.7	-0.09	-1.3
1.00	-0.02	-4.0	-0.03	-3.1	-0.04	-1.4
0.50	-0.02	-8.8	0.03	-6.5	0.04	-5.3
0.25	-0.63	-8.7	-0.62	-4.5	-0.62	-4.0
0.20	0.52	-7.6	0.37	-2.5	0.33	-2.3

Squared poverty gap:						
	Quintiles		Deciles		Ventiles	
Poverty line:	China	Brazil	China	Brazil	China	Brazil
Median x						
2.00	-0.09	2.1	-0.08	-6.1	-0.08	-4.7
1.50	-0.14	-7.0	-0.13	-5.5	-0.14	-4.1
1.00	-0.20	-7.1	-0.18	-4.8	-0.17	-3.6
0.50	-0.52	-6.6	-0.47	-1.3	-0.46	-1.2
0.25	-0.94	2.8	-0.94	14.9	-0.94	12.5
0.20	-1.47	9.3	-1.43	25.6	-1.42	21.6

Table 12. Inequality biases (in percentage points)

Survey:	True	GQ			Beta		
		Quintiles	Deciles	Ventiles	Quintiles	Deciles	Ventiles
Brazil	71.04	-0.63	0.03	0.29	-2.00	-1.16	-0.58
China	38.6	-0.18	-0.11	-0.10	0.07	0.02	-0.04
Tanzania	37.2	0.58	0.56	-0.26	0.83	0.69	n/a
Vietnam	35.0	0.24	0.07	0.04	0.36	0.14	0.03
Nicaragua	45.2	0.22	0.08	0.09	0.19	0.09	0.08