# We forgot the middle class! Inequality underestimation in a changing Sub-Saharan Africa<sup>\*</sup>

F. Clementi<sup>a,†</sup>, A. L. Dabalen<sup>b</sup>, V. Molini<sup>b</sup>, and F. Schettino<sup>c</sup>

<sup>a</sup> University of Macerata, Macerata, Italy
 <sup>b</sup> World Bank, Washington DC, USA
 <sup>c</sup> University of Campania "Luigi Vanvitelli", Naples, Italy

September 8, 2017

#### Abstract

The creation or consolidation of national middle classes and the changes in consumption patterns in many Sub-Saharan African (SSA) countries suggest reconsidering the way welfare and consequently inequality is typically measured in these countries. Using consumption to measure welfare, as it is typically done in many SSA countries, can lead to an important loss of information regarding the real welfare of the top 10-20 percent of the welfare distribution that is generally referred as "middle class" in these countries; this loss of information can lead to a substantial underestimation of inequality. This paper proposes a method capable of correcting the middle-class part of the consumption distribution using information coming from the income distribution of the same surveys. This way, we argue, the inequality measures calculated on the new distribution can more accurately reflect the real welfare distribution in these countries. Preliminary results from 6 SSA countries indicate an increase, compared to original data, of about 20 percent in the Gini index and all the other inequality measures .

**Keywords:** middle class; inequality; Sub-Saharan Africa; consumption expenditure; income; combined distribution **JEL classification:** C46; D31; D63

## **1** Introduction

The two decades between 1995 and 2015 have represented for developing countries a period of fast growth and poverty reduction. Besides reducing poverty, global growth also had a profound

<sup>\*</sup> The authors acknowledge financial support from the World Bank. We thank Federica Alfani (FAO – Food and Agriculture Organization of the United Nations) for excellent assistance with data preparation. We also thank (in random order) Pierella Paci, Francisco H. G. Ferreira, Paolo Verme and Christoph Lakner for comments on an earlier version of the manuscript. Of course, we are the sole responsible for all possible errors the paper may contain.

<sup>&</sup>lt;sup>†</sup> Corresponding author: <u>fabio.clementi@unimc.it</u>.

impact on the social structure of many developing countries, buttressing the creation of a world "middle class" (Milanovic and Yitzhaki, 2002).

The Sub-Saharan Africa (SSA) middle class improvements are more modest than other developing regions, nonetheless there are important progresses. According to the African Development Bank (2011), the continental middle class<sup>1</sup> accounts for 14 percent of total population or about 127 million people (African Development Bank, 2011).

The most prominent characteristic of the middle class group, either in developed (Goldthorpe, 1987) or developing countries (Geithman, 1974) is the type of occupation. Middle class members have in general a formal employment,<sup>2</sup> in either public or private sector, live in urban areas and work in non-farm activities; income from these sources tends to be less volatile and affected by seasonality than agricultural incomes. In SSA, depending on the countries and years, the share of formal employment on total ranges between 10 to 20 percent (International Monetary Fund, 2012; Golub and Hayat, 2014). In Ghana (Honorati and Johansson de Silva, 2016), one of the few countries with repeated and comparable cross sections over two decades, the share of private wage employment on total employment nearly tripled from 6 percent to 16 percent. Adding public wage earners, 22 percent of Ghana's workers in 2012 have a formal job and report stable monthly earnings.

Another distinguishing feature of the middle class group is the propensity to save to save in form of financial assets or money deposited in bank accounts; middle-class households are more likely than poor people to save since marginal propensity to consume declines with higher welfare levels. The World Bank<sup>3</sup> calculates that in SSA the number of depositors with commercial banks increased threefold between 2004 and 2015, from 50 every 1,000 adults to 155 every 100 adults. Likewise, the new African middle class has Western-style consumption patterns (McKinsey, 2016); middle class housheolds diversify their expenditures away from basic needs towards more durable goods (home appliances, computers, smart-phones, cars), luxury goods, entertainment (restaurants, movies, travels), and in some cases properties.

All these elements taken together, suggest that in many SSA countries where welfare is measured via consumption, we miss a substantial piece of information regarding the top 10-20

<sup>&</sup>lt;sup>1</sup> This is defined as the group earning between 4 to 20 USD PPP per capita per day.

<sup>&</sup>lt;sup>2</sup> For low-income SSA countries, formal employment covers public sector employment and wage employment in other sectors (International Monetary Fund, 2012).

<sup>&</sup>lt;sup>3</sup> World Development Indicators, <u>http://data.worldbank.org/data-catalog/world-development-indicators</u>.

percent of the population. Consumption has undoubtedly important advantages when looking at lower income percentiles, typically characterized by volatile incomes and high seasonality, yet it underestimates the actual welfare of the middle class because it does not account for savings (as well as transfers) and does not factor in many non-basic expenditures. Intuitively, this underestimation affects the calculation of various inequality measures, yielding much lower results.

On the other hand, the mere substitution of income as the main welfare indicator for SSA households still looks quite untimely. At the current stage of SSA countries' development, income data are still not representative for the bulk of the households' welfare, because of the prevalence of the informal sector and farm activities.

The aim of the present paper is to overcome these problems by testing and applying a new methodology for reconstructing a more realistic welfare distribution in SSA, and then measuring inequality more accurately. In a nutshell, the idea is to correct the consumption distribution around the top quintile using information coming from the income distirbution from the same survey.

The remainder of the paper is structured as follows. Section 2 motivates our work by discussing reasons why inequality might be underestimated in SSA. Section 3 introduces the data used in the empirical application and explains our approach to inequality measurement. Section 4 analyzes the performance of the proposed method and presents corrected estimates of overall inequality levels. Finally, Section 5 provides a summary and conclusions.

## 2 The welfare puzzle: top incomes, income and consumption

Budget surveys struggle to include the welfare of hard-to-survey populations, in particular the extremely rich. Methods have been proposed to address this issue, the most famous one is to compare top incomes in household surveys with tax records (e.g. Atkinson et al., 2011). In developing countries, where generally is more difficult to obtain this type of information from tax authorities, analysis on top incomes started later but has recently gained momentum (e.g. Leigh et al., 2009, Alvaredo, 2010, and Sanhueza and Mayer, 2011).

While similar in spirit, this paper departs from the top-income literature since our main concern is not the top 1–5 percent but the new middle class which tends to occupy the top 20 percent of the welfare distribution. Therefore, the re-calibration exercise involves a much bigger portion of the distribution. Another important difference pertains the source of information used to

reconstruct part of the distribution. Whereas for top incomes the re-estimation parameters come from another distribution—the tax records—based on a completely different sample, in this exercise we correct the top part of consumption distirbution using information coming from the same sample.

As mentioned in the introduction, the novelty of this paper lays in the way we correct consumption using information form income; it is thus important to understand the relative comparative advantages of the two measures at different points of the welfare distribution.

There are good reasons why many SSA countries, which generally have limited resources for data collection, have focused their attention on getting consumption data right and often disregarded income data collection. Consumption is generally regarded as easier to measure than income in low-income economies (Deaton and Zaidi, 2002). Formal household monetary incomes are mostly constituted by wages and non-labor monetary incomes (such as profits and rents). Yet, most of households in SSA countries earn other forms of monetary incomes, such as those coming from agricultural production (both for selling and for auto-consumption) and from informal activities; these magnitudes are typically easier to capture via the value of consumption rather than via income.

Moreover, monetary incomes in these countries routinely exhibit great seasonal variations (Tarozzi, 2007), while consumption expenditures tend to be naturally smoother (Friedman, 1957). For example, in agricultural economies like most of African countries, income is more volatile and affected by harvest seasons, so that relying on income as an indicator of welfare might under/over-estimate living standards significantly. Finally, consumption tends to be a better measure of permanent welfare, because households can borrow, draw down savings, or get public and private transfers to smooth short-run fluctuations.

When it comes to inequality measurement, however, consumption data show several limitations compared to income. First, while consumption is more informative than income for the bottom of the distribution, since it reflects—in addition to income—welfare transfers, interpersonal transfers and informal income (Meyer and Sullivan, 2004), data on consumption at the very top the distribution could seriously under-estimate welfare because of compelling evidence that the marginal propensity to consume declines as household welfare increases (McCarthy, 1995; Dynan et al., 2004; Jappelli and Pistaferri, 2014).

Second, consumption inequality measures are generally biased downward if the set of goods

in the consumption measure does not include items consumed by the rich (luxury goods, such as vacations, as well as irregularly purchased consumer durables, such as cars). These goods are sometimes not included in surveys or are excluded from the measure of consumption if they are; on the other hand, income, by just measuring the potential claim over items, is not affected by this under-reporting.

Third the inequality measured on consumption shows lower inequality than the data based on income. There is a simple reason for this. There may be people with zero annual income who, for example, finance their current spending out of previously accumulated savings. There are, obviously, no people with zero annual consumption. This makes the distribution according to income more "elongated" around the bottom and thus more unequal (the consumption distribution will be "truncated" at some minimum amount necessary to survive). Also, a similar thing happens at the other end of the distribution. There are many income-rich people who save a part of their income. Thus, their income is greater than their consumption. The high end of the distribution would be also more elongated in the case of income (Milanovic, 2010).

At first glance, therefore, correcting consumption around the top part of the distribution using information coming from income can overcome the limits of measuring inequality just using consumption. There are, however, some important theoretical issues it is worth to touch upon before discussing the proposed methodology: possible alternatives, the income measurement error, and the potential contribution of our method to the top-income analysis.

Regarding alternatives, one might argue that just adding to consumption the reported savings—typically collected in *ad hoc* section of many household budget surveys—could provide an equally reliable estimate of income that relays on real data rather than on re-estimation. There are, however, a number of problems with this method. First, there is a lot of heterogeneity on how the questions on savings are posed. For example, within the sample of housheold surveys we use, for Ghana and Niger there is an explicit question on the amount of savings owned by households either in bank accounts or in informal saving schemes. Therefore, in these two countries it would be feasible to add savings to consumption; in Malawi, Nigeria, Kenya and Uganda's surveys, on the other hand, the exact amount of savings is never asked, hindering this possibility.

Second, even when savings are reported correctly, their addition to consumption might not

5

suffice to reduce the gap with income.<sup>4</sup> Figure 1 illustrates the problem for Ghana.<sup>5</sup> Even adding savings, income still remains shifted to the right. Besides the above-mentioned list of items that are not factored in a typical consumption aggregate—or partially factored in, such as eating outside—there are problems in the food component too. This because, as mentioned before, the consumption aggregates are based on questionnaires intended to capture the consumption pattern of the vulnerable/poor. Therefore, food items not consumed by the vulnerable/poor household are generally not part of the items list:<sup>6</sup> these include imported items, items not typically part of the local diet, luxury food, and so forth.

#### [Figure 1 about here.]

Measurement error for incomes can also be an issue. Evidence suggests that income underreporting grows the higher the income (e.g Hlasny and Verme, 2016, and references therein); precisely, the part of the distribution we want to use for our re-estimation. The answer to this is twofold. As shown in the previous graph and based on the way consumption aggregates are constructed, the under-reporting of consumption in top deciles will be anyway higher than the under-reporting of income. Using the top deciles of income rather than consumption, will certainly not eliminate under-reporting but it will mitigate it compared to consumption. When calculating the inequality measures, this combined distribution will yield probably a lower bound of the real inequality but certainly an estimate much higher than that produced by consumption alone.

Moreover, and this links to the whole issue of top incomes, the proposed method is fully compatible with their estimation methods. The correction we propose can be considered a

<sup>&</sup>lt;sup>4</sup> Differently from other monetary aggregates, checking the accuracy of reported savings might result complicated. For example, to see weather food consumption values can be off, one typically converts quantities into calorie intake. Too high (over 10,000 calories per capita per day) or too low (1,000 calories per capita per day) values signal measurement error problems. To check for salaries accuracy, enumerators often ask for a recent salary slips. Controlling the savings is a bit more complicated, especially if the household is part of one of those informal savings schemes very frequent in Africa; in that case, it is very difficult to provide an official statement of the amount deposited.

<sup>&</sup>lt;sup>5</sup> The figure is a plot of the *complementary distribution function* for each variable, showing the proportion of values greater than each value—i.e. the complement of the cumulative distribution function. The observed probabilities are plotted on a doubly logarithmic scale, which is natural to use when focusing on the top part of the distribution because it accentuates the upper tail (see e.g. Clementi, 2016). Furthermore, in order to circumvent the scale difference between household incomes and consumption expenditures, the former have been median-adjusted by a multiplicative shift to yield identical centers of the consumption and income distributions. For more information on the data variables we have selected for plotting and how they have been pre-processed, see Section 3.1.

<sup>&</sup>lt;sup>6</sup> Sometimes the respondent is asked, in addition to the listed items, to report other food items generally consumed. Since households might report very different additional items, it is very difficult standardizing this information and add it to the consumption aggregates. Therefore, even if some items are not completely excluded, the way they are reported does not allow to improve the consumption aggregates of these richer households.

preliminary exercise that adjusts the distribution before further correcting it with the top-income methods. As mentioned before, we aim at improving the accuracy of the top 10–20 percent of the distribution, while top-income methods typically affect the top 1 percent. The ideal distribution we have in mind to estimate more accurately inequality is one that corrects the top quintile of the consumption distribution using income and further corrects the top 1 percent using tax records.

## **3** Data and methodological essentials

### 3.1 Household income and expenditure surveys in SSA

The present paper uses household-level data taken from several sources. For income, we use yearly data obtained from the Rural Income Generating Activities (RIGA) database, a collaborative effort of the Food and Agriculture Organization (FAO) of the United Nations, the World Bank, and the American University.<sup>7</sup> It is composed of a series of constructed variables about rural and urban income-generating activities created from the original consumption data sources. In particular, we will focus on the household-level income aggregate data set (RIGA-H), which includes a comprehensive measure of household income presenting aggregated data on different income sources, such as crop and livestock production, household enterprises, wage employment, transfers, and non-labor earnings.<sup>8</sup>

Data on households' consumption expenditure come instead from the original budget surveys compiled by national statistical bureaus and the World Bank, which can be easily linked to each country data set in the RIGA database. Specifically, we consider here the following household budget surveys: the Ghana "Living Standards Survey", 2005; the Kenya "Integrated Household Budget Survey", 2005; the Malawi "Integrated Household Survey", 2011; the Niger "National Survey on Household Living Conditions and Agriculture", 2011; the Nigeria "Living

<sup>&</sup>lt;sup>7</sup> <u>http://www.fao.org/economic/riga/riga-database/en/.</u>

<sup>&</sup>lt;sup>8</sup> The household income aggregates and their components included in the RIGA database closely follows the definition given by the International Labour Organization (ILO, 2003), which considers as income receipts those that (*i*) recur regularly, (*ii*) contribute to current economic well-being, and (*iii*) do not arise from a reduction in net worth (Carletto et al., 2007). These three criteria are embodied in each of the components of income; as such, irregular payments such as lottery earnings or inheritances, investments and savings, and the value of durables are not included in the RIGA definition and measure of income. Furthermore, costs are also taken into account to ensure that the final income aggregate is net of costs, as opposed to gross (which could overestimate the income a household actually has at his or her disposal). So far, the only reported cost, which is subtracted during the aggregation process, has been income tax— i.e. the contribution to social security and health system (Quiñones et al., 2009).

Standards Survey", 2004; the Uganda "National Household Survey", 2005.<sup>9</sup> The main consumption variable that is used in the paper is the household total annual expenditure on food and non-food items.

Before undertaking the empirical analysis, both the income and consumption variables were spatially and temporally deflated to 2005 national currency units and expressed per capita. Furthermore, observations with negative and zero incomes were excluded from the analysis, because some indices of inequality are defined only for positive values.<sup>10</sup> Table 1 presents distributional statistics for the consumption and income variables used in this study. Compared to income, consumption expenditure typically produces lower estimates of inequality, independently on the measure that one considers—the Gini coefficient, the mean logarithmic deviation (MLD), or the Theil index. As mentioned in the previous section, this is to be expected and can be explained by a declining marginal propensity to consume and by the fact that consumption surveys tend to understate the spending on durables at the top. Instead, an argument for using consumption rather than income is that data on the former are often of a higher quality in developing and emerging economies and are less vulnerable to idiosyncratic shocks, as households tend to smooth their consumption over time. Because estimates of inequality will be biased if computed using any single one of these variables, what is needed to obtain consistent estimates of inequality is a combination of the information coming from the consumption and income data. The following subsection appeals to multiple-imputation methods in order to achieve this.

[Table 1 about here.]

## **3.2** A multiple-imputation approach to inequality measurement

Given data, our approach to inequality measurement is adapted from earlier work by Jenkins et al. (2011), who proposed a parametric multiple-imputation method to measure income inequality with right-censored (top-coded) data, and goes through the following steps.

First, by means of model selection and goodness-of-fit techniques, we select the best fitting

<sup>&</sup>lt;sup>9</sup> Notice that the RIGA project covers more countries (e.g. Ethiopia, Madagascar, and Tanzania) and provides data sets that are sometimes more recent than those used in the present analysis. However, limited coverage of the population and issues of accuracy caused us to focus only on the six countries (and years) mentioned in the text.

<sup>&</sup>lt;sup>10</sup> Accordingly, the sampling weights of households— used in all calculations— have been re-calibrated in such a way that estimates from the samples after deletion of non-positive records are forced to fit the initial population-level information on the households' geographical location and area of residence (rural/urban).

parametric model for the consumption and income distributions of each country. The models that are fitted to micro-data belong to the family of generalized beta distributions introduced by McDonald and Xu (1995a,b), which includes the four-parameter generalized beta II distribution (GB2) with probability density function

$$f(x; a, b, p, q) = \frac{ax^{ap-1}}{b^{ap}B(p,q)[1+(x/b)^a]^{p+q}}, \quad x > 0,$$
(1)

and cumulative distribution function

$$F(x; a, b, p, q) = I_z(p, q), \quad z = (x/b)^a, \quad x > 0,$$
(2)

where  $B(p,q) = \frac{\Gamma(p)\Gamma(q)}{\Gamma(p+q)}$  is the (complete) beta function,  $\Gamma(\cdot)$  is the gamma function, and  $I_z(p,q) = \frac{B(z;p,q)}{B(p,q)}$  is the regularized incomplete beta function—which is the ratio of the incomplete and complete beta functions. All four parameters are positive, with *b* being the scale parameter and *a*, *p*, and *q* being the shape parameters. The GB2 distribution is a flexible functional form incorporating many distributions as special cases. Of these, our interest is drawn to the three-parameter models of Singh and Maddala (1976) and Dagum (1977), which are often used in the income distribution literature and can be obtained as special cases of the GB2 for, respectively, p = 1 and q = 1.<sup>11</sup>

The second stage of our approach uses the model's parameter estimates to derive imputed values for observations above some lower-bound consumption threshold defining (in absolute terms) a minimum middle-class standard of living. Opting for such a definition of the middle class in the context of developing countries seems reasonable for at least two orders of reasons.<sup>12</sup> First, unlike developed countries, we can not use relative welfare measures for defining the middle class, since in developing countries the latter does not often coincide with some function of the distribution's median (i.e. the middle class does not generally occupy the center of the distribution); scholars thus often opt for absolute measures. Second, a further complication one might encounter in developing countries is defining an upper bound. As already anticipated, these countries often

<sup>&</sup>lt;sup>11</sup> For details, see McDonald (1984), McDonald and Xu (1995a,b), Kleiber and Kotz (2003), and McDonald and Ransom (2008). Of particular importance in the current context is the desirable behavior of the GB2 and related distributions in their upper tail, which is heavy in that it decays like a power function as the size variable increases—rather than decaying exponentially fast like, for instance, the log-normal distribution with mildly heavy upper tail. For more on this, see e.g. Kleiber (1996), Schluter and Trede (2002), Kleiber (2008), and Kleiber and Kotz (2003). <sup>12</sup> See e.g. the discussion in Corral Rodas et al. (2017) for Nigeria.

focus their attention on getting consumption data right and disregard income data collection. Since consumption is very accurate in capturing the well-being of poorer people, while it is rather imprecise in capturing that of people living in the upper percentiles, it follows that when defining the middle class in these countries it seems reasonable to opt for a lower-bound threshold (rather than an interval) of the type "middle class and above" and leave the border between middle class and upper class somehow undefined.

For the purposes of this paper, two absolute thresholds are used to define the middle class in SSA, both derived from the African Development Bank (2011):<sup>13</sup>

- 1. Per capita daily consumption greater than \$4 in PPP US dollars, which includes both the lower- and the upper-middle class.
- 2. Per capita daily consumption greater than \$10 in PPP US dollars, which identifies the upper-middle class.

The reason why we work with two thresholds is that using only the second one could lead us to rather conservative estimates of inequality, as the correction of the consumption data in this case typically affects a tiny group of households at the far end of the distribution. Instead, by using also the first threshold we can impute income variability to the original data for a broader group of households, which prevents us from a potentially downward-biased estimation of inequality.

Once the best fitting parametric model for both consumption and income data has been selected, imputed values for observations above a lower-bound consumption threshold can be derived by means of the so-called "inverse transform method". That is, given the fitted GB2, the cumulative distribution function for each observation i above the consumption threshold  $t_c$  is, using standard notation for left-truncated distributions,

$$G(x_i; \hat{\theta}) = u_i = \frac{F(x_i; \hat{\theta}) - F(t_c; \hat{\theta})}{1 - F(t_c; \hat{\theta})}, \quad x_i > t_c, \quad u_i \in [0, 1),$$
(3)

where  $\hat{\theta} = \{\hat{a}_y, \hat{b}_c, \hat{p}_y, \hat{q}_y\}$  is the set of parameter estimates and the subscripts *c* and *y* refer to consumption and income, respectively.<sup>14</sup> Inverting, we have

<sup>&</sup>lt;sup>13</sup> In the empirical application that follows, we convert these thresholds into annualized money amounts and express them in local currency using 2005 PPP conversion factors from the World Bank (http://data.worldbank.org/indicator/PA.NUS.PRVT.PP).

<sup>&</sup>lt;sup>14</sup> Notice that our approach is designed to alter the *shape* of the consumption distribution at the top end, but not its scale. That is why the set of parameter estimates used for imputing values above the consumption threshold includes the shape parameter estimates for the income distribution of each country  $(\hat{a}_y, \hat{p}_y, \text{ and } \hat{q}_y)$  and the estimated scale

$$G^{-1}(u_i;\hat{\theta}) = x_i = F^{-1}\{u_i[1 - F(t_c;\hat{\theta})] + F(t_c;\hat{\theta});\hat{\theta}\}, \quad x_i > t_c, \quad u_i \in [0,1).$$
(4)

Thus, a value of  $x_i$  for each observation above the consumption threshold is generated by substituting into this expression a value of  $u_i$  that is equal to a random draw from a standard uniform distribution.<sup>15</sup>

The combination of imputed values for observations above the consumption threshold and observed expenditures for those lying below produces a partially synthetic data set for each country to which we apply complete-data methods to estimate inequality statistics such as the Gini coefficient, the MLD, and the Theil index. In the last step, repetition of the process a large number of times (to control for the randomness of each partially-imputed data set) produces M synthetic data sets for each country and, correspondingly, M sets of inequality estimates which we combine using the averaging rule proposed by Reiter (2003). That is, supposing that inference is required about some scalar measure of inequality Q, and indexing the partially synthetic data sets by j = 1, ..., M, one can estimate Q by using

$$\bar{q}_M = \frac{1}{M} \sum_{j=1}^M \hat{q}_j,\tag{5}$$

which is the simple average of the point estimates  $\hat{q}_j$  that are derived using complete-data methods from each of the *M* partially synthetic data sets. In the next section, we report estimates based on M = 999.

#### 4 **Results**

This section is divided in two parts. In the first, we discuss the models' diagnostics looking at whether the proposed parametric distributions fit the original data for both income and consumption. In the second, we calculate a set of distributional indicators on the original consumption data and compare them to the corrected ones.

parameter for the consumption distribution  $(\hat{b}_c)$ .

<sup>&</sup>lt;sup>15</sup> In Equations (3) and (4), the values of the GB2 cumulative distribution function at the truncation point  $t_c$ ,  $F(t_c; \hat{\theta})$ , and those for each x above the consumption threshold,  $F(x_i; \hat{\theta})$ , are estimated by inserting parameter estimates into Equation (2). The cumulative distribution functions in the cases of the Singh-Maddala and Dagum distributions are given by simpler expressions and can be found, for instance, in Kleiber and Kotz (2003, ch. 6).

#### 4.1 Parameter estimation, model selection and goodness of fit

All generalized beta models considered in this paper were fitted to consumption and income distributions using maximum likelihood estimation. For fitting models to data, we used Stata's programs developed by Jenkins (1999, 2007, 2014). These programs maximize the log-likelihood numerically and estimate parameter variance using the negative inverse Hessian. A number of distributional measures implied by fitted models, and their associated standard errors computed using the delta method, were also obtained using the Stata's commands developed by the author.

Tables 2 and 3 present our estimates of models' parameters together with their standard errors, the values of log-likelihood (ln*L*) at last iteration, and model selection criteria such as the Akaike (Akaike, 1973) and Bayesian (Schwarz, 1978) information criteria (AIC and BIC).<sup>16</sup> In order to compare the fit of the GB2 model and its nested alternatives (the Singh-Maddala and Dagum), we also give the results of likelihood ratio tests for the fitted models. The likelihood ratio statistics takes the form

$$2\left[\ln L(\hat{\theta}_{\rm U}) - \ln L(\hat{\theta}_{\rm R})\right] \sim \chi^2(h),\tag{6}$$

where  $\ln L(\hat{\theta}_U)$  and  $\ln L(\hat{\theta}_R)$  are, respectively, the log-likelihood values corresponding to the unconstrained (GB2) and nested or restricted models (Singh-Maddala and Dagum),  $\hat{\theta}$  is the set of estimated parameters, and *h* is the difference in the number of parameters in the two compared models (equal to 1 in our setting). The differences between GB2 and its nested alternatives can be thus compared using a chi-square ( $\chi^2$ ) distribution with one degree of freedom. In the tables, asterisks are placed next to the likelihood ratio values if the improvement gained in adding a further parameter is of practical significance.<sup>17</sup>

The results of model selection for consumption distributions, presented in Table 2, suggest that the GB2 model is a better fit to data in all countries except Niger, where the Dagum model is as good as the GB2. For income data, the results in Table 3 are somewhat mixed. The GB2 is

<sup>&</sup>lt;sup>16</sup> The expressions for the log-likelihood of the GB2 and its nested models (the Singh-Maddala and Dagum) are given in Kleiber and Kotz (2003). Model selection criteria will select, when comparing models with the same number of parameters, the model with the smallest  $l = -\ln L$  according to the formula  $(2 \times l) + (d \times k)$ , where k represents the number of parameters in the fitted model and d = 2 for the usual AIC or  $d = \ln N$  (N being the number of observations) for the so-called BIC. Hence, when comparing models fitted by maximum likelihood to the same data, the smaller the AIC or BIC the better the fit. When comparing models using the log-likelihood criterion, the larger the lnL the better the fit.

<sup>&</sup>lt;sup>17</sup> The critical value of the  $\chi^2(1)$  distribution is 3.84 at the 5% level.

clearly the best model for Ghana, Malawi, Nigeria, and Uganda income distributions, whereas for Kenya the Singh-Maddala seems to be as good as the GB2. A similar conclusion applies to Niger, but in this case the Dagum model fits the data better than the alternatives. In general, the GB2 model gives the best fit to both consumption and income data in 4 out of 6 analyzed countries (i.e. Ghana, Malawi, Nigeria, and Uganda). For Niger, the three-parameter Dagum fits the observed consumption and income distributions better than any of the alternative models, whereas for the 2005 Kenyan survey on household income the Singh-Maddala has to be preferred to the GB2 in practical applications due to its smaller number of parameters. However, although the fit of the GB2 model was not quite as good as for the Singh-Maddala in the case of Kenyan incomes, but very good nonetheless, given the need of working with a single imputation model we shall assume in the following that the distribution of household income in Kenya is described by the fourparameter GB2 model, which is also the preferred one for parametric modeling of the country's consumption distribution.

#### [Table 2 about here.]

#### [Table 3 about here.]

Goodness of fit of the functional forms chosen according to model selection methods is assessed graphically using the Cox-Snell residuals (Cox and Snell, 1968)

$$\hat{e}_i = -\ln\hat{S}(x_i;\hat{\theta}), \quad i = 1, \dots, N,$$
(7)

where  $\hat{S}(x_i; \hat{\theta}) = 1 - \hat{F}(x_i; \hat{\theta})$  denotes the estimated survival or complementary cumulative distribution function for each *x*. If the model is good, these residuals should behave like a sample from an exponential distribution with parameter 1. A plot of the ordered estimated  $\hat{e}_i$  (i = 1, ..., N) against the quantiles of a one-parameter exponential distribution should therefore be roughly a straight line with slope 1 (e.g. Elandt-Johnson and Johnson, 2003).<sup>18</sup>

The Cox-Snell residuals are shown, respectively, in Figure 2 for the best fitting consumption models and in Figure 3 for the functional forms best describing income data. The results indicate that, apart from some noisiness by the most extreme observations, the fit of the selected models was quite good for the majority of countries along the entire distribution range.

<sup>&</sup>lt;sup>18</sup> See Quintano and D'Agostino (2006) for a more detailed methodological explanation of the use of Cox-Snell residuals in parametric income distribution modeling. For another application, see also Betti et al. (2008).

The two most notable exceptions are the Dagum distribution for Niger consumption data and the GB2 for Ugandan household incomes, where the upper tails of the distributions present deviations for these assumptions that appear in the plots as data that do not fit the straight line. However, only about 2 percent of observations at higher quantiles deviated from the 45-degree line in Figures 2(d) and 3(f), which is negligible for overall fit. We can therefore conclude that the Dagum model for Niger and the GB2 for the rest of countries considered in this study can be used as theoretical models for describing the empirical distribution of both consumption and income.

#### [Figure 2 about here.]

#### [Figure 3 about here.]

Goodness of fit of the selected models is also evaluated by comparing the sample values of distributional indicators reported in Table 1 with their counterparts implied by the fitted models see the last four columns of Tables 2 and 3.<sup>19</sup> Specifically, in Figures 4 and 5 the comparison relies on checking for overlap between 95 percent confidence intervals of theoretical and sample indicators to draw conclusions about the accuracy of selected distributional statistics deduced by parameter estimates. The results suggest that for most of the indices (i.e. the mean, the Gini coefficient, the MLD, and the Theil index) the best fitting models produce theoretical values that are quite often in a close agreement with the corresponding sample values—the respective confidence intervals overlap in a way that let us exclude that the predicted values and the actual sample estimates of chosen indicators can be considered different. The most notable exceptions are the theoretical estimates implied by the best fitting GB2 model for Ugandan incomes, which differ significantly from the corresponding sample estimates. This fact could reflect the documented poor performance of the GB2 at the top of the Ugandan income distribution, where there is a systematic departure of empirical observations from the theoretical predictions of the assumed specification. However, the results for the nested three-parameter Singh-Maddala and Dagum distributions (not shown here but available on request from the authors) are even worse, especially for higher quantiles. This explains why we shall keep using the GB2 distribution for imputing observations in the top part of the Ugandan welfare distribution.

[Figure 4 about here.]

<sup>&</sup>lt;sup>19</sup> The analytical expressions for all indices considered here, which are functions of the estimated parameters of the GB2 and its nested distributions, can be found *inter alia* in Kleiber and Kotz (2003, ch. 6) and Jenkins (2009).

## 4.2 Distributional indicator results

Table 4 and Figure 6 display the simulation results by country using the \$4 and \$10 middle class thresholds. As discussed in the methodological section, with the \$10 threshold the correction of consumption data applies to a smaller group of households than with the \$4 threshold (compare shares in columns 4 and 9 of Table 4). The estimated inequality using the \$4 threshold is clearly higher since more information is taken from the income distribution and, as discussed before, income tends to have higher variability than consumption.

[Table 4 about here.]

#### [Figure 6 about here.]

For example, in Ghana, where according to the \$4 line the middle-class group would account for about 27 percent of the population, the correction of consumption for this group would lead to a Gini of 0.46 from 0.42 in the original data (compare columns 5 in Tables 4 and 1). On the other hand, when using the \$10 line, only 4 percent of the households will see their consumption corrected and the obtained Gini is 0.43 (column 10 in Table 4.). Likewise, in all analyzed countries, the two thresholds define an upper and lower bound for the simulated Gini, where the upper bound is obtained from the \$4 line and the lower is obtained from the \$10 line (Figure 6).

Figure 7 displays the impact of the correction on the Nigerian consumption data. In the first quadrant, the correction is applied on the middle-class group defined by the \$4 threshold, whereas in the second it is applied using the \$10 line. Correcting consumption implies using from the middle-class thresholds onwards (vertical dashed line) the parametrized tail derived from the corresponding income distribution (blue line) that—in Nigeria like in all the other considered countries—lays above that of the consumption distribution. The difference is very clear when cutting the consumption distribution at \$4 (left quadrant), but less pronounced when using the \$10 line (right quadrant). As a consequence, when using the \$4 line, we re-estimate a bigger chunk of the distribution and introduce in this way more variability than in the case of \$10, leading to a bigger increase of the Gini index (see Table 4).

[Figure 7 about here.]

## 5 Concluding remarks

The creation or consolidation of national middle classes and the changes in consumption patterns in many SSA countries suggest reconsidering the way inequality is typically measured in these countries. Specifically, the use of consumption as the main welfare measure can lead to an important loss of information regarding the real welfare of the top 10-20 percent of the welfare distribution and, as consequence to a substantial underestimation of inequality.

The present paper develops a methodology that re-estimates the top part of the consumption distribution using information from the the income distribution obtained from the same sample. The new distribution, we argue, can more accurately reflect the real welfare distribution in these countries and yield more precise estimates on inequality.

Using household-level data from Ghana, Kenya, Malawi, Niger, Nigeria and Uganda obtained from the Rural Income Generating Activities (RIGA) database and corresponding consumption aggregates estimated on the same survey, we re-estimate the top part of consumption distribution. For this purpose, we adapt a parametric multiple-imputation method (Jenkins et al., 2011) originally used to measure income inequality with right-censored (top-coded) data.

The re-estimation proceeds in three steps. First, we select the best fitting parametric model for the consumption and income distributions of each country. Second, we use the model's parameter estimates to derive imputed values for observations above two lower-bound consumption threshold typically used to define the middle class status in SSA: per capita daily consumption greater than \$4 in PPP US dollars and per capita daily consumption greater than \$10 in PPP US dollars. Finally, the combination of imputed values for observations above the consumption threshold and observed expenditures for below produces a partially synthetic data set on which we estimate inequality statistics such as the Gini coefficient, the MLD, and the Theil index.

Results show that in all six countries inequality increases substantially; depending on the threshold used, Gini on average increases by 20 percent compared to the original figures. In terms of levels, it is important to note that in four out of the six countries analyzed, inequality levels reach those of the traditionally unequal countries of the Southern cone (South Africa and neighbors). Further research is clerly needed, possibly having more waves of data for single country to gauge at trends. Nontheless, there are two preliminary conclusions we can already draw based on this

outcome.

First, against a general narrative on inequality as not being a major issue in Sub-Saharan Africa (see, *inter alia*, Pinhovskiy and Sala-i-Martin, 2014), our results indicate that inequality should become a central problem in SSA economic politicy debates as much as it has been in Latin America and in South Africa.

Second, these new inequality figures seem more in line with both development theory findings and what the present structure of these economies would suggest. In the first case, literature has pointed at the existence during the economic take-off of large productivity gaps between agriculture and non-agricultural sectors (e.g. Lewis, 1955, Kuznets, 1971, and Gollin et al. 2014). In SSA, where the agricultural sector still employs arund 50 percent (or more) of the labor force and where the production's aim is mainly the subsistence, this gap is expected to translate into relevant income differentials. Gap that has also an important spatial dimension since many SSA economies appear to be highly dualistic with longstanding spatial differences in terms of human capital, infrastructures and economic opportunities (Molini and Paci, 2015; World Bank, 2016).

## References

- African Development Bank. *The Middle of the Pyramid: Dynamics of the Middle Class in Africa*. Market brief, African Development Bank, Abidjan (Ivory Coast), 2011. Available at: <u>https://www.afdb.org/fileadmin/uploads/afdb/Documents/Publications/The%20Middle%2</u> <u>0of%20the%20Pyramid\_The%20Middle%20of%20the%20Pyramid.pdf</u>.
- H. Akaike. Information theory and an extension of the likelihood ratio principle. In B. N. Petrov and F. Csaki, editors, *Proceedings of the Second International Symposium of Information Theory*, pages 257–281. Akademiai Kiado, Budapest, 1973.
- F. Alvaredo. The rich in Argentina over the twentieth century: 1932–2004. In A. B. Atkinson and T. Piketty, editors, *Top Incomes: A Global Perspective*, pages 253–298. Oxford University Press, Oxford, 2010.
- A. B. Atkinson, T. Piketty, and E. Saez. Top incomes in the long run of history. *Journal of Economic Literature*, 49:3–71, 2011.

- G. Betti, A. D'Agostino, and A. Lemmi. Fuzzy monetary poverty measures under a Dagum income distributive hypothesis. In D. Chotikapanich, editor, *Modeling Income Distributions and Lorenz Curves*, pages 225–240. Springer, New York, 2008.
- G. Carletto, K. Covarrubias, B. Davis, M. Krausova, and P. Winters. *Rural Income Generating Activities Study: Methodological Note on the Construction of Income Aggregates*. Prepared for the Rural Income Generating Activities (RIGA) project of the Agricultural Development Economics Division, Food and Agriculture Organization, 2007. Available at: <a href="http://ftp.fao.org/docrep/fao/010/ai197e/ai197e00.pdf">http://ftp.fao.org/docrep/fao/010/ai197e/ai197e00.pdf</a>.
- F. Clementi. Heavy-tailed distributions for agent-based economic modelling. In A. Caiani, A. Russo, A. Palestrini, and M. Gallegati, editors, *Economics with Heterogeneous Interacting Agents: A Practical Guide to Agent-Based Modeling*, pages 157–190. Springer International Publishing AG, Cham, 2016.
- P. A. Corral Rodas, V. Molini, and G. O. Siwatu. *No Condition is Permanent: Middle Class in Nigeria in the Last Decade*, forthcoming on Journal of Development Studies.
- D. R. Cox and E. J. Snell. A general definition of residuals. *Journal of the Royal Statistical Society, Series B (Methodological)*, 30:248–275, 1968.
- C. Dagum. A new model of personal income distribution: specification and estimation. *Economie Appliquée*, 30:413–436, 1977.
- A. Deaton and S. Zaidi. Guidelines for Constructing Consumption Aggregates for Welfare Analysis. LSMS Working Paper 135, World Bank, Washington DC, 2002. Available at: <u>https://openknowledge.worldbank.org/handle/10986/14101</u>.
- K. E. Dynan, J. S. Skinner, and S. P. Zeldes. Do the rich save more? *Journal of Political Economy*, 112:397–444, 2004.
- R. C. Elandt-Johnson and N. L. Johnson. Survival Models And Data Analysis. John Wiley & Sons, Inc., Hoboken NJ, 3<sup>rd</sup> edition, 2003.
- M. Friedman. A Theory of the Consumption Function. Princeton University Press, Princeton NJ, 1957.
- D. T. Geithman. Middle class growth and economic development in Latin America. American

Journal of Economics and Sociology, 33:45–58, 1974.

- J. H. Goldthorpe. *Social Mobility and Class Structure in Modern Britain*. Clarendon Press, Oxford, 1987.
- D. Gollin, D. Lagakos, and M. E. Waugh. Agricultural productivity differences across countries. *American Economic Review*, 104:165–170, 2014.
- S. S. Golub and F. Hayat. *Employment, Unemployment and Underemployment in Africa*. Working Paper 2014/014, UNU-WIDER, Helsinki, 2014. Available at: <u>https://www.wider.unu.edu/publication/employment-unemployment-and-</u> <u>underemployment-africa</u>.
- V. Hlasny and P. Verme. Top incomes and the measurement of inequality in Egypt. *World Bank Economic Review*, DOI:<u>https://doi.org/10.1093/wber/lhw031</u>, 2016.
- M. Honorati and S. Johansson de Silva. *Expanding Job Opportunities in Ghana*. World Bank, Washington DC, 2016.
- ILO. *Resolution 1: Resolution Concerning Household Income and Expenditure Statistics*. Adopted at the 17<sup>th</sup> International Conference of Labour Statisticians, December, Geneva. Available at: <u>http://www.ilo.org/global/statistics-and-databases/standards-and-guidelines/resolutions-adopted-by-international-conferences-of-labour-statisticians/WCMS\_087503/lang--en/index.htm.</u>
- International Monetary Fund. *Sub-Saharan Africa: Maintaining Growth in an Uncertain World*. Regional economic outlook, International Monetary Fund, Washington DC, 2012. Available at: https://www.imf.org/external/pubs/ft/reo/2012/afr/eng/sreo1012.htm.
- T. Jappelli and L. Pistaferri. Fiscal policy and MPC heterogeneity. *American Economic Journal: Macroeconomics*, 6:107–136, 2014.
- S. P. Jenkins, R. V. Burkhauser, S. Feng, and J. Larrimore. Measuring inequality using censored data: a multiple-imputation approach to estimation and inference. *Journal of the Royal Statistical Society. Series A (Statistics in Society)*, 174:63–81, 2011.
- S. P. Jenkins. Fitting Singh-Maddala and Dagum distributions by maximum likelihood. *Stata Technical Bulletin*, 48:19–25, 1999.

- S. P. Jenkins. gb2fit: Stata Module to fit Generalized Beta of the Second Kind Distribution by Maximum Likelihood. Statistical Software Components Archive S456823, 2007. <u>http://ideas.repec.org/c/boc/bocode/s456823.html</u>.
- S. P. Jenkins. Distributionally-sensitive inequality indices and the GB2 income distribution. *Review of Income and Wealth*, 55:392–398, 2009.
- S. P. Jenkins. gb2lfit: Stata Module to fit Generalized Beta of the Second Kind Distribution by Maximum Likelihood (Log Parameter Metric). Statistical Software Components Archive S457897, 2014. https://ideas.repec.org/c/boc/bocode/s457897.html.
- C. Kleiber and S. Kotz. *Statistical Size Distributions in Economics and Actuarial Sciences*. John Wiley & Sons, New York NY, 2003.
- C. Kleiber. Dagum vs. Singh-Maddala income distributions. *Economics Letters*, 53:265–268, 1996.
- C. Kleiber. A guide to the Dagum distributions. In D. Chotikapanich, editor, *Modeling Income Distributions and Lorenz Curves*, pages 97–117. Springer, New York, 2008.
- S. Kuznets. *Economic Growth of Nations. Total Output and Production Structure*. Harvard University Press, Cambridge MA, 1971.
- A. Leigh and P. Van der Eng. Inequality in Indonesia: what can we learn from top incomes? *Journal of Public Economics*, 93:209–212, 2009.
- W. A. Lewis. The Theory of Economic Growth. Allen & Unwin, London, 1955.
- J. McCarthy. Imperfect insurance and differing propensities to consume across households. Journal of Monetary Economics, 36:301–327, 1995.
- J. B. McDonald and M. R. Ransom. The generalized beta distribution as a model for the distribution of income: estimation of related measures of inequality. In D. Chotikapanich, editor, *Modeling Income Distributions and Lorenz Curves*, pages 147–166. Springer, New York, 2008.
- J. B. McDonald and Y. J. Xu. A generalization of the beta distribution with applications. *Journal of Econometrics*, 66:133–152, 1995a.
- J. B. McDonald and Y. J. Xu. Errata. Journal of Econometrics, 69:427-428, 1995b.

- J. B. McDonald. Some generalized functions for the size distribution of income. *Econometrica*, 52:647–665, 1984.
- McKinsey. *Lions on the Move II: Realizing the Potential of Africa's Economies*. Report, McKinsey Global Institute, New York NY, 2016. Available at: <u>http://www.mckinsey.com/global-themes/middle-east-and-africa/lions-on-the-move-realizing-the-potential-of-africas-economies</u>.
- B. D. Meyer and J. X. Sullivan. The effects of welfare and tax reform: the material well-being of single mothers in the 1980s and 1990s. *Journal of Public Economics*, 88:1387–1420, 2004.
- B. Milanovic and S. Yitzhaki. Decomposing world income distribution: does the world have a middle class? *Review of Income and Wealth*, 48:155–178, 2002.
- B. Milanovic. The Haves and the Have-Nots: A Brief and Idiosyncratic History of Global Inequality. Basic Books, New York NY, 2010.
- V. Molini and P. Paci. *Poverty Reduction in Ghana: Progress and Challenges*. Technical report, World Bank, Washington DC, 2015. Available at https://openknowledge.worldbank.org/handle/10986/22732.
- M. Pinkovskiy and X. Sala-i-Martin. Africa is on time. *Journal of Economic Growth*, 19:311–338, 2014.
- C. Quintano and A. D'Agostino. Studying inequality in income distribution of single-person households in four developed countries. *Review of Income and Wealth*, 52:525–546, 2006.
- E. J. Quiñones, A. P. de la O-Campos, C. Rodríguez-Alas, T. Hertz, and P. Winters. *Methodology for Creating the RIGA-L Database*. Prepared for the Rural Income Generating Activities (RIGA) project of the Agricultural Development Economics Division, Food and Agriculture Organization, 2009. Available at: <a href="http://www.fao.org/fileadmin/templates/riga/docs/Country\_survey\_information/RIGA-L\_Methodology.pdf">http://www.fao.org/fileadmin/templates/riga/docs/Country\_survey\_information/RIGA-L\_Methodology.pdf</a>.
- J. P. Reiter. Inference for partially synthetic, public use microdata sets. *Survey Methodology*, 29:181–188, 2003.
- C. Sanhueza and R. Mayer. Top incomes in Chile using 50 years of household surveys: 1957–2007.

Estudios de Economía, 38:169–193, 2011.

- C. Schluter and M. Trede. Tails of Lorenz curves. Journal of Econometrics, 109:151–166, 2002.
- G. E. Schwarz. Estimating the dimension of a model. Annals of Statistics, 6:461–464, 1978.
- S. K. Singh and G. S. Maddala. A function for size distribution of incomes. *Econometrica*, 44:963–970, 1976.
- A. Tarozzi. Calculating comparable statistics from incomparable surveys, with an application to poverty in India. *Journal of Business & Economic Statistics*, 25:314–336, 2007.
- World Bank. Nigeria Poverty Assessment: Poverty Reduction in the Last Decade. Technical report,
   World Bank, Washington DC, 2016. Available at:
   https://openknowledge.worldbank.org/handle/10986/25825.

# Tables

Table 1: Distributional summary statistics for the consumption and income variables used in our analysis.

1000	Year	Households <sup>a</sup>	(	Consum	ption		Income				
Country			$Mean^b$	Gini	MLD	Theil	Mean <sup>b</sup>	$\operatorname{Gini}$	MLD	Theil	
Ghana	2005	7,659	563	0.42	0.31	0.32	98	0.60	0.76	0.70	
Kenya	2005	11,700	38,049	0.52	0.47	0.54	28,495	0.67	0.97	1.04	
Malawi	2011	11,712	38,524	0.47	0.38	0.46	19,010	0.60	0.69	0.82	
Niger	2011	3,729	$214,\!434$	0.31	0.16	0.17	$105,\!140$	0.49	0.45	0.46	
Nigeria	2004	16,922	40,869	0.40	0.27	0.28	60,707	0.65	0.93	1.11	
Uganda	2005	7,199	$450,\!564$	0.47	0.37	0.46	$265,\!182$	0.61	0.72	0.81	

 $^{\rm a}$  Effective number of observations after removal of missing and non-positive values.  $^{\rm b}$  Money amount in local currency, 2005 prices.

Country Model <sup>a</sup> rataneous connacco connacco connacco rataneous connector connector	Predictions <sup>d</sup>			
$\hat{a}_c$ $\hat{b}_c$ $\hat{p}_c$ $\hat{q}_c$ $\ln L$ AIC BIC Mean Gini M	D Theil			
Ghana GB2 1.34 396 2.39 2.28 -54,989 109,985 110,013 563 0.42	31 0.32			
(0.14) (22) $(0.44)$ $(0.41)$				
D 2.25 393 1.0855,005 110,016 110,037 582 0.44	34 0.39			
(0.04) (13) (0.05) (33)*				
SM 2.30 420 - 1.02 -55,007 110,019 110,040 573 0.43	33 0.37			
$(0.04)$ (13) $(0.05)$ $(36)^*$				
Kenya GB2 1.38 14,673 2.27 1.38 -133,307 266,622 266,651 39,400 0.53	51 0.66			
(0.09) $(782)$ $(0.28)$ $(0.15)$				
D $1.71$ $15,619$ $1.62$ - $-133,312$ $266,630$ $266,652$ $41,270$ $0.56$	55 0.79			
(0.02) (576) $(0.07)$ (11)*				
SM 2.26 17.766 - 0.70 -133,338 266,682 266,704 42,900 0.57	59 0.94			
(0.04) (428) $(0.02)$ (63)*				
Malawi GB2 1.34 9.741 4.71 1.56 -132.809 265.627 265.656 38.790 0.47	38 0.49			
(0.11) $(1.415)$ $(1.00)$ $(0.19)$				
D $1.84$ $13.594$ $2.40$ - $-132.818$ $265.643$ $265.665$ $40.280$ $0.49$	42 0.59			
(0.02) (601) (0.14) (18)*				
SM 2.88 18.454 - 0.58 -132.874 265.753 265.775 41.570 0.51	45 0.71			
$(0.05)$ $(338)$ $(0.02)$ $(129)^*$				
Niger GB2 2.29 115,537 3.03 1.34 -48,047 96,101 96,126 215,200 0.31	16 0.18			
(0.32) $(11,566)$ $(0.89)$ $(0.27)$				
D $2.81$ $126,824$ $2.07$ $ -48,048$ $96,102$ $96,120$ $216,700$ $0.32$	16 0.20			
(0.06) $(5,840)$ $(0.18)$ $(2)$				
$\mathrm{SM} \qquad 4.18 \qquad 146,805 \qquad - \qquad 0.61 \qquad -48,059 \qquad 96,125 \qquad 96,144 \qquad 218,100 \qquad 0.32 \qquad - \qquad 0.61 \qquad -48,059 \qquad - \qquad 0.61 \qquad -48,059 \qquad - \qquad 0.61 \qquad - \qquad 0.6$	17 0.21			
$(0.12) (3,237) (0.03) (26)^*$				
Nigeria GB2 1.04 33,094 3.75 3.98 -193,859 387,726 387,757 40,820 0.40	27 0.28			
(0.10) $(1,914)$ $(0.63)$ $(0.68)$				
D 2.36 30,366 1.04193,940 387,885 387,908 42,430 0.42	31 0.35			
(0.03) (641) (0.03) (161)*				
SM 2.30 33,613 - 1.13 -193,933 387,872 387,895 41,460 0.41	29 0.32			
(0.03) (741) $(0.04)$ (148)*				
Uganda GB2 1.36 89.272 6.04 1.47 -99.233 198.475 198.503 458.200 0.48	39 0.52			
(0.15) $(22.145)$ $(1.99)$ $(0.23)$				
D 1.80 133.130 3.0799.237 198.481 198.501 474.700 0.50	42 0.62			
(0.03) (8,631) (0.26) (8)*				
SM 3.14 201,493 - 0.51 -99,276 198,558 198,579 502,000 0.53	48 0.82			
(0.07) (4.255) $(0.02)$ (85)*				

Table 2: Maximum likelihood estimation of generalized beta models for consumption distributions.

<sup>a</sup> GB2 = generalized beta II; D = Dagum; SM = Singh-Maddala.
 <sup>b</sup> Numbers in parentheses: estimated standard errors.
 <sup>c</sup> Numbers in parentheses: likelihood ratio statistics.
 <sup>d</sup> Analytic values obtained by substituting the estimated parameters into the relevant expressions; the formulas for the generalized beta II, Dagum and Singh-Maddala distributions can be found in Kleiber and Kotz (2003, ch. 6) and Jenkins (2009).

 $^{\ast}$  Denotes significance at the 5% level.

Country	Modela		Parameter estimates <sup>b</sup>				rison fit statist	ics <sup>c</sup>	Predictions <sup>d</sup>			
country	modd	$\hat{a}_y$	$\hat{b}_y$	$\hat{p}_y$	$\hat{q}_y$	$\ln L$	AIC	BIC	Mean	Gini	MLD	Theil
Ghana	GB2	0.72	133	2.16	4.09	-42,132	84,271	84,299	100	0.61	0.77	0.75
	D	(0.08)	(33)	(0.38)	(1.05)	40 166	04 990	01 950	100	0.00	0.07	1.90
	D	(0.03)	(3)	(0.03)	-	-42,100 (69)*	64,559	64,559	122	0.08	0.97	1.50
	SM	1.21	79	(0.00)	1.55	-42.147	84.300	84.321	106	0.63	0.84	0.93
	(000.000)	(0.02)	(6)		(0.09)	(31)*	- ,	,				0.00
Kenya	GB2	1.06	15,938	1.16	1.54	-129,137	258,283	258,312	30,420	0.69	1.04	1.29
		(0.07)	(1,007)	(0.11)	(0.17)							
	D	1.38	13,864	0.82	10 <u>10</u> 2	-129,146	258,297	258,320	35,500	0.74	1.19	1.91
		(0.02)	(513)	(0.03)	10 Stells	$(17)^*$			Income Repaired	1000.020-00.020	(1) Presidente	
	SM	1.17	15,373		1.31	-129,139	258,283	258,305	31,490	0.71	1.07	1.43
		(0.02)	(797)		(0.05)	(3)						
Malawi	GB2	2.25	10,661	0.61	0.68	-125,144	250,296	250,326	20,300	0.62	0.75	1.12
		(0.15)	(271)	(0.05)	(0.06)				2			
	D	1.72	11,187	0.85	—	-125,154	250,313	250,335	18,850	0.59	0.68	0.87
		(0.02)	(314)	(0.03)		$(19)^*$						
	$\mathbf{SM}$	1.56	10,714	100 - 100 100-100	1.11	-125,162	250,329	250,352	19,000	0.60	0.68	0.87
		(0.02)	(357)		(0.04)	$(35)^*$						
Niger	GB2	1.97	89,747	0.80	1.12	-46,550	93,109	93,134	106,400	0.50	0.46	0.50
0		(0.22)	(5,262)	(0.12)	(0.19)	•	1999 <b>-</b>					
	D	2.12	87,474	0.73	<u> </u>	-46,551	93,107	93,126	107,300	0.50	0.47	0.52
		(0.06)	(3,552)	(0.04)		(0)						
	$\mathbf{SM}$	1.68	93,458	-	1.44	-46,551	93,109	93,128	105,400	0.49	0.45	0.47
		(0.04)	(5,929)		(0.10)	(2)						
Nigeria	GB2	1.54	58,798	0.59	1.27	-199,988	399,984	400,015	56,370	0.62	0.86	0.84
		(0.08)	(2,583)	(0.04)	(0.11)							
	D	1.79	53,570	0.49		-199,992	399,991	400,014	58,060	0.63	0.89	0.93
		(0.03)	(1,207)	(0.01)		$(8)^{*}$						
	$\mathbf{SM}$	1.04	78,950	100 to 100	2.41	-200,016	400,039	400,062	55,030	0.61	0.83	0.76
		(0.01)	(4,341)		(0.10)	$(56)^*$						
Uganda	GB2	2.20	112,962	0.70	0.60	-95,809	191,626	191,653	343,000	0.70	0.97	1.96
0.41-002-002-04		(0.19)	(3,680)	(0.08)	(0.07)							
	D	1.54	115,801	1.11	10 N	-95,818	$191,\!642$	191,663	287,000	0.64	0.80	1.19
		(0.03)	(4,835)	(0.05)		$(18)^{*}$						
	SM	1.69	110,287		0.85	-95,814	191,634	$191,\!654$	308,000	0.67	0.86	1.44
		(0.03)	(4, 198)		(0.03)	$(10)^*$						

Table 3: Maximum likelihood estimation of generalized beta models for income distributions.

<sup>a</sup> GB2 = generalized beta II; D = Dagum; SM = Singh-Maddala. <sup>b</sup> Numbers in parentheses: estimated standard errors. <sup>c</sup> Numbers in parentheses: likelihood ratio statistics.

<sup>d</sup> Analytic values obtained by substituting the estimated parameters into the relevant expressions; the formulas for the generalized beta II, Dagum and Singh-Maddala distributions can be found in Kleiber and Kotz (2003, ch. 6) and Jenkins (2009).

\* Denotes significance at the 5% level.

			\$4 per day					\$10 per day					
Country	Model	Threshold <sup>a</sup>	Middle class size (%)	$\operatorname{Gini}^{\mathrm{b}}$	$\mathrm{MLD}^{\mathrm{b}}$	Theil <sup>b</sup>	Threshold <sup>a</sup>	Middle class size $(\%)$	$\operatorname{Gini}^{\mathrm{b}}$	$\mathrm{MLD}^{\mathrm{b}}$	Theil <sup>b</sup>		
Ghana	GB2	645	27.22	0.46	0.37	0.42	$1,\!613$	4.04	0.43	0.33	0.37		
Kenya	GB2	29,015	38.56	0.60	0.65	0.96	72,537	11.56	0.59	0.63	0.95		
Malawi	GB2	78,259	8.41	0.53	0.49	0.83	195,646	1.69	0.51	0.46	0.75		
Niger	Dagum	330,811	12.55	0.36	0.22	0.31	827,027	1.03	0.33	0.18	0.23		
Nigeria	GB2	74,745	10.97	0.47	0.38	0.51	186,863	0.81	0.41	0.30	0.35		
Uganda	GB2	894,536	9.18	0.60	0.65	1.28	$2,\!236,\!341$	1.65	0.54	0.52	0.99		

Table 4: Inequality estimates from partially synthetic data sets, by country, definition of the middle class, and index.

<sup>a</sup> Annualized money amount in local currency, 2005 prices and PPP. <sup>b</sup> Average of the point estimates derived from each of the R = 999 partially synthetic data sets.

## Figures



Figure 1: Log-log complementary distributions of household consumption, consumption plus savings, and income for Ghana.



Figure 2: Cox-Snell residuals of generalized beta models for consumption distributions.



Figure 3: Cox-Snell residuals of generalized beta models for income distributions.



**Figure 4:** Comparison between the sample values of chosen distributional indicators and their counterparts implied by the best fitting consumption models.



**Figure 5:** Comparison between the sample values of chosen distributional indicators and their counterparts implied by the best fitting income models.



**Figure 6:** Inequality estimates derived from repetition of the imputation process R = 999 times, by index and country. The height of the bars is the level of inequality estimated using the original consumption data. The top cap of the spikes denotes the multiple-imputation point estimate derived using a minimum threshold of \$4 per day, whereas the bottom cap shows the estimate for each of the three indices derived using an absolute definition of the middle class with per capita daily consumption greater than \$10.







