

Pudney, Stephen; Francavilla, Francesca

**Working Paper**

## Income mis-measurement and the estimation of poverty rates: an analysis of income poverty in Albania

ISER Working Paper Series, No. 2006-35

**Provided in Cooperation with:**

Institute for Social and Economic Research (ISER), University of Essex

*Suggested Citation:* Pudney, Stephen; Francavilla, Francesca (2006) : Income mis-measurement and the estimation of poverty rates: an analysis of income poverty in Albania, ISER Working Paper Series, No. 2006-35, University of Essex, Institute for Social and Economic Research (ISER), Colchester

This Version is available at:

<https://hdl.handle.net/10419/92057>

**Standard-Nutzungsbedingungen:**

Die Dokumente auf EconStor dürfen zu eigenen wissenschaftlichen Zwecken und zum Privatgebrauch gespeichert und kopiert werden.

Sie dürfen die Dokumente nicht für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, öffentlich zugänglich machen, vertreiben oder anderweitig nutzen.

Sofern die Verfasser die Dokumente unter Open-Content-Lizenzen (insbesondere CC-Lizenzen) zur Verfügung gestellt haben sollten, gelten abweichend von diesen Nutzungsbedingungen die in der dort genannten Lizenz gewährten Nutzungsrechte.

**Terms of use:**

*Documents in EconStor may be saved and copied for your personal and scholarly purposes.*

*You are not to copy documents for public or commercial purposes, to exhibit the documents publicly, to make them publicly available on the internet, or to distribute or otherwise use the documents in public.*

*If the documents have been made available under an Open Content Licence (especially Creative Commons Licences), you may exercise further usage rights as specified in the indicated licence.*



# **Income Mis-Measurement and the Estimation of Poverty Rates**

**An Analysis of Income Poverty in Albania**

**Stephen Pudney and Francesca Francavilla**

ISER Working Paper  
2006-35

**Acknowledgement:** This work was supported by the Italian National Research Centre (CNR) and the UK Economic and Social research Council through project “Fertility and poverty in developing countries” (award no. RES000230462) and the ULSC and MiSoC Research Centres (award nos. H562255004 and RES518285001).

Readers wishing to cite this document are asked to use the following form of words:

**Pudney, Stephen and Francavilla, Francesca (July 2006) ‘Income Mis-measurement and the Estimation of Poverty Rates. An Analysis of Income Poverty in Albania’, ISER Working Paper 2006-35. Colchester: University of Essex.**

The on-line version of this working paper can be found at <http://www.iser.essex.ac.uk/pubs/workpaps/>

The Institute for Social and Economic Research (ISER) specialises in the production and analysis of longitudinal data. ISER incorporates

- MISOC (the ESRC Research Centre on Micro-social Change), an international centre for research into the lifecourse, and
- ULSC (the ESRC UK Longitudinal Studies Centre), a national resource centre to promote longitudinal surveys and longitudinal research.

The support of both the Economic and Social Research Council (ESRC) and the University of Essex is gratefully acknowledged. The work reported in this paper is part of the scientific programme of the Institute for Social and Economic Research.

Institute for Social and Economic Research, University of Essex, Wivenhoe Park,  
Colchester. Essex CO4 3SQ UK  
Telephone: +44 (0) 1206 872957 Fax: +44 (0) 1206 873151 E-mail: [iser@essex.ac.uk](mailto:iser@essex.ac.uk)  
Website: <http://www.iser.essex.ac.uk>

© July 2006

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system or transmitted, in any form, or by any means, mechanical, photocopying, recording or otherwise, without the prior permission of the Communications Manager, Institute for Social and Economic Research.

## ABSTRACT

We investigate the reliability of estimated income poverty profiles for Albanian survey data. We find evidence that a significant number of households with low reported incomes have relatively high living standards and are consequently misclassified as poor. We extend the theory of contaminated distributions to incorporate direct measures of well-being as indicators of data contamination, and develop a new nonparametric approach for constructing bounds on conditional poverty rates. We find very large upward biases in measured income poverty under the assumption of independence between living standards and the misreporting propensity, but a wide range of uncertainty under more general conditions.

## NON-TECHNICAL SUMMARY

Empirical analyses of income poverty are often used to guide the design of welfare policies such as social assistance programmes and employment subsidies. The typical poverty study is based on survey data covering a large number of randomly-selected households, whose members are interviewed in detail about their incomes and other circumstances. This information is first used to classify households as 'poor' or 'non-poor' in relation to a poverty line, then statistical methods are used to analyse the way that poverty varies between different social groups. Groups which are identified as particularly vulnerable to poverty then become very important for anti-poverty policy.

An accepted problem with this type of analysis is the accurate survey measurement of household income. Researchers have sometimes reported significant numbers of survey respondents with implausibly low measured incomes, in relation to other indicators of the household's standard of living. We investigate this issue, using data from recent cross-section and panel surveys of Albanian households. We find very strong evidence that measured income greatly understates true living standards for a substantial group of households below the poverty line in the 2002 Albanian Living Standards Measurement Survey. These problem cases are identified in two ways: through a discrepancy between the position of the household in the measured income distribution and its position in the distribution of measured consumption expenditure; and through a discrepancy between its income and other non-financial indicators of living standards (such as the number of durables owned or the household's own subjective assessment of well-being).

Using a smaller panel dataset for 2002 and the later years 2003-4 which have no consumption information, we find that the income measurement problem is most serious in 2002 and that there is only a weak tendency for households who 'misreported' their incomes in 2002 also to show signs of mis-measurement in later years. Thus, these difficulties with income measurement do not appear to be very persistent over time, at the level of individual households.

The paper also develops a new method of taking account of this measurement problem in making poverty analyses; this is applied to the 2002 cross-section survey. We find that, under the (strong) assumption that income mis-measurement is unrelated to living standards, 'true' poverty rates are clearly much lower than measured rates. However, if this assumption is relaxed, to allow the possibility that understatement of income is mainly confined to households with high living standards, then we can only draw a much weaker conclusion - that there is a very wide range of uncertainty associated with measured poverty rates, and measured poverty *may* therefore be much too high.

## 1 Introduction

An enduring concern in the literature on poverty and income dynamics is the quality of income data used as a basis for the measurement of poverty. Particularly in the context of developing and transition countries which have large informal and non-market sectors, consumption is often preferred to income for this purpose. There are good reasons to believe that consumption expenditure is observed with less error than income for poor households and that expenditure data are less vulnerable to relatively unimportant short-term fluctuations. However, consumption data are often not available or only partially available, particularly in panel surveys that allow the study of change over time. This is true of most panel studies in developed countries, including the PSID and SIPP in the USA, the BHPS in the UK, SOEP in Germany and HILDA in Australia. It is also the case for all but the first wave of the Albanian panel used here.

Given that there is often no alternative but to use income data for poverty analysis, it is important to try to understand the nature of measurement problems and think about the likely impact of these problems on conventional poverty measures. A further possibility is to devise methods for adjusting poverty measures to reduce the biases caused by income mis-measurement.

There are two main strands in the statistical literature on income measurement error, based on distinct models of the mis-measurement process. One adopts the classical errors-in-variables approach, which assumes that measured income (or its logarithm) is the sum of true income and a zero-mean measurement error drawn from some continuous distribution. Under standard specifications, this approach implies that essentially every income observation is affected by measurement error to some extent. Chesher and Schluter (2002) use this framework, deriving small- $\sigma$  approximations to the bias in a range of poverty measures. An alternative approach, used in the derivation of robust estimators for data contaminated by outliers (Huber, 1964), is to assume that only a proportion of observations is significantly affected by mis-measurement. This is the approach adopted by Cowell and Victoria-Feser (1996) in their study of robust estimation of inequality indices.

Our analysis is closest to the latter approach, since the contamination model appears more consistent with our income data. However, we generalise the approach

considerably by introducing into the analysis other variables which are informative for the contamination process. The most promising variables to use for this purpose are direct indicators of well-being, of the kind which are often used in the construction of deprivation indices. However, it should be noted that we use these variables as indicators of cases where measured income may be misleading, not as the empirical counterpart of a broader, multi-dimensional concept of poverty. The enduring debate about the meaning of poverty and the value and interpretation of deprivation indices as measures of welfare (see Sen, 1985; Ringen, 1988, Ravallion and Lokshin, 2001; McKay, 2004; Berthoud *et. al.*, 2004), is therefore not directly relevant to our analysis.

The paper has five objectives. Firstly, we investigate the evidence for income mis-measurement and draw some conclusions about its nature, using a 2002 dataset from Albania. An advantage of the Albanian dataset is that it contains a panel element and, in its first wave, provides data on income, consumption expenditure and other direct indicators of well-being. Subsequent waves have no consumption information. Our second objective is to evaluate the power of direct well-being measures as a source of identifying information on the incidence of mis-measurement and to use those variables to assess the persistence of mis-measurement through time, using the panel observations. A third aim is to extend the contamination model of income mis-measurement to allow the use of external indicators of error and to develop a nonparametric approach to the adjustment of sample poverty rates. This is distinct from the Cowell and Victoria-Feser (1996) approach, which uses robust estimators of parametric models and assumes mean-preserving contamination. A final objective is to give an empirical indication of the nature and range of uncertainty associated with income poverty rates, and the conditions under which the use of well-being indicators can resolve this uncertainty. We begin by describing the Albanian datasets used in this study.

## **2 Data**

Our analysis is based on two overlapping datasets relating to the period 2002-4. Our main focus is on the large cross-section Albanian Living Standards Measurement Survey (ALSMS) carried out in 2002 (see Aassve *et. al.* (2005) and Carletto and

Zeza, (2004) for further details) and this is supplemented by a (balanced) panel dataset following a subset of the ALSMS sample in 2003 and 2004 (see Azzarri *et. al.*, 2006).

## **2.1 The Albanian Living Standards Measurement Survey (ALSMS)**

The ALSMS was first conducted under the auspices of the World Bank in 2002 as a cross-section sample of 3600 households. The Republic of Albania is divided geographically into 12 Prefectures. The latter are divided into Districts which are, in turn, divided into Cities and Communes. The Communes contain all the rural villages and the very small cities. The sampling frame was divided in four regions (strata), Coastal Area, Central Area, and Mountain Area, and Tirana and these four strata were further divided into major cities, other urban, and other rural cities and villages were divided into Enumeration Areas (EAs), which formed the basis for the LSMS sampling frame. The sample was drawn from 450 EAs and, in each of these, eight households were selected. Although probabilities within strata were (approximately) equal, probabilities varied greatly between the strata. Notably, the mountain region was heavily over-represented and the Central Rural region was under-represented in the sample. In order to obtain correct estimates the data need to be weighted.

An individual is defined as a household member if he or she was not away from the household for more than six months. Lodgers, hired workers and servants are not included. The head of the household is considered a household member if he or she had been away less than 12 months, rather than the 6 month limit for anyone else absent. The survey is of the standard LSMS format (Ghosh and Glewwe, 1995, 2000). It includes information on consumption expenditure, income, health, education, and employment. The household questionnaire included additional modules on migration, fertility history, dwelling, utilities, durables, subjective poverty, agriculture and non-farm enterprises.

## **2.2 The Albanian panel**

The Albanian panel survey sample was conducted by the Albania Institute of Statistics (INSTAT) with technical assistance from the World Bank and the Institute for Social



and Economic Research of the University of Essex (ISER). Households in the panel were selected from those interviewed in the 2002 LSMS and were re-interviewed in 2003 and 2004. The panel followed approximately half the LSMS households and was designed to provide a nationally representative sample of households and individuals within Albania. The balanced panel we use in our analysis includes 1,682 households. Like the ALSMS, the panel includes information on education, health, employment, migration, fertility, dwelling, income and utilities but, for reasons of cost, does not include consumption expenditure in waves 2 and 3. This prevents the use of a consumption-based analysis of poverty dynamics. Other elements of the questionnaire were redesigned at waves 2 and 3, giving a generally lower degree of detail, particularly for some components of income.

### **3 The distribution of income, consumption expenditure and deprivation**

#### **3.1 Definitions of variables**

The income variable for the first wave was constructed by a World Bank team (Carletto and Zezza, 2004) and we adopt their variable without modification. It represents total monthly income of the household and includes labour market income (wages, in-kind salaries and job-related bonuses), income from non-agricultural business, agricultural income, private and public transfers and other income such as rents, inheritance, gambling. The income definition includes the imputed value of self-produced consumption goods. Income variables for waves two and three were constructed in line with the same income definition. However, there were important changes in the income questionnaire at each wave, which generates potential discontinuities in income measurement between all waves. The income variables for waves one and three are comparable in principle but there is an important difference in implementation because of a reduction in questionnaire detail on business income after wave one. The income variable for wave two is not directly comparable with those for waves one and three, because questions on some income components (non-public transfers and the residual ‘other income’ category) were omitted from the questionnaire in that year. Thus, it is important to avoid attaching too much importance to changes in absolute income levels between waves, especially for

changes involving wave 2. For this reason, we prefer to use relative measures based on change in the percentile position in the within-year income distribution.

*Our consumption variable includes food and non-food expenses (clothing, household supplies for cleaning, tobacco, household articles, entertainment, services, etc.) and utilities (electricity, gas, telephone services, water and fuels) but excludes payments of rent and durable goods in order to avoid problems of dynamic adjustment. The consumption variable also includes the imputed value of self-produced food.*

*Income and consumption variables are deflated to 2002 prices, using the Consumer Price Index, as published by the IMF.<sup>1</sup> We also adjusted for the large spatial variations in price levels, using a Paasche index (at the level of primary sampling units) constructed by the World Bank for wave one and described in the survey documentation. We are not able to recalculate this spatial price index for waves two and three since consumption data were not collected at those waves. Nevertheless it seems reasonable to assume that geographical price differentials did not change greatly over the three years we consider.*

A wide range of well-being indicators exist in the ALSMS. Questions on the description of household dwellings, utilities and durables and on subjective poverty would be crucial for any study of household deprivation but our aim is different: to develop indicators of the incidence of income misreporting rather than represent a wider concept of poverty or deprivation. We use four questions that are present in the 2002 survey: the number of durables owned from a list of eleven possibilities; general satisfaction with current circumstances; adequacy of present consumption of food; and a self-assessment of the household's distributional position. These variables are defined in detail in the appendix.

### **3.2 Evidence on income contamination**

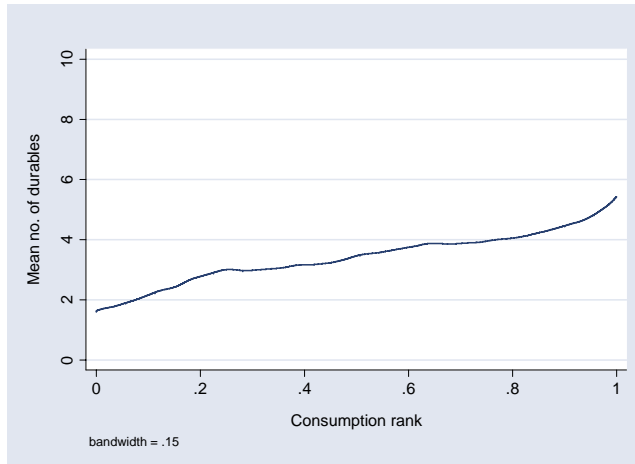
If income and consumption are good indicators of the resources of the household, then we would expect to observe a monotonically increasing statistical relationship between the standard of living and either the income or consumption rank. Figures 1 and 2 illustrate this relationship graphically, using each of the four indicators of the household's standard of living. They show nonparametric locally-weighted regressions of each indicator on the consumption or income rank. These regressions use a moving window of  $0.15n$  observations, where  $n$  is sample size.

---

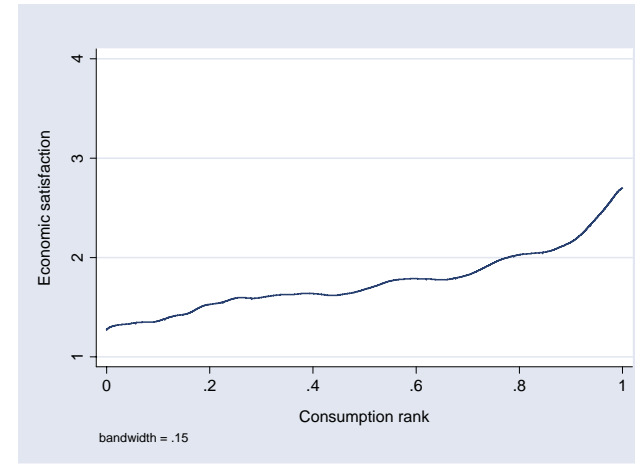
<sup>1</sup> *World Economic Outlook, September 2004.*

Figure 1 confirms the expected monotonic relationship between living standards and consumption, whichever indicator is used for the former. The evidence for income is quite different, since there is a strong negative gradient in the bottom income decile. This is clear evidence of the existence of a group of households whose measured incomes (but not consumption expenditures) are contaminated by a large negative measurement error, with respect to any normal concept of a standard of living. We refer to such errors as income contamination. Note that it may be possible in some cases for income to be contaminated in our sense despite being accurately measured in some particular accounting sense. We are concerned here with any large deviation of measured income from a 'meaningful' indication of living standards.

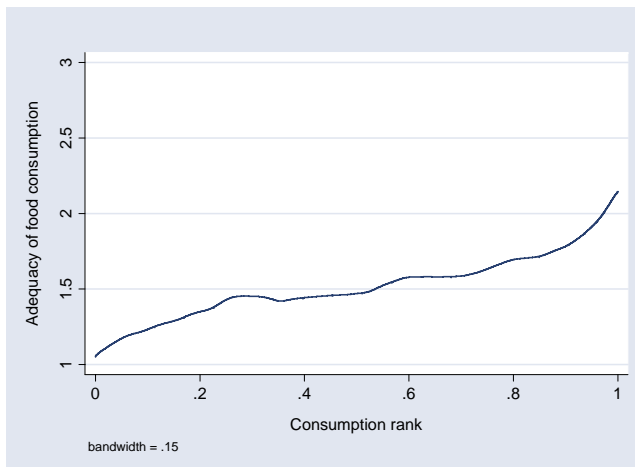
Income contamination is not a phenomenon unique to Albania or to developing countries. For instance, Berthoud *et. al.* (2004) and Saunders (2005) report similar, albeit weaker, phenomena in UK and Australian data respectively. Note that the negative gradient for low incomes and positive gradient for high incomes is not what we would expect to observe if income were contaminated by a simple additive measurement error: an attenuated but monotonic relationship would be typical of that case.



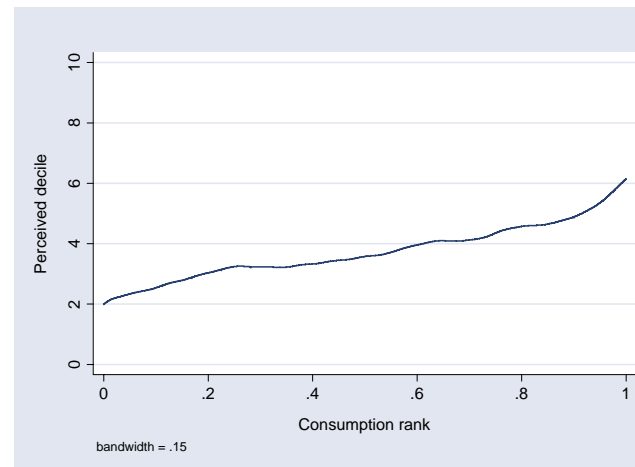
(a) Ownership of durables



(b) General satisfaction

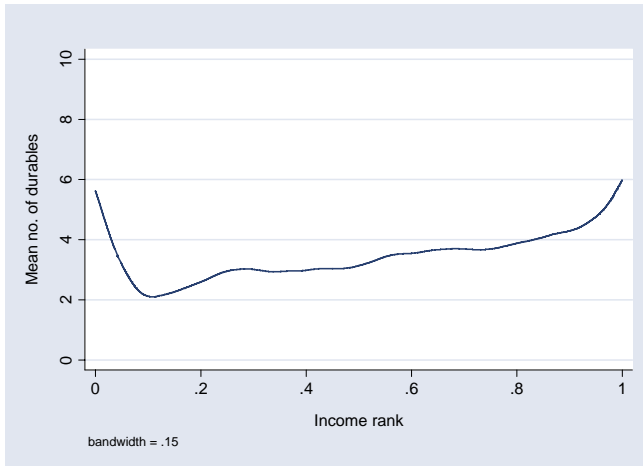


(c) Adequacy of food consumption

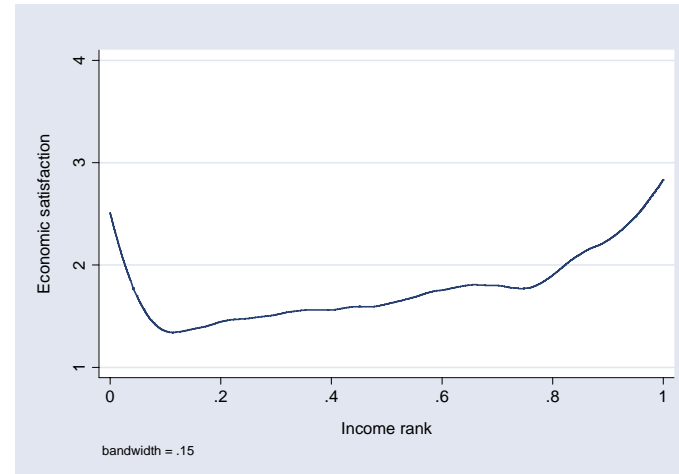


(d) Perceived place in living standards distribution

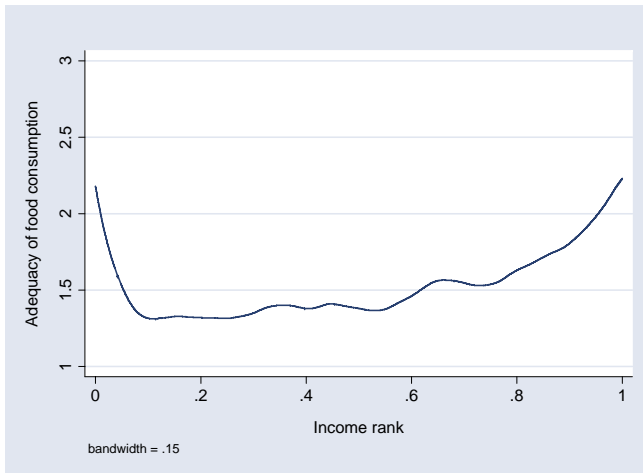
**Figure 1** Nonparametric locally-weighted regressions of living standards on consumption rank



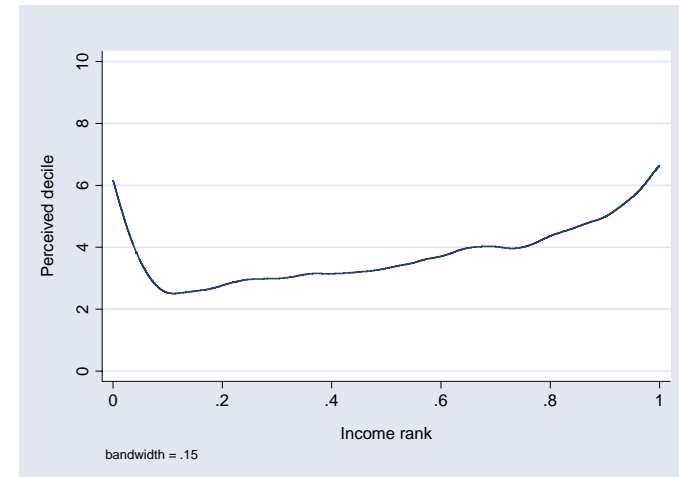
(a) Ownership of durables



(b) General satisfaction



(c) Adequacy of food consumption



(d) Perceived place in living standards distribution

**Figure 2** Nonparametric locally-weighted regressions of living standards on income rank

Table 1 shows that the perverse discrepancies in standard of living between the first and second income deciles are to be found among people in the top six consumption deciles. Thus the nature of the income measurement problem is such that there are large displacements into the bottom income decile, with no converse displacement in the top decile. Conventional measurement error models would not produce this effect.

**Table 1** Ownership of durables by income and consumption deciles  
(Mean number of durables owned in Roman type; standard errors in Italics)

		Income decile										Total
		1	2	3	4	5	6	7	8	9	10	
Consumption decile	1	1.5 <i>0.2</i>	1.5 <i>0.1</i>	2.2 <i>0.2</i>	2.1 <i>0.2</i>	2.3 <i>0.2</i>	1.8 <i>0.2</i>	3.2 <i>0.6</i>	2.2 <i>0.6</i>	3.3 <i>0.3</i>	-	1.9 <i>0.1</i>
	2	1.8 <i>0.2</i>	2.0 <i>0.2</i>	2.5 <i>0.2</i>	2.8 <i>0.2</i>	2.6 <i>0.2</i>	2.4 <i>0.3</i>	3.2 <i>0.2</i>	2.7 <i>0.3</i>	2.9 <i>0.5</i>	1.8 <i>0.3</i>	2.4 <i>0.1</i>
	3	3.0 <i>0.3</i>	2.8 <i>0.2</i>	3.0 <i>0.2</i>	2.8 <i>0.2</i>	3.1 <i>0.2</i>	3.4 <i>0.2</i>	3.3 <i>0.3</i>	3.1 <i>0.3</i>	2.9 <i>0.4</i>	3.3 <i>0.8</i>	3.0 <i>0.1</i>
	4	2.4 <i>0.4</i>	2.4 <i>0.3</i>	3.3 <i>0.2</i>	2.9 <i>0.3</i>	2.9 <i>0.2</i>	3.5 <i>0.2</i>	3.6 <i>0.3</i>	3.1 <i>0.3</i>	3.2 <i>0.3</i>	3.6 <i>0.4</i>	3.1 <i>0.1</i>
	5	3.4 <i>0.3</i>	2.6 <i>0.2</i>	2.9 <i>0.2</i>	2.9 <i>0.2</i>	2.7 <i>0.2</i>	3.8 <i>0.2</i>	3.6 <i>0.2</i>	3.5 <i>0.2</i>	3.2 <i>0.3</i>	4.7 <i>0.3</i>	3.2 <i>0.1</i>
	6	3.7 <i>0.3</i>	3.0 <i>0.4</i>	3.0 <i>0.3</i>	3.7 <i>0.3</i>	3.1 <i>0.2</i>	3.5 <i>0.2</i>	4.1 <i>0.2</i>	3.7 <i>0.2</i>	4.0 <i>0.2</i>	4.2 <i>0.3</i>	3.6 <i>0.1</i>
	7	4.4 <i>0.5</i>	3.1 <i>0.3</i>	3.5 <i>0.3</i>	3.7 <i>0.3</i>	3.6 <i>0.5</i>	3.4 <i>0.2</i>	3.9 <i>0.2</i>	3.8 <i>0.2</i>	4.3 <i>0.2</i>	4.6 <i>0.2</i>	3.9 <i>0.1</i>
	8	4.1 <i>0.4</i>	2.9 <i>0.6</i>	3.4 <i>0.4</i>	3.2 <i>0.3</i>	3.5 <i>0.3</i>	3.8 <i>0.2</i>	3.7 <i>0.2</i>	4.1 <i>0.2</i>	4.3 <i>0.2</i>	4.8 <i>0.2</i>	3.9 <i>0.1</i>
	9	4.8 <i>0.3</i>	3.0 <i>0.8</i>	3.8 <i>0.4</i>	3.3 <i>0.4</i>	4.2 <i>0.3</i>	4.1 <i>0.3</i>	3.6 <i>0.2</i>	4.3 <i>0.2</i>	4.2 <i>0.2</i>	5.1 <i>0.2</i>	4.3 <i>0.1</i>
	10	5.3 <i>0.2</i>	6.0 <i>1.2</i>	4.0 <i>0.8</i>	3.0 <i>0.9</i>	4.1 <i>0.4</i>	3.9 <i>0.3</i>	3.6 <i>0.4</i>	4.1 <i>0.3</i>	4.9 <i>0.2</i>	5.2 <i>0.1</i>	4.7 <i>0.1</i>
Total	3.1 <i>0.1</i>	2.2 <i>0.1</i>	2.9 <i>0.1</i>	2.9 <i>0.1</i>	3.0 <i>0.1</i>	3.4 <i>0.1</i>	3.6 <i>0.1</i>	3.7 <i>0.1</i>	4.1 <i>0.1</i>	4.8 <i>0.1</i>	3.4 <i>0.03</i>	

### 3.3 Characteristics of households with questionable income responses

Various household characteristics are associated with the incidence of income contamination. The source from which income is derived is particularly relevant. For example, if we include only transfers and earnings from employment in the income definition, the perverse relationship between durables ownership and income decile shown in Figure 2 and Table 1 largely disappears or becomes statistically insignificant. However, other categories of income (which we refer to collectively as self-employment income) are important in Albania, so their

exclusion from our analysis is not feasible. Figure 3 demonstrates this, through a regression of the household's share of total income from self-employment or enterprise, on its consumption rank. It shows greater dependence on self-employment in the lower parts of the consumption distribution.



**Figure 3** Nonparametric locally-weighted regression of the share of income from self-employment/enterprise on consumption rank

Table 2 gives an indication of the characteristics of households for which there is a large discrepancy between their positions in the consumption and income distributions. We estimate a (median) regression of the difference between the income and consumption ranks on the household characteristics listed in Table 2, as a descriptive device. A negative coefficient in this regression indicates a tendency for the corresponding characteristic to be associated with understatement of income relative to consumption expenditure. Major factors associated with understatement of income are: education of the household head and spouse, the absence of a spouse and, particularly, a high share of income from self-employment or enterprise. Large, multi-generation households are least associated with income contamination. There are no strong associations with age or rural location.

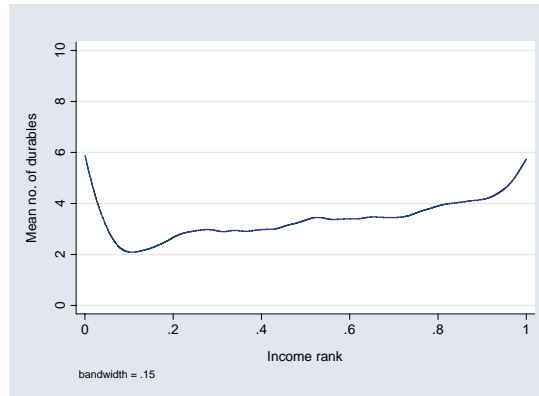
**Table 2** The relationship between the income-consumption discrepancy and household characteristics ( $n = 3,521$ )

Covariate	Coefficient	t-ratio
Age of household head / 10	-0.0085	2.35
No. of men in household	0.029	3.89
No. of women in household	0.026	3.91
Years of education of household head	-0.004	3.04
Years of education of spouse of household head	-0.006	3.69
No spouse of household head present	-0.084	4.73
Household has self-employment/enterprise income	0.031	2.18
Share of self-employment/enterprise income	-0.173	11.55
Missing self-employment data	-0.076	0.91
No. of generations in household	0.042	4.03
Household size	0.018	5.16
Farming household	-0.012	0.91
Intercept	0.019	0.57

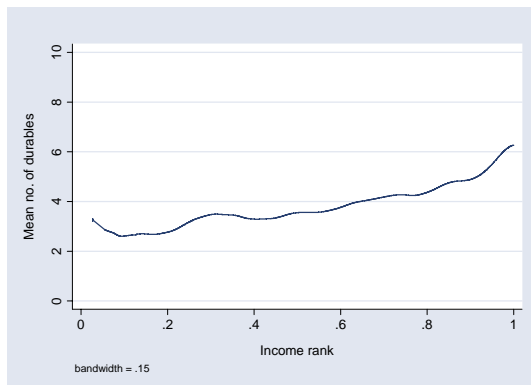
### 3.4 The persistence of income contamination through time

Panel data are particularly valuable for analysing poverty dynamics and the persistence of poverty through time. We have found significant evidence of income contamination at wave 1 of the ALSMS, to an extent that could make a substantial difference to conclusions from standard poverty research. The temporal nature of this contamination process is obviously important if we are interested in poverty dynamics as well as static poverty measurement. Our analysis of this issue is complicated by the fact that consumption expenditure is only observed in 2002. However, by repeating the nonparametric regressions of indicators of living standards on income rank we can confirm that the evidence for income contamination is strong in wave1, as Figure 3(a) shows for the durables-count indicator. The perverse pattern is less evident in waves 2 and 3. However, there remains enough evidence of income contamination in the bottom decile for there to be serious concern about the measurement of poverty from income data at each wave.

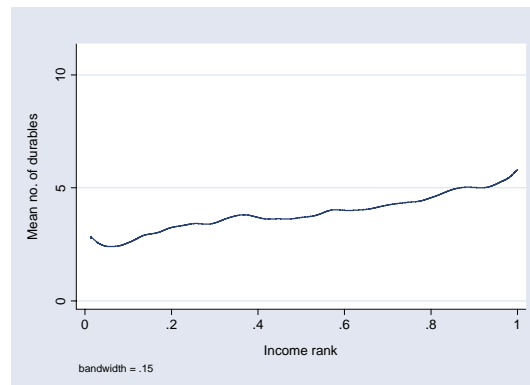




(a) Wave 1 (2002)



(b) Wave 2 (2003)



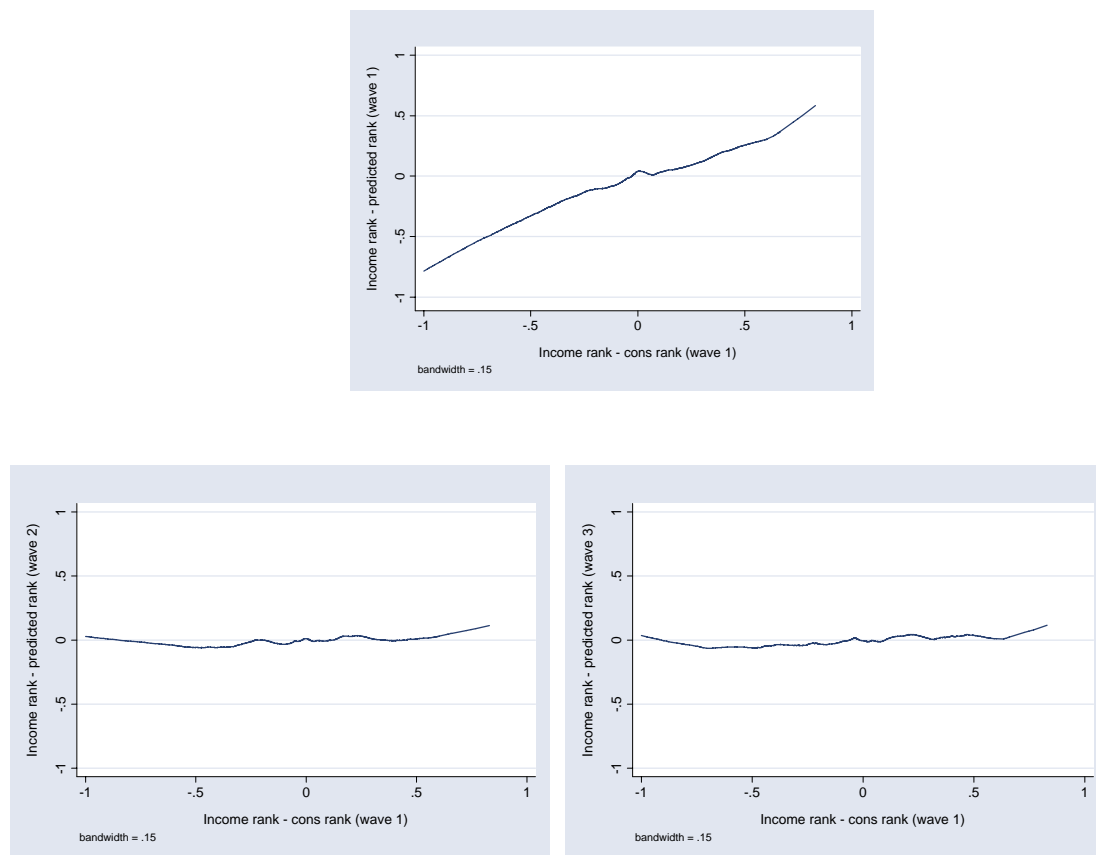
(c) Wave 3 (2004)

**Figure 3** Nonparametric locally-weighted regressions of durables ownership on income rank (panel sample, 2002-4)

We are interested in the persistence of contamination through time. Is it possible to say that, once a household reports ‘contaminated’ income, it tends to continue to do so in subsequent periods? We have two different indicators of possible contamination. The first (available only at wave 1) is the difference between the income and consumption rank. A second (available at all waves) is the residual from a wave-specific median regression model which predicts the income rank conditional on durables ownership and household characteristics. A negative value of either indicator suggests under-recording of income during the year in question. To investigate the persistence of income contamination through time, we estimate the relationship between the latter indicator at each wave and the value of the former indicator as observed at wave 1. We have already found a strong positive relationship between the two indicators in the full ALSMS and we would expect to see this also in wave 1 of the panel sample. The strength of the relationship at waves 2 and 3 of the panel sample will then give a picture of the degree of within-household persistence of income contamination through time.

We first summarise the relationship graphically in Figure 3, using nonparametric locally-weighted regression.

The results are striking. We see the anticipated strong positive relationship between the two contamination indicators at wave 1, so the concordance of the two indicators is confirmed in the panel sample. The patterns we observe at waves 2 and 3 are quite different, and show a positive but very slight relationship over most of the relevant range.<sup>2</sup>



**Figure 4** Nonparametric locally-weighted regressions of income rank at wave  $t$  on the difference between income and consumption ranks at wave 1 (panel sample)

<sup>2</sup> Note that the interval  $[-0.56, 0.46]$  covers 90% of the sample range of values for the difference between the income and consumption ranks, so the decreasing sections of the relationship for waves 2 and 3 over  $[-1, -0.5]$  is of no consequence.

The further analysis summarised in Table 4 allows for other household circumstances, using a random-effects regression for each of the four indicators of living standards, of the following form:

$$D_{it} = \alpha_t(Y_{i1} - C_{i1}) + \beta Y_{it} + \gamma X_{it} + U_i + V_{it} \quad (1)$$

where:  $D_{it}$  is the living standards indicator for household  $i$  in year  $t$ ;  $(Y_{i1}-C_{i1})$  is the difference between the income and consumption ranks at wave 1, which has a time-varying coefficient;  $Y_{it}$  is the current income rank;  $X_{it}$  is a set of other household characteristics and year dummies;  $U_i$  is an unobserved household effect; and  $V_{it}$  is a residual. These regressions are descriptive devices rather than formal models, summarising the relationship between indicators of living standards and household resources and circumstances. The significant positive coefficients for the year 2 and 3 dummies reflects the general improvement in the first and fourth indicators of living standards that took place between the 2002 and 2003 survey interviews. Note that the fourth indicator is based on responses to a question which explicitly asks for an assessment of well-being relative to the rest of the population and would therefore be expected to show no trend in the average.

The term  $\alpha_t(Y_{i1} - C_{i1})$  summarises the extent to which households with a large income-consumption discrepancy at wave 1 tend also to have high living standards relative to measured income in year  $t$ . If  $\alpha_t < 0$ , then a household with low measured income (relative to consumption) in year 1 tends also to have high living standards relative to measured income in year  $t$ . We find a strong, highly significant, negative value for  $\alpha_1$ , which confirms the pattern observed in Figure 3(a) for year 1. There are smaller negative estimates for  $\alpha_2$  and  $\alpha_3$ , which are not uniformly significant. This in turn confirms the pattern observed in Table 4. Thus, after controlling for other observable characteristics and unobserved household effects, there is only weak evidence of persistence of the income contamination observed at wave 1.

**Table 4** Random-effects regression estimates of living standards models  
(t-ratios in parentheses)

Covariate	Living standards indicator			
	Durables ownership	General satisfaction	Adequacy of food consumption	Perceived distributional position
$(Y_{it} - C_{it}) \times \text{year 1}$	-1.088 (7.63)	-0.435 (6.72)	-0.368 (7.10)	-1.012 (7.16)
$(Y_{it} - C_{it}) \times \text{year 2}$	-0.407 (2.92)	-0.073 (1.15)	-0.012 (0.24)	-0.377 (2.72)
$(Y_{it} - C_{it}) \times \text{year 3}$	-0.241 (1.73)	-0.048 (0.76)	-0.004 (0.08)	-0.175 (1.27)
Current income rank	1.242 (13.80)	0.806 (17.35)	0.532 (14.20)	1.829 (18.34)
Age of household head / 10	0.145 (4.99)	0.041 (3.48)	0.029 (3.17)	0.104 (3.97)
No. of men in household	0.132 (2.89)	-0.018 (0.90)	-0.007 (0.42)	0.009 (0.20)
No. of women in household	0.196 (4.08)	0.030 (1.46)	0.015 (0.91)	0.147 (3.19)
Years of education of household head	0.088 (9.34)	0.012 (3.22)	0.010 (3.49)	0.045 (5.39)
Years of education of spouse of household head	0.050 (4.68)	0.018 (4.11)	0.010 (2.92)	0.047 (4.78)
No spouse of household head present	0.203 (1.86)	-0.072 (1.52)	-0.014 (0.38)	-0.095 (0.90)
Household has self-employment/enterprise income	-0.129 (2.15)	0.129 (4.29)	0.139 (5.74)	0.256 (3.93)
Share of self-employment/enterprise income	0.199 (3.24)	0.036 (1.15)	0.007 (0.28)	-0.013 (0.20)
Missing self-employment data	-0.138 (1.53)	0.167 (3.54)	0.128 (3.34)	0.269 (2.67)
No. of generations in household	0.233 (4.12)	0.035 (1.36)	-0.003 (0.15)	0.136 (2.36)
Household size	-0.051 (1.65)	-0.007 (0.53)	0.009 (0.83)	-0.054 (1.84)
Farming household	-0.219 (3.60)	0.030 (1.00)	0.027 (1.11)	0.021 (0.32)
Wave 2	0.452 (11.52)	0.047 (2.04)	0.189 (10.03)	0.439 (9.25)
Wave 3	0.519 (13.17)	0.098 (4.30)	0.166 (8.80)	0.473 (9.96)
Intra-household correlation <sup>2</sup>	0.597	0.286	0.239	0.361

## **4 A nonparametric decontamination method for the income distribution**

We have established that there is a substantial group of non-deprived households for whom measured income greatly understates their economic resources and welfare. The effectiveness of consumption and deprivation indicators as a means of identifying cases with a high probability of income error suggests their use in some adjustment mechanism for the contaminated empirical income distribution. The most satisfactory way of designing such a mechanism is to derive it from an explicit model of the contamination process.

There are two principal classes of measurement error model that have been used in the applied literature on income distribution. One assumes that measurement error is continuously distributed and distorts measured income via a relationship of known functional form, typically additive or multiplicative. Chesher and Schluter (2002) work with this type of model and use small- $\sigma$  approximations to assess the bias in various inequality and poverty measures. Another approach uses a mixture model to generate distorted outliers. This is used by Cowell and Victoria-Feser (1996) in an analysis of the sensitivity of estimated inequality measures to outliers. Our findings for Albania appear more consistent with the latter view of contamination. However, we depart from the Cowell and Victoria-Feser (1996) analysis in four ways. Firstly, we use extraneous information on deprivation indicators to give additional information on the contamination process. Secondly, we do not make their assumption of mean-preserving contamination, which is clearly inappropriate in this case, since the dominant form of distortion appears to be occasional cases of very large understatement. Thirdly, we allow the contamination probability to vary with household characteristics, rather than being constant. Finally, we use a nonparametric approach to estimation rather than robust estimation of a parametric model.

### **4.1 A contamination model**

Let  $D$  be a well-being or inverse deprivation indicator and let  $X$  be a vector of household characteristics.  $\tilde{Y}$  is observed income, which may be contaminated by error in some cases. Define  $Y$  to be the true level of income and  $Y^*$  to be the income level that would be reported in the event of contamination. The error process is then assumed to be as follows:

$$\tilde{Y} = \begin{cases} Y & \text{if } Z = 0 \\ Y^* & \text{if } Z = 1 \end{cases} \quad (2)$$

where  $Z$  is a binary latent variable such that:

$$\Pr(Z = 1 \mid X, D, Y, Y^*) = \pi(X) \quad (3)$$

Note that assumption (3) asserts that the incidence of error is independent of both the true and contaminated levels of income, and of deprivation. It is a considerable generalisation of the Cowell and Victoria-Feser (1996) model but is nevertheless a strong assumption which we reconsider in section 4.4.

Let  $F_{y|x,d}(\tilde{Y} \mid X, D)$  be the conditional distribution function of measured income, which can, in principle, be estimated consistently from a sufficiently large set of sample data. Under assumption (3), this observed contaminated income distribution has the following form:

$$F_{y|x,d}(\tilde{Y} \mid X, D) = [1 - \pi(X)]G_{y|x,d}(\tilde{Y} \mid X, D) + \pi(X)H_{y|x}(\tilde{Y} \mid X) \quad (4)$$

where  $G_{y|x,d}(\cdot)$  and  $H_{y|x}(\cdot)$  are the conditional distribution functions of true and contaminated incomes respectively and where we have assumed that the contaminated income figure is independent of true living standards. We postpone to section 4.4 consideration of more general structures where the contamination process may depend on living standards,  $D$ , as well as household characteristics,  $X$ .

## 4.2 Exact identification

Our aim is to identify the conditional headcount poverty rate,  $G_{y|x}(T \mid X)$ , where  $T$  is a poverty threshold. This is clearly impossible without further restrictions. The key to identification is to use a combination of income and deprivation variables to identify cases with a very high probability of error, in other words, a region of  $(Y, X, D)$ -space where the density  $g_{y|x,d}(Y \mid X, D)$  is close to zero. We proceed in two stages.

First, assume that it is possible to specify a non-trivial level of well-being at which we can say *a priori* that the household cannot be classed as poor. Then there is a (high) level of well-being,  $d_0$ , such that  $\Pr(Y \leq T \mid X, D = d_0) = 0$  and  $\Pr(X, D = d_0) > 0$ . The contamination model (4) implies that  $F_{y|x,d}(T \mid X, d) - F_{y|x,d}(T \mid X, d_0)$  can be written  $[1 - \pi(X)](G_{y|x,d}(T \mid X, d) - G_{y|x,d}(T \mid X, d_0))$  and thus:

$$\begin{aligned} \sum_d P(D = d | X) F_{y|x,d}(T | X, d) - F_{y|x,d}(T | X, d_0) &= F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0) \\ &= [1 - \pi(X)] (G_{y|x}(T | X) - G_{y|x,d}(T | X, d_0)) \end{aligned}$$

Consequently:

$$G_{y|x}(T | X) = \frac{F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)}{1 - \pi(X)} + G_{y|x,d}(T | X, d_0) \quad (5)$$

The relation (5) identifies  $G_{y|x}(T | X)$  provided we can find a well-being level  $d_0$  such that  $G_{y|x,d}(T | X, d_0) = 0$  and provided that the contamination rate  $\pi(X)$  can itself be identified.

The second stage of the identification strategy gives information about  $1 - \pi(X)$ . For any  $X$ , choose a (middle-) income threshold  $y$  and a pair of well-being levels  $d_1, d_2$  where  $d_1 < d_2$ . Note that  $y, d_1$  and  $d_2$  may depend on the conditioning value  $X$ . Then (3) implies:

$$\begin{aligned} F_{y|x,d}(y | X, d_1) - F_{y|x,d}(y | X, d_2) &= [1 - \pi(X)] [G_{y|x,d}(y | X, d_1) - G_{y|x,d}(y | X, d_2)] \\ &\leq 1 - \pi(X) \end{aligned} \quad (6)$$

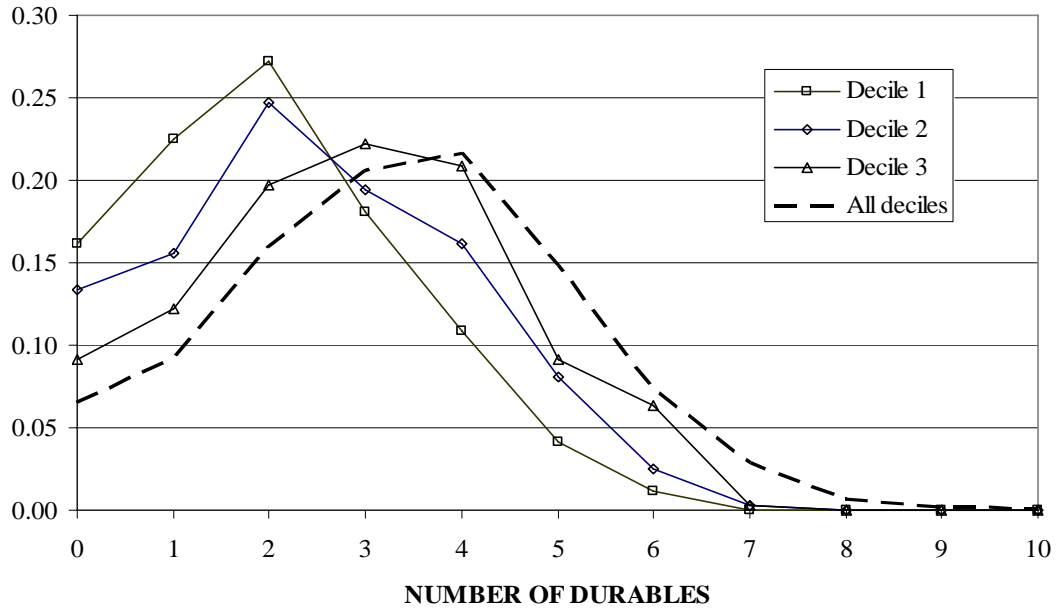
If, for any  $X$ , it is possible to find  $(y, d_0, d_1, d_2)$  which satisfy  $G_{y|x,d}(y | X, d_2) = G_{y|x}(T | X, d_0) = 0$  and  $G_{y|x}(y | X, d_1) = 1$ , then (6) is an equality and  $G_{y|x}(T | X)$  is exactly identified by the ratio:

$$G_{y|x}(T | X) = \frac{F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)}{F_{y|x,d}(y | X, d_1) - F_{y|x,d}(y | X, d_2)} \quad (7)$$

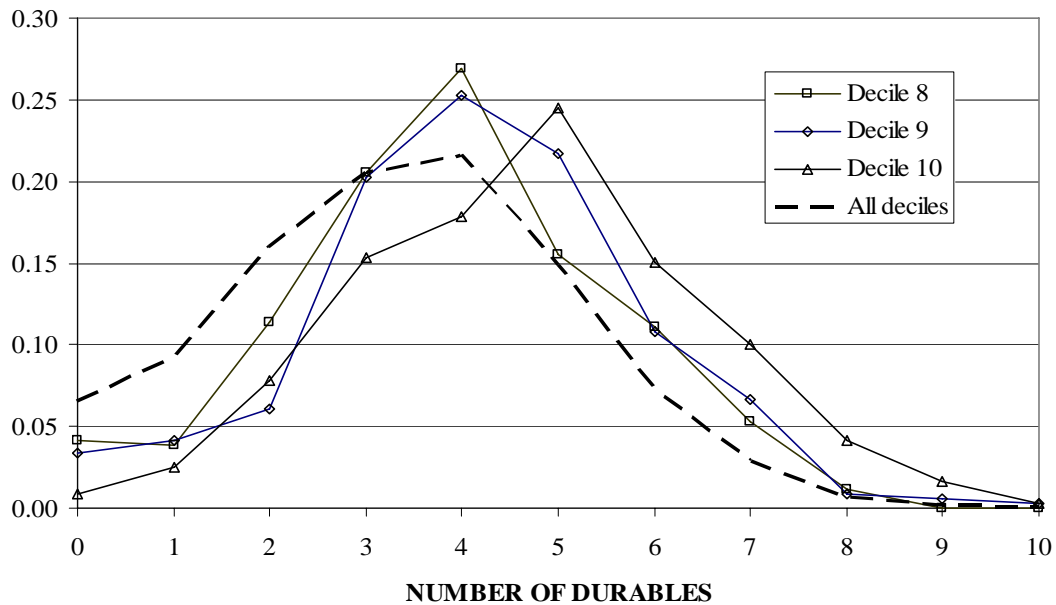
But is this a reasonable assumption? Since we have consumption expenditure data for one wave of the Albanian panel, it is possible to use the 2002 data to investigate the choice of  $(y, d_0, d_1, d_2)$ . Henceforth, we use the durables ownership variable as the well-being indicator, since it is consistent with the three subjective indicators but gives a finer categorisation of living standards. Figures 5 and 6 show the distributions of numbers of durables owned by households in the bottom three and top three deciles of the consumption distribution. There are two important conclusions. Firstly, among households in the bottom three consumption deciles, virtually none owns 7 or more of the listed durables. Therefore,  $d_0 \geq 7$  appears to be a good choice and  $d_0 \geq 6$  a possible alternative. Thus  $d_0$  refers to a composite category of ownership levels, rather than a single level.

Secondly, the situation is less clear in the upper part of the consumption distribution. Around 6% of the households in the top three deciles own no more than one of the listed durables and only in the top decile is the proportion as low as 3.5%. Consequently, with the

durables-based well-being index, it is not possible to find a triple  $(y, d_1, d_2)$  which will reliably identify  $1-\pi(X)$  exactly through (6).



**Figure 5** Number of durables owned in the bottom three deciles of the 2002 Albanian consumption distribution



**Figure 6** Number of durables owned in the top three deciles of the 2002 Albanian consumption distribution



### 4.3 Bounds for the poverty rate

Exact decontamination of the income distribution is not possible if the construction of the deprivation scale and its distribution in the population are such that the required values ( $y$ ,  $d_0$ ,  $d_1$ ,  $d_2$ ) do not exist. However, we can still calculate bounds on the true income poverty rate and these bounds will be useful as an indicator of the uncertainty associated with measured income poverty.

First note that  $F_{y|x,d}(T | X, d_0) = \pi(X)H_{y|x}(Y | X) \leq \pi(X)$ . Now consider a bound based on (5) and assume  $d_0$  is such that  $G_{y|x,d}(T | X, d_0) = 0$ . Choose a value  $\tilde{\pi}(x)$  to minimise the criterion  $\|G_{y|x}(T | X) - F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)/(1 - \pi(X))\|$  subject to  $F_{y|x,d}(T | X, d_0) = \pi(X)H_{y|x}(Y | X)$ , with  $\pi(X)$  and  $H_{y|x}(Y | X)$  both constrained to lie in the unit interval. The solution to this minimisation problem is  $\tilde{\pi}(x) = F_{y|x,d}(T | x, d_0)$ , giving the following sharpest possible lower bound on  $G_{y|x}(T | x)$ :

$$L(x) = \frac{F_{y|x}(T | x) - F_{y|x,d}(T | x, d_0)}{1 - F_{y|x,d}(T | x, d_0)} \quad (8)$$

Its sample analogue is:

$$\hat{L}(x) = \frac{\#(Y \leq T, X = x) / \#(X = x) - \#(Y \leq T, X = x, D = d_0) / \#(X = x, D = d_0)}{1 - \#(Y \leq T, X = x, D = d_0) / \#(X = x, D = d_0)} \quad (9)$$

where  $\#(A)$  denotes the sample frequency of any specified event  $A$ . A standard error for this estimated bound can be constructed using the usual first-order large-sample approximation. The bound  $L(x)$  will be binding if our specification of  $d_0$  is accurate and the contaminant distribution is contained below the poverty line in the sense that  $H_{y|x}(T | x) = 1$ . It will lie far below the unadjusted poverty rate  $F_{y|x}(T | x)$  if the proportion of high-welfare households with and low income,  $F_{y|x,d}(T | x, D = d_0)$ , is large. Note that  $\min\{\pi(X), H_{y|x}(Y | X)\} \geq F_{y|x,d}(T | X, d_0)$ , so a large value of the latter necessarily implies large values for the contamination rate and the poverty rate among contaminated observations.

Now consider the upper bound on  $G_{y|x,d}$ . Given the nature of the apparent bias induced by data contamination, it is reasonable to assume that the contaminant distribution is left-

shifted, in the sense that  $H_{y|x}(T | X) \geq G_{y|x,d}(T | X, D)$  for all  $X, D$ . Then the unadjusted poverty rate is itself an upper bound on the true poverty rate  $G_{y|x}(T | x)$ . The sample analogue of this bound is:

$$\hat{G}_{y|x}(T | X) = \frac{\#(Y \leq T, X = x)}{\#(X = x)} \quad (10)$$

An alternative upper bound is based on the following ratio, which can be estimated consistently from sample data:

$$U = \frac{F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)}{F_{y|x,d}(y | X, d_1) - F_{y|x,d}(y | X, d_2)} \quad (11)$$

where  $(y, d_1, d_2)$  are arbitrary values satisfying  $d_1 < d_2$ . Equation (3) implies:

$$U = \frac{G_{y|x}(T | X) - G_{y|x,d}(T | X, d_0)}{G_{y|x,d}(y | X, d_1) - G_{y|x,d}(y | X, d_2)} \quad (12)$$

$U$  is an upper bound on  $G_{y|x}(T | x)$  provided :

$$\frac{G_{y|x,d}(T | X, d_0)}{G_{y|x}(T | X)} \leq 1 - G_{y|x,d}(y | X, d_1) + G_{y|x,d}(y | X, d_2) \quad (13)$$

The right-hand side of (13) is the sum of the probabilities of having income above  $y$  given durables ownership of  $d_1$  and below  $y$  at ownership level  $d_2$ . It is necessarily satisfied if  $d_0$  is well-chosen so that  $G_{y|x,d}(T | x, d_0) = 0$ , but will also hold more generally if the poverty rate among high-ownership households is much lower than the overall poverty rate.

Provided (13) is satisfied for all feasible  $(y, d_1, d_2)$ , the upper bound can be made as tight as possible by choosing  $y, d_1, d_2$  to minimise  $U$  as defined by (11). In practice, the minimisation would need to be constrained to values for which there are adequate numbers of observations to estimate  $F_{y|x,d}$  with acceptable precision. Thus, we propose:

$$\hat{U}(x) = \min_{y, d_1, d_2} \frac{F_{y|x}(T | x) - F_{y|x,d}(T | x, D \geq d_0)}{F_{y|x,d}(y | x, D = d_1) - F_{y|x,d}(y | x, D = d_2)} \quad (14)$$

In the sample analogue, we would also include a constraint to control statistical precision:

$$\#(X = x, D \geq d_j) \geq N_{\min} \quad j = 0, 1, 2 \quad (15)$$

where  $N_{\min}$  is the minimum acceptable number of observations for estimating cell-specific poverty rates; we use  $N_{\min} = 50$  in the application. We also restrict  $d_1 \leq 4$  and  $d_2 \geq 6$  and  $y \in \{Q_3 \dots Q_7\}$ , where  $Q_j$  is the  $j$ th income decile point.

We calculate the lower bound on the poverty rate for  $d_0 = 6$  and 7, for a poverty line defined as 60% of the median of (unadjusted) income and two alternative equivalence scales. The results are summarised in Table 3. They are very striking. The bounds are narrow but lie far below the sample poverty rate, suggesting the existence of a very large upward bias in unadjusted poverty rates. The overall estimated poverty rate is reduced from around 30% to around 10%. In the urban sample it is reduced from 25-28% to 5-9%, and from 41-43% to 16-22% in the rural sample. The reason for this large size of the adjustment is that  $F_{y|x,d}(T | x, d_0)$  is large: almost 20% of households in the high ownership group report incomes below the poverty line. Therefore, this very large bias is not an artefact of the method – there is a large conflict between the standard of living indicators and the income data which is very evident in the data.

**Table 3** Bounds on the poverty rate  
(Poverty line = 60% of median measured equivalised income)

Equivalence scale	$d_0 = 6+$		$d_0 = 7+$		Sample poverty rate $\hat{U}$
	$\hat{L}$	$\hat{U}$	$\hat{L}$	$\hat{U}$	
<i>Whole sample</i>					
Per capita scale	.108 (.020)	.116 (-)	.105 (.034)	.115 (-)	.312 (.008)
Root household size	.105 (.019)	.107 (-)	.095 (.033)	.098 (-)	.287 (.008)
<i>Non-farm households</i>					
Per capita scale	.047 (.021)	.050 (-)	.044 (.035)	.049 (-)	.247 (.009)
Root household size	.048 (.020)	.049 (-)	.033 (.035)	.035 (-)	.228 (.008)
<i>Farming households</i>					
Per capita scale	.235 (.061)	.253 (-)	.223 (.125)	.242 (-)	.472 (.016)
Root household size	.235 (.057)	.242 (-)	.267 (.108)	.276 (-)	.433 (.015)

#### 4.4 Bounds under endogenous income misreporting

An objection to the previous analysis is that it assumes a uniform rate of misreporting across the distribution of well-being (conditional on  $X$ ) and may therefore overestimate  $\pi(X)$  for most of the income distribution. If, instead, it is principally those with high living standards

who tend to understate their incomes, then the bounds  $L(X)$  and  $U(X)$  will be excessively low.

Under these weaker conditions, (3) is replaced by :

$$F_{y|x,d}(\tilde{Y} | X, D) = [1 - \pi(X, D)]G_{y|x,d}(\tilde{Y} | X, D) + \pi(X, D)H_{y|x}(\tilde{Y} | X, D) \quad (16)$$

In this case, the differenced poverty rate is :

$$F_{y|x,d}(T | X, D) - F_{y|x,d}(T | X, d_0) = [1 - \pi(X, D)]G_{y|x,d}(T | X, D) + \left\{ \pi(X, D)H_{y|x}(T | X, D) - \pi(X, d_0)H_{y|x}(T | X, d_0) \right\} \quad (17)$$

This will tend to understate the true poverty rate  $G_{y|x}(T | X)$  excessively whenever the last bracket in (17) is negative: in other words, when the “misreporting as poor” rate is much greater for high-welfare people than for low-welfare people.

Consider first the lower bound  $L(x)$  defined by (8) and assume that misreporting is biased towards the highest durables ownership group, in the sense that  $\pi(X, D)$  satisfies:

$$\pi(X, D) = \begin{cases} \pi_- & D < d_0 \\ \pi_+ & D = d_0 \end{cases} \quad (18)$$

where  $\pi_- < \pi_+$ . In this case, after some manipulation,  $L(x)$  can be written:

$$L(x) = G(T | x) + \frac{\pi_- (H_{y|x}(T | x) - G_{y|x}(T | x)) + \pi_+ H_{y|x,d}(T | x, d_0) G_{y|x}(T | x)}{(1 - \pi_+) H_{y|x,d}(T | x, d_0)} - \frac{H_{y|x,d}(T | x, d_0) (\pi_+ (1 - G_{d|x}(d_0 | x)) + \pi_- G_{d|x}(d_0 | x))}{(1 - \pi_+) H_{y|x,d}(T | x, d_0)} \quad (19)$$

This is no longer necessarily a lower bound, since the last two ratio terms in (19) are positive and negative respectively. If the conditional probability of  $D = d_0$  is sufficiently large and  $\pi_-$  sufficiently low,  $L(x)$  will lie above the true poverty rate  $G_{y|x}(T | X)$ .

For an alternative approach, consider the sample poverty rate:.

$$F_{y|x}(T | X) = \sum_{D=0}^{d_0-1} G_{d|x}(D | X) F_{y|x,d}(T | X, D) + G_{d|x}(d_0 | X) F_{y|x,d}(T | X, d_0) \\ = \sum_{D=0}^{d_0-1} G_{d|x}(D | X) \left\{ (1 - \pi_-) G_{y|x,d}(T | X, D) + \pi_- H_{y|x,d}(T | X, D) \right\} + \pi_+ G_{d|x}(d_0 | X) H_{y|x,d}(T | X, d_0) \quad (20)$$

Note that the last term of (20) is equal to  $F_{y|x,d}(T | X, D)$ , and we can write:

$$F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0) = G_{y|x}(T | X) + \pi_- \left\{ H_{y|x}(T | X) - G_{d|x}(d_0 | X) H_{y|x,d}(T | X, d_0) - G_{y|x}(T | X) \right\} \quad (21)$$

The term in braces in (21) is the difference between the poverty rate in misreported income (with the component from the top ownership level removed) and the overall true poverty rate. It is multiplied by the misreporting rate,  $\pi$ , for lower-ownership households and is expected to be positive but of moderate size. Thus  $F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)$  is expected to be an upper bound. It is binding if, outside the top ownership category, there is no misreporting or if it has no impact on the poverty rate.

The natural sample analogue of  $F_{y|x}(T | X) - F_{y|x,d}(T | X, d_0)$  is:

$$\hat{U}^*(x) = \frac{\#(Y \leq T, X = x, D < d_0^*)}{\#(X = x)} \quad (22)$$

This estimate is given in Table 4. It lies only slightly below the crude sample poverty rate and thus leaves a wide margin of uncertainty within which the true poverty rate lies.

**Table 4** Bounds on the poverty rate  
(Poverty line = 60% of median measured equivalised income)

Equivalence scale	Upper bound $\hat{U}^*$		Sample poverty rate $\hat{U}$
	$d_0 = 6+$	$d_0 = 7+$	
<i>All households</i>			
Per capita scale	.289 (.008)	.304 (.008)	.312 (.008)
Root household size	.266 (.007)	.280 (.007)	.287 (.008)
<i>X = Urban households</i>			
Per capita scale	.219 (.008)	.237 (.008)	.247 (.009)
Root household size	.203 (.008)	.219 (.008)	.228 (.008)
<i>X = Rural households</i>			
Per capita scale	.461 (.016)	.470 (.016)	.472 (.016)
Root household size	.423 (.015)	.431 (.015)	.433 (.015)

How do we interpret these findings? Firstly, there is a large group of income observations which are clearly misleading indicators of economic welfare in the sense that they are contradicted by data from the same households on consumption expenditure, durables ownership and subjective assessments of well-being. Specification of an explicit contamination model as a description of the income distribution, together with a plausible  $a$

*priori* assumption about the impossibility of true poverty for high-consumption households, allows us to put bounds on the true poverty rate. If we feel able to assume that the contamination process is, at least approximately, independent of the true standard of living, then there is evidence of a very strong upward bias in a conventional estimate of the poverty rate. On the other hand, if we make the extreme assumption that the contamination process is essentially confined to the very top of the distribution of living standards, then the bias is very much smaller. The main finding from this is the extremely high degree of uncertainty associated with any measure of poverty constructed from these income data.

## **5 Conclusions**

This study has had three aims. One is to use survey data on income, consumption and well-being from Albania in 2002 to assess the scale and nature of income misreporting. This has revealed strong evidence that the bottom tail of the income distribution is contaminated by misleading income figures for some relatively high-welfare households. Analysis of a subset of these households, for whom a sequence of income measures is available for two further years, suggests that this perverse pattern is less evident in 2003-4, but there remains sufficient evidence of income contamination in the bottom decile to give cause for concern about the accuracy of conventional poverty analyses.

Secondly, we have developed the statistical theory of contaminated distributions to incorporate information on an ancillary variable which is informative for the income misreporting process and proposed a nonparametric method of analysis based on an *a priori* specification of a level of well-being below which poverty is seen as impossible. The choice of this level was motivated by our analysis of conflict between Albanian consumption and income data.

Our third aim was to quantify the degree of uncertainty associated with estimation of poverty rates in the Albanian case by using the method to place bounds on the true poverty rate for 2002. These bounds turn out to be very sensitive to the assumed relationship between the propensity to misreport and standard of living. If we assume that misreporting is independent of living standards, the bounds are narrow and indicate very large upward biases in standard survey-based income poverty rates. However, if misreporting is largely confined to high-welfare households, the bounds are very wide. In the absence of more direct

information on the incidence of misleading income reports, this leaves us with the general conclusion that income poverty measures are subject to a great deal of uncertainty.

If income variables are to be used for the study of poverty, our analysis emphasises the importance of also using ancillary variables to check whether a significant number of apparently low-income households are misclassified because of misreported or otherwise misleading income data. Well-being indicators provide a very fruitful way of making this check but there is a need for the design of these indicators to be as informative as possible. A count of durable ownership is reasonably effective in a country as poor as Albania, which still has low levels of durables ownership. However, such variables would be less useful in developed countries with high rates of ownership throughout the welfare distribution. In these countries, subjective assessments of well-being may work better, but they need to be carefully designed to give good discrimination, by showing adequate variation across the distribution of true economic welfare.

## References

- Aassve, A., Engelhardt, H., Francavilla, F., Kedir, A., Kim, J., Mealli, F., Mencarini, L., Pudney, S. E., Fuernkranz-Pskawetz, A. (2005). Poverty and Fertility in Less Developed Countries: a comparative analysis. *Working Papers of the Institute for Social and Economic Research*, paper 2005-13. Colchester: University of Essex.
- Azzarri, C., Carletto, G., Davis, B. and Zezza, A. (2006). Monitoring poverty without consumption data, *Eastern European Economics* **44**, 59-82.
- Berthoud, R., Bryan, M. and Bardasi, E. (2004). *The dynamics of deprivation: the relationship between income and material deprivation over time*. London: Department for Work and Pensions, Research report no. 219.
- Carletto, G. and Zezza, A. (2004). Being poor, feeling poorer: combining objective and subjective measures of welfare in Albania. Working paper, FAO.
- Chesher, A. and Schluter, C. (2002). Welfare measurement and measurement error, *Review of Economic Studies* **69**, 357-378.
- Cowell, F. A. and Victoria-Feser, M.-P. (1996). Robustness properties of inequality measures, *Econometrica* **64**, 77-101.
- Grosh, M. and Glewwe, P. (1995). *A Guide to Living Standards Surveys and Their Data Sets*. Washington: The World Bank, LSMS Working Paper #120.

- Grosh, M.E. and Glewwe, P. (eds.) (2000): *Designing Household Survey Questionnaires for Developing Countries: Lessons from 15 Years of the Living Standards Measurement Study. Vols 1, 2 and 3*, Washington: World Bank.
- Huber, P. J. (1964). Robust estimation of a location parameter, *Annals of Mathematical Statistics* **35**, 73-101.
- McKay, S. (2004). Poverty or preference: what do 'consensual deprivation indicators' really measure?, *Fiscal Studies* **25**, 201-224.
- Ravallion, M. and Lokshin, M. (2001). Identifying welfare effects from subjective questions, *Economica* **68**, 335-357.
- Sen, A. (1985). *Commodities and Capabilities*, Amsterdam: North Holland.
- Ringen, S. (1988). Direct and indirect measures of poverty, *Journal of Social Policy* **17**, 351-365.
- Saunders, P. (2005). Researching social policy: trends, tragedies and triumphs, plenary address to the 9<sup>th</sup> Australian Social Policy Conference, University of New South Wales, July 2005.



## **Appendix: Indicators of well-being in the ALMS**

*In the paper four well-being indicators are used. The first is the number of durables owned from a list of 11 possibilities. The indicator has been created using the following question in the ALMS present in all the three waves:*

*“How many of the following items does your household own?”*

*The question in the survey is repeated for 24 items. For the creation of our first well-being indicator we generated a variable equal to the number of each type of durables owned by the household focusing on the following 11 items: Colour TV, Video Player, CD Player, Video Camera, Refrigerator, Freezer, Washing Machine, Dishwasher, Computer, Satellite Dish, Car. The indicator ranges from 0 for households who own none of the listed durables to 11 for households owning at least one of each type.*

The second indicator from responses to the following question on household satisfaction with current economic circumstances:

*“How satisfied are you with your current situation?”*

Responses are recoded as (4) “fully satisfied”; (3) “rather satisfied”; (2) “less than satisfied”; “not at all satisfied” (1). The “don’t know” and “refuse to answer” categories are treated as missing responses and excluded from the analysis.

The third indicator is based on responses to the following question on the adequacy of current food consumption:

*“Would you consider the current level of food consumption of your family as ...”*

The selected response is recoded as: (3) “more than adequate”; (2) “just adequate”; (1) “less than adequate”. “Don’t know” and “refuse to answer” are again treated as missing responses.

The fourth indicator reflects the respondent’s own assessment of the household’s current position in the distribution of living standards. It is coded as a direct response to the following question:

*“Imagine a 10-step ladder where on the bottom, the first step, stand the poorest people, and on the highest step, the TENTH, stand the rich. On which step are you today*