

## IS INEQUALITY UNDERESTIMATED IN EGYPT? EVIDENCE FROM HOUSE PRICES

BY ROY VAN DER WEIDE\*, CHRISTOPH LAKNER AND ELENA IANCHOVICHINA

*The World Bank*

### Abstract

Household income surveys often fail to capture top incomes, which leads to an underestimation of income inequality. A popular solution is to combine the household survey with data from income tax records, which has been found to result in significant upward corrections of inequality estimates. Unfortunately, tax records are unavailable in many countries, including most of the developing world. In the absence of data from tax records, this study explores the feasibility of using data on house prices to estimate the top tail of the income distribution. In an application to Egypt, where estimates of inequality based on household surveys alone are low by international standards, the study finds strong evidence that inequality is indeed being underestimated by a considerable margin. The Gini index of household per capita income for urban Egypt is found to increase from 39 to 52 after correcting for the missing top tail.

**JEL Codes:** D31

**Keywords:** inequality, Egypt, house prices, top incomes

### 1. INTRODUCTION

Estimates of income inequality are conventionally derived from household income and expenditure surveys. Due to the sizeable cost of collecting accurate data on household standards of living, the sample size of these surveys generally constitutes less than half a percent of the total population. Unfortunately, the rich are often missing or under-covered, either due to non-response or under-reporting of income or both; see the recent literature on top income shares (e.g. Atkinson *et al.*, 2011). Surveys still permit accurate estimation of median income and measures of poverty, even when data on top incomes are poor or are missing altogether, since the rich make up a small percentage of the total population. For the estimation of income inequality however, having good data on top incomes is crucial.

*Note:* The authors wish to thank Guoliang Feng and Youssouf Kiendrebeogo for excellent research assistance, and Gabriel Lara Ibarra for help with the household survey data. We would like to thank Facundo Alvaredo, Francisco Ferreira, Alan Fuchs, Nadine Ghobrial, Vladimir Hlasny, Aart Kraay, Peter Lanjouw, Branko Milanovic, Thomas Piketty, Martin Ravallion, Paolo Verme, an anonymous referee, and participants of the World Bank workshop on the Arab Inequality Puzzle and the IARIW-CAPMAS Conference “Experiences and Challenges in Measuring Income, Wealth, Poverty and Inequality in the Middle East and North Africa” for useful comments. We are grateful to the U.K. Department for International Development for financial assistance through its Strategic Research Program.

\*Correspondence to: Roy van der Weide, 1818 H Street NW, Washington, DC 20433, USA (rvanderweide@worldbank.org).

A remedy that has gained considerable traction recently is to estimate the top tail of the income distribution using data from income tax records. This estimate of the top tail can then be combined with an estimate of the bottom part from the household survey to obtain an estimate of the complete income distribution (Atkinson, 2007; Alvaredo, 2011; Alvaredo Londoño Vélez, 2013; Diaz-Bazan, 2014; Anand and Segal, 2015).<sup>1</sup> Income tax records denote the ideal source of data as far as top incomes are concerned. For lower incomes, tax records may be less reliable; here, the household income survey arguably denotes the ideal data. When household survey and tax data are combined in this way, the Gini index for (i) the United States (U.S.) in 2006 increases from 59 to 62 (Alvaredo, 2011), (ii) Colombia in 2010 from 55 to 59 (Alvaredo *et al.*, 2013), (iii) Korea in 2010 from 31 to 37 (Kim and Kim, 2016), (iv) Ecuador in 2011 from 44 to 49 (Cano, 2015), and (v) Chile in 2013 from 52 to 59 (World Bank, 2016).

For all the pros of income tax records, the availability of the data is unfortunately rather limited, particularly in developing and emerging economies. The World Wealth and Income Database (Alvaredo *et al.*, 2017), for example, includes no countries from the Middle East and North Africa region, while we were able to find one study of top income shares in the region using fiscal data.<sup>2</sup> Furthermore, data derived from tax records are less useful in places where tax evasion is more pervasive, as is the case in many developing countries. It should also be noted that combining household survey data and tax records is not without complications, because the two data sources use different income definitions (disposable versus taxable) and have different units of analysis (households versus tax units, which could be individuals).

In the absence of data from tax records, this study explores the feasibility of using data on house prices to estimate the top tail of the income distribution. Market house price data can often be obtained more easily and, most importantly, tend to be available in the public domain, in contrast to tax administration data, which are subject to important confidentiality concerns and require cooperation from governments. Also, house sellers have no incentive to understate the value of their homes, in contrast to the income that taxpayers report on their tax returns.

Using house prices as an alternative to income tax records demands two methodological innovations to the study of top incomes. First, we will not be observing actual household income (as is the case with tax record data) but, rather, a predictor of income. Second, a database with house price listings is generally not obtained using a particular sampling design. Therefore, the data are not guaranteed to provide a nationally representative sample, they will arguably be biased toward large urban centers. We will propose workable solutions to both of these challenges that will hopefully contribute to a wider use of this approach.

<sup>1</sup>Diaz-Bazan (2014) generalizes the method of Atkinson (2007) and Alvaredo (2011) by allowing for a more general choice of the cutoff level for joining up the distributions. Morelli *et al.* (2015) review the literature attempting to combine household surveys and tax data in rich countries.

<sup>2</sup>Assouad (2015) estimates top income shares in Lebanon using individual tax records. She finds a high level of inequality, with a top 1 percent income share of 13 percent in 2012 (compared with 8 percent in Spain or 19 percent in the U.S.). She also highlights a number of concerns over the reliability of the national accounts data and tax evasion.

Note that the methodology is not restricted to the use of house prices; it can be applied to any database containing predictors of top incomes.

We illustrate our approach with an empirical application to Egypt, which provides a good testing ground for our method. In addition to being a major Arab country, inequality in Egypt is of considerable interest not least because it has been cited as one of the factors behind the Egyptian revolution (Hlasny and Verme, 2016). Estimates of inequality based on household surveys suggest that inequality is low in Egypt and that it has declined in the past decade to an expenditure-based Gini of around 31 in 2009.<sup>3</sup> Using house prices to capture top incomes, we find that inequality may be significantly underestimated in Egypt. The Gini of household per capita income for urban Egypt in 2009 is estimated at 51.8, compared to a survey-only estimate of 38.5. Our results are in contrast with other studies using different methods of adjusting for top incomes in Egypt (Hlasny and Verme, 2016), which report a more modest effect.<sup>4</sup> Their correction, however, does not consult a second source of data. If the main problem is that high-income earners are simply missing from the survey, then no adjustment that relies solely on the survey will resolve the downward bias in estimates of inequality. The only way to obtain a meaningful correction is to bring in a second source of data that carries the necessary information on top incomes and hence will permit for the consistent estimation of income inequality. This reasoning is shared by Alvaredo and Piketty (2014), who similarly argue that the household survey data by themselves are insufficient to estimate top incomes in Egypt. While they make an appeal for making data on income tax records available, we propose to work with house price data instead. It should be noted that relying on predictors of top incomes rather than actual incomes derived from tax records is not without caveats. For example, we need to make assumptions about the functional form of the relationship between the house price and household income, and about the functional form of the upper tail of the house price distribution. In addition, it is assumed that one house constitutes one household and that all houses are domestically owned. Therefore, in cases where tax record data are available, these should undoubtedly be considered first. However, we certainly believe that our approach provides more reliable estimates of inequality than estimates obtained using survey data alone. The perfect should not be the enemy of the good.

This paper is related to a number of other studies which have tried to correct household surveys for the problem of missing or underreported top incomes.<sup>5</sup> Korinek *et al.* (2006) exploit geographic variation in response rates to correct for selective non-response in the U.S. Lakner and Milanovic (2016) exploit the gap

<sup>3</sup>Source: PovcalNet, accessed October 31, 2015.

<sup>4</sup>The Gini coefficient of household expenditure per capita in 2009 increases from 30.7 to 31.8 in the preferred specification, which is found to be statistically significant at the 5 percent level.

<sup>5</sup>Recently, the EU-SILC survey in some countries began using register-based information (including tax records) for some questions (Jäntti *et al.*, 2013). This is, of course, preferable to any ex post combination of these different data sources, as we use in this paper. In the year after the introduction of the register data, the Gini index for France increased from 39 to 44, which is consistent with the previously used household data underestimating top incomes (Burricaid, 2013).

between household surveys and national accounts to adjust the top end of the income distribution.<sup>6</sup>

This paper is structured as follows. The methodology is presented in Section 2. In Section 3, we introduce the data used in the empirical application to Egypt. The empirical application itself is presented in Section 4, and Section 5 concludes. Finally, the Annex (in the Online Supporting Information) presents the results from a small validation of our methodology in a controlled setting.

## 2. METHODOLOGY

### 2.1. Combining Income Survey with Top Income Data

The objective is to estimate the level of income inequality for a given population. We will refer to database 1 (DB-1) as the primary data source for the estimation of inequality. It is assumed that top incomes are mostly missing from this database. Database 2 (DB-2), which we will refer to as the secondary data source, primarily contains data on top incomes but generally not on lower incomes. Estimates of income inequality will be biased if computed using any single one of these databases. It takes a combination of the two to obtain consistent estimates of inequality. DB-1 commonly represents a household income or expenditure survey. For DB-2 researchers often look at tax record data, as is discussed in the introduction.

Let us denote household income by  $y$  and its cumulative distribution function by  $F(y)$ . Let  $\tau$  denote the income threshold above which we will refer to incomes as “top incomes,” and let  $\lambda$  measure the share of the population enjoying a top income; that is,  $\lambda = \Pr[Y > \tau] = 1 - F(\tau)$ . It is assumed that DB-1 permits a consistent estimator for  $F_1(y) = \Pr[Y \leq y | Y \leq \tau]$ , and that DB-2 permits a consistent estimator for  $F_2(y) = \Pr[Y \leq y | Y > \tau]$ . By the same token, it is assumed that DB-1 does not permit a consistent estimator for  $F_2(y)$ , while DB-2 does not permit a consistent estimator for  $F_1(y)$ . Suppose also that an estimate of  $\lambda$  is available.<sup>7</sup> Given estimates of  $F_1(y)$ ,  $F_2(y)$ , and  $\lambda$ , an estimator for the complete income distribution function  $F(y)$  can be obtained as follows:

$$(1) \quad F(y) = \begin{cases} (1-\lambda)F_1(y), & y \leq \tau, \\ (1-\lambda) + \lambda F_2(y), & y > \tau. \end{cases}$$

Given  $F(y)$ , any measure of income inequality can readily be computed. Alternatively, one may appeal to the subgroup decomposition of one’s inequality measure of choice, which would by-pass the need for evaluating the income distribution for the population ( $F(y)$ ).<sup>8</sup> We have two subgroups, those with income below  $\tau$

<sup>6</sup>See also the study on global interpersonal inequality by Anand and Segal (2015), who append for every country the estimated top 1 percent share to the household survey distribution. The latter is assumed to represent the bottom 99 percent. For the majority of countries, the top 1 percent share is predicted from a cross-country regression using the top 10 percent share in the household survey.

<sup>7</sup>It is generally assumed that DB-2 contains the total number of units (i.e. households or tax units) whose income is above  $\tau$ . Combined with the total population, this yields an estimator for  $\lambda$ .

<sup>8</sup>Jenkins (2017) provides a categorization of different approaches to addressing top income under-reporting in survey data. Our approach of combining inequality indices derived from DB-1 and DB-2 falls under Category C.

(subgroup 1) and those with income above  $\tau$  (subgroup 2). Let  $P_k$  denote the population share of subgroup  $k$ , and let  $S_k$  denote their corresponding income shares; that is,  $S_k = P_k \mu_k / \mu$ , where  $\mu_k$  and  $\mu$  measure average income in subgroup  $k$  and the total population, respectively. Note that  $P_1 = 1 - \lambda$  and  $P_2 = \lambda$ . Let us also define  $S_1 = 1 - s$  and by extension  $S_2 = s$ . It can be verified that income inequality as measured by the Gini coefficient satisfies the following decomposition (see, e.g., Alvaredo, 2011):

$$(2) \quad Gini = P_1 S_1 Gini_1 + P_2 S_2 Gini_2 + S_2 - P_2$$

$$(3) \quad = (1 - \lambda)(1 - s) Gini_1 + \lambda s Gini_2 + s - \lambda,$$

where  $Gini_k$  measures the Gini coefficient for population subgroup  $k$ . A similar decomposition can be obtained for the mean-log-deviation  $MLD$  (see, e.g., Shorrocks, 1980):

$$(4) \quad MLD = P_1 MLD_1 + P_2 MLD_2 + P_1 \log \left( \frac{P_1}{S_1} \right) + P_2 \log \left( \frac{P_2}{S_2} \right)$$

$$(5) \quad = (1 - \lambda) MLD_1 + \lambda MLD_2 + (1 - \lambda) \log \left( \frac{\mu}{\mu_1} \right) + \lambda \log \left( \frac{\mu}{\mu_2} \right)$$

$$(6) \quad = (1 - \lambda) MLD_1 + \lambda MLD_2 + \log(\mu) - \log(\mu_1^{1-\lambda} \mu_2^\lambda),$$

and for the Theil index  $T$  (see, e.g., Shorrocks, 1980):

$$(7) \quad T = S_1 T_1 + S_2 T_2 + S_1 \log \left( \frac{S_1}{P_1} \right) + S_2 \log \left( \frac{S_2}{P_2} \right)$$

$$(8) \quad = (1 - s) T_1 + s T_2 + (1 - s) \log \left( \frac{\mu_1}{\mu} \right) + s \log \left( \frac{\mu_2}{\mu} \right)$$

$$(9) \quad = (1 - s) T_1 + s T_2 + \log(\mu_1^{1-s} \mu_2^s) - \log(\mu),$$

where  $MLD_k$  and  $T_k$  measure the mean-log-deviation and Theil index for population subgroup  $k$ , respectively. Note that the between-group inequality components of both the mean-log-deviation and the Theil index equal the difference between the arithmetic- and the geometric-mean income levels. They differ only in the weights used in the geometric mean; the mean-log-deviation weighs the subgroup means by their population shares, whereas the Theil index weighs them by their incomes shares.

An inspection of the three subgroup decompositions tells us that the Theil index will be most sensitive to the top tail of the income distribution.<sup>9</sup> To illustrate the significance of the top tail to total inequality, consider the limit where the population share of top income earners tends to zero ( $\lambda \rightarrow 0$ ) while their income share tends to some positive value ( $s > 0$ ). It can readily be seen that the

<sup>9</sup>Hence it is expected that any efforts made to fix the top tail of the income distribution by bringing in complementary data (top income database) will be rewarded the most by the Theil index.

between-group inequality component of the Gini coefficient tends to  $s > 0$  in that case, while the within-group inequality among top income earners tends to zero; that is,  $G \rightarrow (1-s)Gini_1 + s$ . It follows that the between-group inequality component for the mean-log-deviation tends to  $\log(1-s)^{-1}$ , while also here (as with the Gini) the within-group inequality among top earners tends to zero (yet it does not discount the contribution of within-group inequality among non-top earners); that is,  $MLD \rightarrow MLD_1 - \log(1-s)$ . The Theil index stands out as the only of the three inequality measures for which the within-group inequality among top earners does not vanish (i.e. it makes a positive contribution to total inequality) while the between-group inequality component tends to infinity (when  $\mu_2$  tends to infinity as  $\lambda \rightarrow 0$  while  $s > 0$ ).

## 2.2. An Alternative to Top Income Data: Challenges

As stated in the previous section, DB-2 (the top income database) typically takes the form of tax record data. These data have at least two advantages: (1) they directly observe realized incomes (which makes the estimation of  $F_2(y)$  or any income statistics such as inequality among top earners rather straightforward); and (2) they provide a count of the number of top income earners, which makes for a straightforward estimation of  $\lambda$ . A key disadvantage of tax record data is that it is often difficult to obtain access to them. Moreover, they are more likely to be available in developed countries with good quality data systems in place, and less likely to be available in developing countries.

This paper explores the feasibility of using an alternative to tax record data that is more readily available. The empirical application presented in Section 4 considers data on house prices compiled from publicly available real estate property listings as the alternative.<sup>10</sup> The advantage of these data is that their availability extends to developing countries. The flip side is that they also introduce a number of key methodological challenges due to the fact that the alternative database (a) observes predictors of income, not actual incomes, and (b) need not constitute a proper sample, so that it is unclear what population is being represented by the data.

The following two subsections aim to provide workable solutions to these two challenges that will hopefully contribute to a wider application of this approach.

### 2.2.1. A Database of Predictors of Top Incomes

Let us first focus on the challenge posed by observing a predictor of household income rather than actual income. Consider the following assumption.

**Assumption 1.** *Suppose that household income per capita can be described by*

$$(10) \quad \log(Y_h) = m(x_h; \beta) + \varepsilon_h$$

$$(11) \quad = \beta_0 + \beta_1 \log(x_h) + \varepsilon_h,$$

<sup>10</sup>Alternatively, one could, for example, also look to data on mortgages or credit card statements. However, this approach may not be feasible in countries with underdeveloped or non-existent mortgage markets.

where  $x_h$  denotes the predictor of household income per capita,  $\varepsilon_h$  denotes a zero expectation error term, subscript  $h$  indicates the household, and where  $\beta$  denotes a vector of model parameters.

The assumption of a log-linear model is motivated by ease of exposition and by the fact that it fits our empirical data remarkably well. This assumption can be relaxed, however, by accommodating flexible functional forms for  $m(x_h; \beta)$  if the data call for it. In our application, the value of the household's house (or rental value) will serve as the predictor  $x_h$ .

To obtain some intuition for the implications of Assumption 1 it may be helpful to verify what it implies for the relationship between the expenditure share on housing and household income. Note, however, that we are concerned with predicting household income per capita rather than household income. Let us abstract away from this distinction, for this thought experiment only, by considering the household size fixed (so that it is absorbed in the constant  $\beta_0$ ). It can be verified that the assumed functional form implies that the expenditure share on housing is a convex declining function of income when  $\beta_1 > 1$ . The expenditure share is constant for  $\beta_1 = 1$  and an increasing function of income for  $0 < \beta_1 < 1$ . More specifically, it is a concave increasing function for  $\beta_1 \in (\frac{1}{2}, 1)$ , a linear increasing function for  $\beta_1 = \frac{1}{2}$ , and a convex increasing function for  $\beta_1 \in (0, \frac{1}{2})$ . Despite its simplicity, the log-linear assumption permits a reasonable degree of flexibility in how the expenditure share on housing varies with income. Our prior would be that  $\beta_1 > 1$ , which is consistent with the empirical evidence that is available for the Engel curve on housing expenditure (see, e.g., Larsen, 2014).  $\beta_1 > 1$  also ensures that the expenditure share stays below 1 when incomes tend to extreme values.

Let  $F_\varepsilon(e; \sigma)$  denote the distribution function of  $\varepsilon_h$  with unknown parameter vector  $\sigma$ . We will assume that  $\varepsilon_h$  is identically distributed across households, although this assumption can easily be relaxed. Note that the unknown parameter vectors  $\beta$  and  $\sigma$  both have to be estimated. In our empirical application, where the value of housing is considered as a predictor of income, the two can be estimated using the household income survey, since it includes both data on household incomes and data on the value of housing.

It will be convenient to define  $n(\tau, y)$  as the number of households with income between  $\tau$  and  $y$ ,  $n(\tau)$  as the number of households with income exceeding  $\tau$ , and  $n$  as the total number of households in the population. For ease of exposition, we will ignore the fact that the data may constitute a sample with sampling weights.  $F_2(y)$  ( $=\Pr[Y \leq y | Y > \tau]$ ) and  $\lambda$  ( $=\Pr[Y > \tau]$ ) are seen to solve

$$(12) \quad F_2(y) = \frac{n(\tau, y)}{n(\tau)},$$

$$(13) \quad \lambda = \frac{n(\tau)}{n}.$$

When DB-2 does not contain data on household incomes but data on a predictor of household incomes instead, we have that  $n(\tau, y)$  and  $n(\tau)$  can no longer be observed with certainty and so they have to be estimated. Consider first an estimator for  $n(\tau)$ :

$$\begin{aligned} \hat{n}(\tau) &= \sum_h E[1(Y_h > \tau)|x_h] \\ &= \sum_h E[1(m(x_h; \beta) + \varepsilon_h > \log \tau)|x_h] \\ &= \sum_h \Pr[\varepsilon_h > \log \tau - m(x_h; \beta)] \\ &= \sum_h (1 - F_\varepsilon(\log \tau - m(x_h; \beta); \sigma)), \end{aligned}$$

where  $1(a > b)$  denotes the indicator function that equals 1 if  $a > b$  and 0 otherwise. In practice of course  $\beta$  and  $\sigma$  will have to be replaced with their respective estimators  $\hat{\beta}$  and  $\hat{\sigma}$ . Similarly, an estimator for  $n(\tau, y)$  can be obtained:

$$\begin{aligned} \hat{n}(\tau, y) &= \sum_h E[1(\tau < Y_h \leq y)|x_h] \\ &= \sum_h E[1(m(x_h; \beta) + \varepsilon_h \leq \log y)|x_h] - E[1(m(x_h; \beta) + \varepsilon_h \leq \log \tau)|x_h] \\ &= \sum_h \Pr[\varepsilon_h \leq \log y - m(x_h; \beta)] - \Pr[\varepsilon_h \leq \log \tau - m(x_h; \beta)] \\ &= \sum_h F_\varepsilon(\log y - m(x_h; \beta); \sigma) - F_\varepsilon(\log \tau - m(x_h; \beta); \sigma). \end{aligned}$$

Given  $\hat{n}(\tau, y)$  and  $\hat{n}(\tau)$ , we may construct the estimators  $\hat{F}_2(y) = \hat{n}(\tau, y) / \hat{n}(\tau)$  and  $\hat{\lambda} = \hat{n}(\tau) / n$ . Combined with the estimator for  $F_1(y)$ , which is estimated using DB-1 (i.e. the household income survey), we have all we need to estimate  $F(y)$  (see equation (1)), the income distribution for the complete population. This, in turn, is all we need to compute any inequality measure of choice.

No assumptions have been made about the distribution of  $x_h$  at this point. Let us assume that the top end of the distribution of  $x_h$  can be described by a Pareto distribution.

**Assumption 2.** Let  $G_2(x)$  denote the distribution function of  $x$  conditional on  $x > x_0$ . It is assumed that  $G_2(x)$  follows a Pareto distribution with shape parameters  $\alpha$ :

$$G_2(x) = 1 - \left(\frac{x}{x_0}\right)^{-\alpha}.$$

For ease of exposition, let us also assume that the income threshold  $\tau$  is set sufficiently high that  $Y > \tau$  implies  $X > x_0$ .



**Assumption 3.**

$$\Pr[Y \leq y | Y > \tau] = \Pr[Y \leq y | Y > \tau, X > x_0].$$

It then follows that top incomes, exceeding the income threshold  $\tau$ , are also Pareto distributed.

**Proposition 4.** *Given Assumptions 1, 2, and 3,  $F_2(y)$  follows a Pareto distribution with shape parameter  $\theta = \alpha/\beta_1$ :*

$$(14) \quad F_2(y) = \Pr[Y \leq y | Y > \tau] = 1 - \left(\frac{y}{\tau}\right)^{-\theta}.$$

**Proof.** By Assumption 3 we have:

$$\Pr[Y \leq y | Y > \tau] = \Pr[Y \leq y | Y > \tau, X > x_0].$$

This is equivalent to the following:

$$(15) \quad \Pr[Y \leq y | Y > \tau, X > x_0] = \frac{\Pr[\tau < Y \leq y | X > x_0]}{\Pr[Y > \tau | X > x_0]}$$

$$(16) \quad = \frac{\Pr[Y \leq y | X > x_0] - \Pr[Y \leq \tau | X > x_0]}{\Pr[Y > \tau | X > x_0]}.$$

Appealing to Assumptions 1 and 2, the term  $\Pr[Y \leq y | X > x_0]$  is seen to solve

$$\begin{aligned} \Pr[Y \leq y | X > x_0] &= \Pr[\exp(\beta_0 + \varepsilon)X^{\beta_1} \leq y | X > x_0] \\ &= \Pr\left[X \leq \exp(-\varepsilon/\beta_1) \left(\frac{y}{\exp(\beta_0)}\right)^{1/\beta_1} \mid X > x_0\right] \\ &= E_\varepsilon \left[ G_2 \left( \exp(-\varepsilon/\beta_1) \left(\frac{y}{\exp(\beta_0)}\right)^{1/\beta_1} \right) \right] \\ &= E_\varepsilon \left[ 1 - \exp(\alpha\varepsilon/\beta_1) x_0^\alpha \left(\frac{y}{\exp(\beta_0)}\right)^{-\alpha/\beta_1} \right] \\ &= 1 - E_\varepsilon[\exp(\alpha\varepsilon/\beta_1)] x_0^\alpha \left(\frac{y}{\exp(\beta_0)}\right)^{-\alpha/\beta_1} \\ &= 1 - y_0^\theta y^{-\theta}, \end{aligned}$$

with  $\theta = \alpha/\beta_1$  and  $y_0 = M_\varepsilon^{1/\theta}(\theta) \exp(\beta_0) x_0^{\beta_1}$ , where  $M_\varepsilon(t)$  denotes the moment generating function of  $\varepsilon$ ; that is,  $M_\varepsilon(t) = E[\exp(t\varepsilon)]$ . By extension, we have  $\Pr[Y \leq \tau | X > x_0] = 1 - y_0^\theta \tau^{-\theta}$ .

Substitution of the expressions for  $\Pr[Y \leq y|X > x_0]$  and  $\Pr[Y \leq \tau|X > x_0]$  into equation (16) yields the following:

$$\frac{\Pr[Y \leq y|X > x_0] - \Pr[Y \leq \tau|X > x_0]}{\Pr[Y > \tau|X > x_0]} = \frac{1 - y_0^\theta y^{-\theta} - (1 - y_0^\theta \tau^{-\theta})}{1 - (1 - y_0^\theta \tau^{-\theta})} = 1 - \tau^\theta y^{-\theta},$$

which completes the proof. ■

Note that  $\theta$  controls the thickness of the top end of the income distribution, which is a key determinant of income inequality; the smaller the value of the tail index  $\theta$ , the larger the proportion of high incomes, and the higher the value of inequality. Under the assumption that top incomes are Pareto distributed, the mean top income level takes on the following form:

$$(17) \quad E[Y|Y > \tau] = \left(\frac{\theta}{\theta - 1}\right)\tau.$$

This mean top income level features in the computation of the top income shares as well as the computation of the between-inequality components.<sup>11</sup>

### 2.2.2 The Population Underlying the Top Income Database Is Unclear

Let us next address the challenge that emerges when the data underlying DB-2 are not necessarily representative of the whole population (i.e. households with incomes exceeding  $\tau$ ). Consider the possibility that DB-2 has “over-sampled” some and “under-sampled” other households among the top earners, such that DB-2 no longer yields a consistent estimator for  $F_2(y)$  unless some corrective efforts are made. This is a rather realistic scenario, as the data may constitute a series of transactions or listing prices rather than a proper sample drawn from the target population. For ease of exposition, we will assume that DB-2 observes actual household incomes and not predictors of income, so that we may focus exclusively on the challenges presented in this section.

We will assume that the data are representative for selected subpopulations and that a representative “sample” can be obtained by anchoring DB-2 to some known population totals. Suppose that the target population can be subdivided into  $D$  districts with  $d=1, \dots, D$  indicating the district. The top income distribution for district  $d$  will be denoted by  $F_{2,d}(y) = \Pr[Y \leq y|Y > \tau, \text{district } d]$ . By extension, let  $F_{1,d}(y) = \Pr[Y \leq y|Y \leq \tau, \text{district } d]$ . Using this notation, the complete income distribution for district  $d$ , denoted  $F_d(y)$ , satisfies the following:

<sup>11</sup>As an alternative to assuming a Pareto distribution for the top tail, and estimating the tail index parameter, one could also appeal to multiple imputation methods (see, e.g. Doudich *et al.*, 2016). This approach might, in fact, be more practical in the event that a more flexible functional form for  $m(x_h; \beta)$  is being considered.

$$(18) \quad F_d(y) = \begin{cases} (1-\lambda_d)F_{1,d}(y), & y \leq \tau, \\ (1-\lambda_d) + \lambda_d F_{2,d}(y), & y > \tau, \end{cases}$$

where  $\lambda_d = \Pr[Y > \tau | \text{district } d]$ . The density functions corresponding to  $F_{1,d}(y)$ ,  $F_{2,d}(y)$ , and  $F_d(y)$  will be denoted by  $f_{1,d}(y)$ ,  $f_{2,d}(y)$ , and  $f_d(y)$ , respectively.

By definition, the distribution of top incomes for the whole population solves

$$(19) \quad F_2(y) = \sum_d F_{2,d}(y) P_{2,d},$$

with  $P_{2,d} = \Pr[Y > \tau, \text{district } d]$ . These mixing probabilities permit the following decomposition:

$$(20) \quad P_{2,d} = \lambda_d \pi_d,$$

where  $\pi_d$  denotes the share of the total population (regardless of income) residing in district  $d$ . We make the following assumption.

**Assumption 5.** *It is assumed that:*

- *The data at hand permit consistent estimation of  $(F_{2,d}, f_{2,d})$  and  $(F_{1,d}, f_{1,d})$  for all  $d$ .*
- *The district population shares  $\{\pi_d\}$  are known.*

That leaves  $\lambda_d = \Pr[Y > \tau | \text{district } d]$  as the only unknown that needs to be estimated. One way to estimate  $\lambda_d$  is to impose the assumption that  $f_d(y)$  is a continuous function.

**Assumption 6.**  *$f_d(y)$  is a continuous function of  $y$ .*

Let  $\hat{f}_{1,d}(\tau)$  and  $\hat{f}_{2,d}(\tau)$  denote the estimators for  $f_{1,d}(\tau)$  and  $f_{2,d}(\tau)$ , respectively. Assumption 5 ensures that these are consistent estimators. The following proposition derives an estimator for  $\lambda_d$  by appealing to Assumption 6.

**Proposition 7.** *Let  $\hat{f}_{k,d}(y)$  denote a consistent estimator for  $f_{k,d}(y)$  for  $k = 1, 2$ . Under Assumption 6,  $\hat{\lambda}_d$  presented below provides a consistent estimator for  $\lambda_d$ :*

$$(21) \quad \hat{\lambda}_d = \frac{\hat{f}_{1,d}(\tau)}{\hat{f}_{1,d}(\tau) + \hat{f}_{2,d}(\tau)}.$$

**Proof.** Evaluation of the first-order derivative of  $F_d(y)$  from equation (18) with respect to  $y$  yields:

$$(22) \quad f_d(y) = \begin{cases} (1-\lambda_d)f_{1,d}(y) & y \leq \tau, \\ \lambda_d f_{2,d}(y) & y > \tau. \end{cases}$$

By Assumption 6,  $f_d(y)$  is continuous in  $y$ , which imposes that  $(1-\lambda_d)f_{1,d}(y) = \lambda_d f_{2,d}(y)$  for  $y = \tau$ . Rearranging the terms in this equality gives us the following solution for  $\lambda_d$ :

$$(23) \quad \lambda_d = \frac{f_{1,d}(\tau)}{f_{1,d}(\tau) + f_{2,d}(\tau)}.$$

The estimator for  $\lambda$  is obtained by replacing  $f_{1,d}(\tau)$  and  $f_{2,d}(\tau)$  with their estimators. Provided that all terms on the right-hand side of equation (23) are consistently estimated, which is guaranteed by Assumption 5, it follows that the estimator for  $\lambda_d$  will be consistent. ■

Finally, note that the subgroup inequality decompositions presented in Section 2.1 can readily be extended to accommodate the subdivision of the top tail into  $D$  districts. (Note that the bottom segment can in principle stay as is; that is, it need not be subdivided into districts.) Let us denote the income share going to the top tail from district  $d$  by  $s_d = P_{2,d}(\mu_{2,d}/\mu)$ , where  $\mu_{2,d} = E[Y|Y > \tau, \text{district } d]$ . Note that the population and income shares corresponding to the bottom segment now solve  $1 - \sum_d \lambda_d$  and  $1 - \sum_d s_d$ , respectively. Similarly, let us denote the Theil index or the mean-log-deviation for the top incomes from district  $d$  by  $T_{2,d}$  and  $MLD_{2,d}$ , respectively. Using this notation, the decomposition of the Theil index and the mean-log-deviation into the  $1+d$  subgroups is seen to solve the following:

$$MLD = (1 - \sum_d \lambda_d)MLD_1 + \sum_d \lambda_d MLD_{2,d} + \log(\mu) - \log\left(\mu_1^{(1 - \sum_d \lambda_d)} \prod_d \mu_{2,d}^{\lambda_d}\right),$$

$$T = (1 - \sum_d s_d)T_1 + \sum_d s_d T_{2,d} + \log\left(\mu_1^{(1 - \sum_d s_d)} \prod_d \mu_{2,d}^{s_d}\right) - \log(\mu).$$

### 3. DATA

This paper uses two different types of datasets: (1) Household Income, Expenditure, and Consumption Survey (HIECS) data; and (2) listings of homes for sale derived from (large) real estate databases. All data used in this study are for Egypt. The HIECS is from 2008/9. The house price data are slightly more recent, covering the period from early 2014 to 2015, and come from two different real estate firms. Details are given below.

#### 3.1. Egyptian Household Income, Expenditure, and Consumption Survey

The Egypt HIECS 2008/9 was conducted by the Central Agency for Public Mobilization and Statistics (CAPMAS). We were given a 50 percent sample of

the survey (approximately 24,000 observations).<sup>12</sup> Throughout the paper, our welfare aggregate is household income per capita, adjusted for temporal differences in prices by deflating nominal values by a monthly price index.<sup>13</sup>

In most developing countries, consumption expenditure is the welfare aggregate used in poverty and inequality measurement. Compared to income, consumption expenditure produces lower estimates of inequality, especially at the top. This 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 (e.g. for the U.S., Aguiar and Bils, 2015). An argument for using consumption instead of income is that data on the former are often of a higher quality in developing and emerging economies, especially at the bottom tail, and are less vulnerable to idiosyncratic noise as households tend to smooth their consumption over time. For their study of top incomes in Egypt, Hlasny and Verme (2016) use both income and consumption expenditure as their welfare measures, and they find the HIECS income data to be of good quality. Given that our paper concentrates on the top tail, where consumption data are arguably less accurate (e.g. consumption of durables), we therefore opt to use income instead of consumption data. Our earlier working paper (van der Weide *et al.*, 2016) provides results using consumption expenditure for comparison.

As discussed in detail in Verme *et al.* (2014), inequality in Egypt as assessed from household surveys is low and has even declined in the decade before the 2011 revolution. The Gini coefficient of consumption expenditure declined by around two percentage points, from 32.8 in 2000 to 30.8 in 2009.<sup>14</sup> Our paper tests whether the low estimate in 2009 is robust to replacing the top tail of the income distribution with an estimate that is obtained using a combination of household survey and house price data.

### 3.2. Real Estate Data

In late 2014/early 2015, we obtained data on houses and apartments for sale from two Egyptian real estate firms: Betak-online and Bezaat.<sup>15</sup> The two rank among the larger real estate firms whose listing databases can be accessed online; analogous to Redfin and Zillow in the U.S. The data differ in detail, but a listing typically consists of the asking price, the location (the city or a further subdivision), and the date when it was listed. Interviews with the Ministry of Housing in Cairo confirmed that the listing price provides a good approximation to the actual sales price.<sup>16</sup> We keep listings classified as houses, apartments, flats, or

<sup>12</sup>Hlasny and Verme (2016) were able to access the 100 percent sample on site at CAPMAS. Our data access was provided by the Economic Research Forum's Open Access Micro Data Initiative (OAMDI, 2014).

<sup>13</sup>We use CAPMAS's monthly price index for all items for urban Egypt throughout the paper. Note that spatial price differences are not accounted for. Researchers often adjust for urban-rural price differences. Such an adjustment is inconsequential for this paper, given that we only cover urban areas. For a recent discussion of challenges with real income and consumption measurement, see, for example, van Veelen (2002) and van Veelen and van der Weide (2008).

<sup>14</sup>Source: PovcalNet, accessed October 31, 2015.

<sup>15</sup>The URLs are, respectively, <http://www.betakonline.com> and <http://www.bezaat.com>.

<sup>16</sup>For our purposes, it is sufficient that the actual price is proportional to the listing price.

TABLE 1  
THE ANNUAL INCOME OF THE TOP EARNERS IN EGYPT (USD, 2009 PRICES)

	Median
Top 5% household income survey	13,737
Top 1% household income survey	27,187
Top 0.5% household income survey	35,740
CEO total pay	68,970
CFO total pay	54,563
Top 1% household survey + house price data	32,628

villas, since these refer to private housing. There are a number of other types of listings which we exclude, the three largest groups being land, shop, and chalet.

The model that relates the value of the house to household income (per capita) is estimated using the household survey data, which report (imputed) rents, not property prices.<sup>17</sup> We will be assuming that rent and sale (or listing) prices are proportional to each other, which is sufficient for our needs.

The data from Betak-online and Bezaat, respectively, cover the periods January 2014 to November 2014 and November 2014 to January 2015. To adjust for temporal price fluctuations, all house price data are expressed in January 2014 prices, using the monthly price index. While the household survey is from 2009, compared to 2014 for the real estate data, there is no real need to express the values in prices for the same year; that is, to inflate the 2009 incomes to 2014 prices or to deflate the house prices to 2009 prices. Instead, we will be assuming that the Pareto tail index associated with the top tail of the income distribution is stable over the 2009–14 period.

### 3.3. *Does the Household Survey Indeed Omit the Rich?*

One way of illustrating whether the household data underrepresent the top part of the distribution is to compare some of the characteristics of the top 1 percent in the household survey with those of senior Egyptian executives. We would expect senior executives to be close to this part of the income distribution. Table 1 reports median incomes for a number of top groups in the household survey, as well as median salaries for CEOs and CFOs.<sup>18</sup> The data on executive pay are drawn from Payscale, an online information company providing current information on salary, benefits, and compensation by type of job, location, and other characteristics.<sup>19</sup>

<sup>17</sup>The rent variable combines actual rents with imputed rents for owner-occupiers as well as for households paying a reduced rent or housed for free.

<sup>18</sup>A number of differences between these datasets complicate the comparison. For instance, the household survey reports disposable income (i.e. all income sources, after taxes), while executive compensation most likely refers to gross earnings. Furthermore, the comparison assumes that executives reside in single-earner households and have no other sources of income.

<sup>19</sup>If anything, we expect the data on executive pay to be on the conservative side. The senior executives surveyed by Payscale are either chief executive officers (CEOs) or chief financial officers (CFOs) in Egyptian firms. CEOs and CFOs have the highest reported median compensation among survey participants. The total compensation of senior executives refers to 2015. The values in the table are deflated and converted from EGP into USD using annual average inflation and exchange rate data from the World Bank's World Development Indicators.

It is clear from Table 1 that the highest incomes in the household income survey are considerably below the earnings of senior executives. Similarly, in Vietnam, the top salaries recorded in the household survey are less than half of average executive salaries obtained from corporate salary surveys (World Bank, 2014). In the case of Argentina, Alvaredo (2010) finds that while the tax data have almost 700 observations with incomes exceeding 1 million USD, there are none in the Argentine household survey. In a comparison of 16 Latin American household surveys, the ten richest households have incomes similar to a managerial wage, which is arguably substantially smaller than the incomes of top capital owners (Székely and Hilgert, 1999). The last row in Table 1 reports an estimate of median household income among the top 1 percent that is obtained by combining household survey and house price data.<sup>20</sup> It can be seen that this yields a correction toward the CFO/CEO salaries.

#### 4. EMPIRICAL APPLICATION

This section presents our empirical application to Egypt. As outlined in the methodology section, we combine data on household incomes with data on house prices. The household incomes are obtained from the 2008/9 Egypt Household Income, Expenditure, and Consumption Survey (HIECS), which is also used for Egypt's official estimates of poverty and inequality. The house prices represent listing prices for houses that have been put up for sale via two large real estate firms operating in Egypt. We use the real estate database to estimate the top end, defined as the top 5 percent, of the income distribution. The "bottom" 95 percent of the income distribution is estimated using the HIECS.

The following practical decisions and assumptions are made: (a) we restrict the analysis to urban Egypt only (this can be extended to apply to all of Egypt under the assumption that rural households do not rank in the top of the income distribution in Egypt); (b) it is assumed that house price quotes are proportional to (imputed) rental values (as the household income survey contains data on rents only, and we rely on the survey to identify the relationship between house value and household income);<sup>21</sup> (c) it is assumed that the Pareto tail index of the income distribution has been stable between 2008/9 (the time of the survey) and 2014/15 (the time of the house price database); (d) it is assumed that one house constitutes one household (the fact that top income households could be associated with multiple houses may lead us to underestimate inequality) and that all houses are domestically owned; and (e) we will only be using house price data for Cairo and Alexandria to estimate the top tails of their respective income distributions. For the rest of urban Egypt, the entire income distribution will be estimated using the HIECS. The latter decision is motivated by the fact that: (i) the lion-share of the

<sup>20</sup>For this purpose, we assume a household size of 3 to convert per capita incomes to household incomes and we set the income cutoff above which the income distribution is assumed to be Pareto at the 99th percentile. We use the house price data, combined with the household survey data, to estimate the corresponding Pareto tail index. For the main empirical application presented in Section 4, the cutoff above which the income distribution is estimated using a combination of household survey and house price data is set at the 95th percentile of the income distribution.

<sup>21</sup>Under a non-arbitrage condition, house prices and rental values are expected to move in parallel (see, e.g. Himmelberg *et al.*, 2005).

TABLE 2  
THE NUMBER OF OBSERVATIONS USED

Subgroup	Database		
	Betak-online	Bezaat	HIECS
Cairo	5,970	8,502	2,592
Alexandria	1,338	2,021	1,400
Urban Egypt			10,763

“rich” that are missing or whose incomes are understated in the HIECS arguably reside in either Cairo or Alexandria; and (ii) the real estate markets are most developed in Cairo and Alexandria, such that the coverage and the quality of the house price data are highest for these two cities (henceforward, we will refer to these as “districts”).

Table 2 provides some basic statistics on the number of observations available to us.<sup>22</sup> For the house price databases, we only counted observations above the median house price value (which practically coincides with the mode of the house price density). Since we are interested in the top tail behavior of the house price distribution, we do not use the lower house price values.

The following sections proceed with the empirical application, which combines the household income survey and the house price data. A validation of our methodology in a controlled setting where only the survey data are used can be found in the Annex.

#### 4.1. Pareto Tail Index Estimated on Income Survey Data

This subsection presents first estimates of the Pareto tail index of Cairo’s and Alexandria’s income distributions by using household survey data only. These estimates will serve as a reference point. Under the assumption of Pareto distributed top tails, we have that  $1 - F_2(y) = (\frac{y}{\tau})^{-\theta}$ . Rearranging terms yields:

$$(24) \quad \log(y) = \log(\tau) - \frac{1}{\theta} \log(1 - F_2(y)).$$

If this assumption holds true, a plot of  $\log(y)$  against  $-\log(1 - F_2(y))$  should reveal a linear relationship with a slope parameter equal to  $\frac{1}{\theta}$ . Figure 1 provides this plot, using the top 5 percent of the household income data from the HIECS.<sup>23</sup> For the majority of data points, a linear relationship seems to provide a reasonable fit. A deviation from linearity can be observed however toward the far

<sup>22</sup>The region variable in HIECS refers to the 27 governorates of Egypt, with the survey samples being representative of the entire governorates. The Cairo governorate is entirely urban, while the urban population share is 98.8 percent in the Alexandria governorate (CAPMAS, 2015). Therefore, our study of urban Egypt covers the overwhelming majority of both governorates. The regional variables available in the house price databases allows us to construct the same governorate identifier that is used by the HIECS. Also here, the data cover the entire governorates, with both capital and non-capital cities being represented in the house price databases.

<sup>23</sup>The 95th percentile in the urban per capita income distribution for 2009 is estimated at 13,400 EGP (in January 2007 prices).



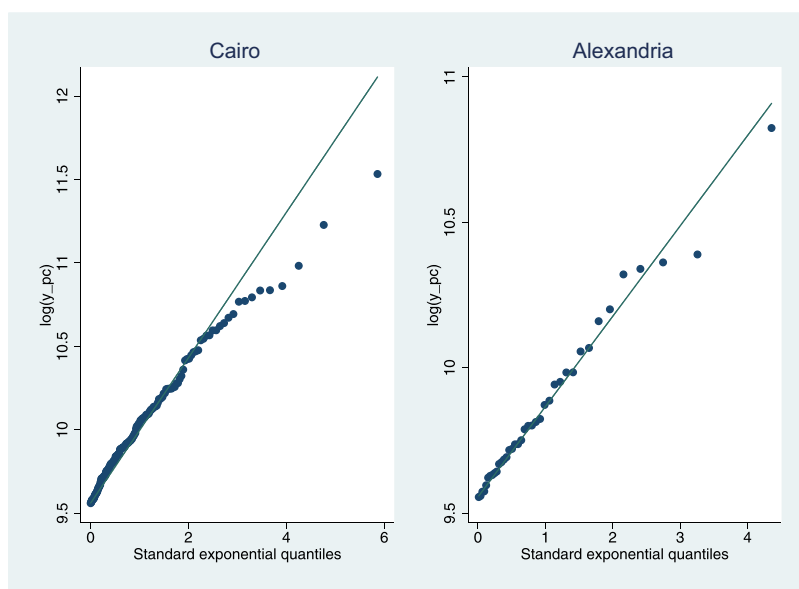


Figure 1. The Pareto Quantile Plot for Household Income per capita (Household Survey) [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

end of the income spectrum, where the slope appears to fall. Consequently, we should expect estimates of  $\theta$  to come out higher if we were to increase the income threshold above which observations were included.

Figure 2 plots the maximum-likelihood (ML) estimates of  $\theta$  for different values of the number of top observations used, ranging from the top 15 percent (85th percentile and up) to the top 5 percent of income observations (95th percentile and up). The gray area indicates the 95 percent confidence interval, which is seen to widen as the number of observations is reduced. It is also confirmed that for both Cairo and Alexandria the tail index is estimated to be higher at higher income thresholds (i.e. when the number of observations is reduced toward the top end), which is consistent with what we observed in Figure 1. The dotted line indicates the median level of the tail index (taken over all estimates within the plotted range). These will serve as our benchmark estimates of  $\theta$ .

It can also be observed that the HIECS estimates the top tail of the income distribution to be heavier (lower tail index) in Cairo than in Alexandria. Put differently, top income shares and income inequality are estimated to be highest in Cairo, which is arguably what one would expect. Relative ordering put aside, the question is whether the tail indices are being overestimated; that is, whether the thickness of the top tails are being underestimated. The next subsection will address this question by consulting data on house prices.

#### 4.2. Estimating the Tail Index Using Both Income and House Price Data

We will go through the following steps in order to estimate the Pareto tail index  $\theta$  by combining data on household income from the HIECS with data on house prices.

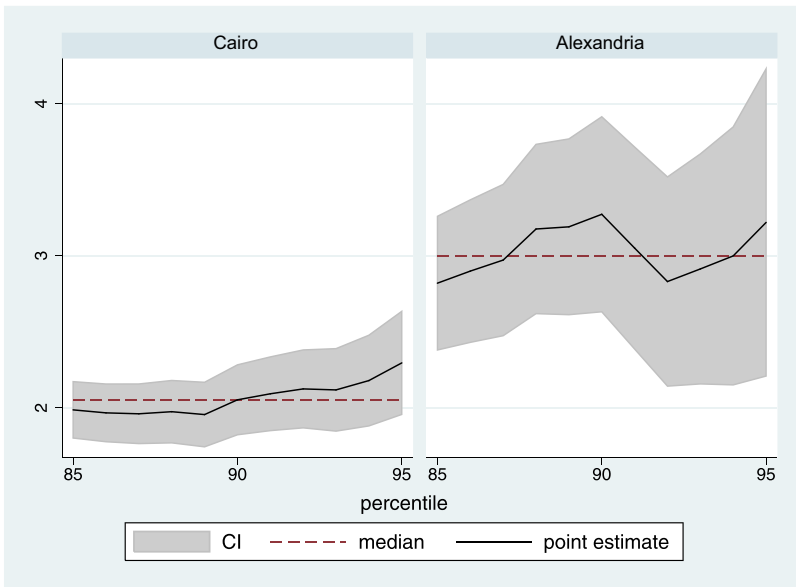


Figure 2. Pareto Tail Index Estimates for Household Income per capita (Household Survey) [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

First, we estimate the tail index associated with the top end of the house price distributions in Cairo and Alexandria, which we denoted  $\alpha$  (see Assumption 2). Next, we estimate the model from Assumption 1 that provides a link between house prices and household incomes, where it is particularly parameter  $\beta_1$  in which we are interested. With the estimators  $\hat{\alpha}$  and  $\hat{\beta}_1$  in hand, for Cairo and Alexandria separately, we apply Proposition 4 and obtain  $\hat{\theta}_{mix} = \hat{\alpha} / \hat{\beta}_1$  as an alternative estimator for  $\theta$ .

Figure 3 plots  $\log(x)$  against  $-\log(1 - G_2(x))$ , analogous to Figure 1 but now using data on house prices (i.e.  $x$  denotes the listing price of a house). This plot uses the top 5 percent of above median value house prices from the respective house price databases (Betak-online and Bazaar). While a linear model appears to fit the data reasonably well, which supports the Pareto assumption, a deviation from linearity can be observed toward the top of the house price distribution. This non-linearity at the top is also observed for the household income data from the HIECS (see Figure 1), albeit more pronounced for the house price data. The pattern is most noticeable for Cairo.

Figure 4 gives us an idea of the range of values that  $\alpha$  might attain by plotting estimates of the tail index as we vary the database and the number of top observations used for estimation. Note that this figure is analogous to Figure 2. We omitted the confidence intervals in this case as they are small in comparison to the differences observed between the databases. The dotted line indicates our estimate of  $\alpha$ ; it is obtained as the median value of  $\hat{\alpha}$  obtained over the two databases and between the percentiles 75 and 92 (i.e. between the top 25 and 8 percent). In the case of Alexandria, the estimate roughly corresponds to a range where  $\hat{\alpha}$  is found to level off. For Cairo, it proved harder to find such a range. Our estimator is arguably on the conservative side in this case; our data appear to

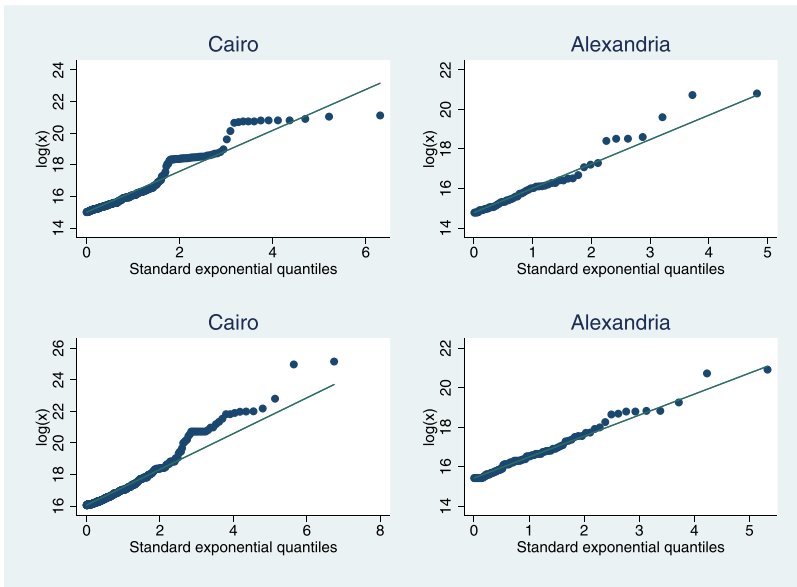


Figure 3. The Pareto Quantile Plot for House Prices (Real Estate Data): Top Row, Betak-online; Bottom Row, Bezaat [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

indicate that the tail index for Cairo is more likely to be lower than higher. In other words, if anything, we may be slightly underestimating the top income share (and hence inequality) for Cairo. Obviously, where we draw the line for  $\hat{\alpha}$  is to a

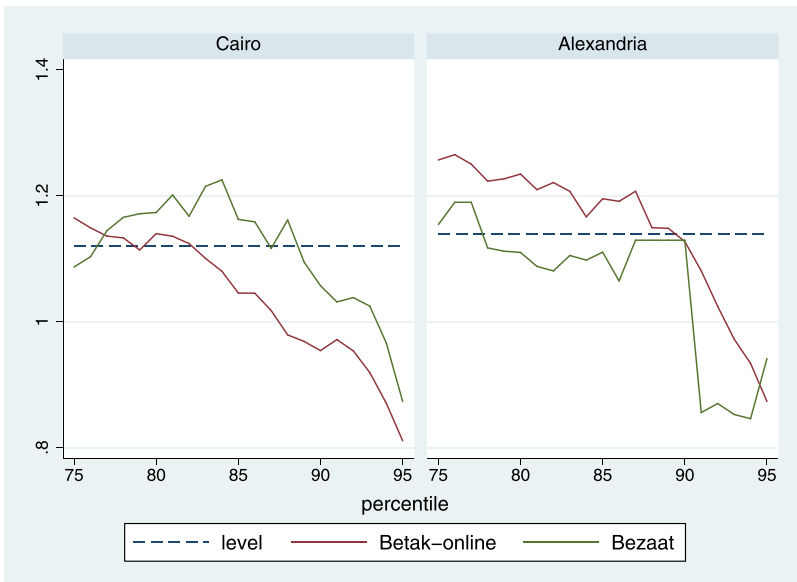


Figure 4. Pareto Tail Index Estimates for House Prices (Real Estate Data) [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

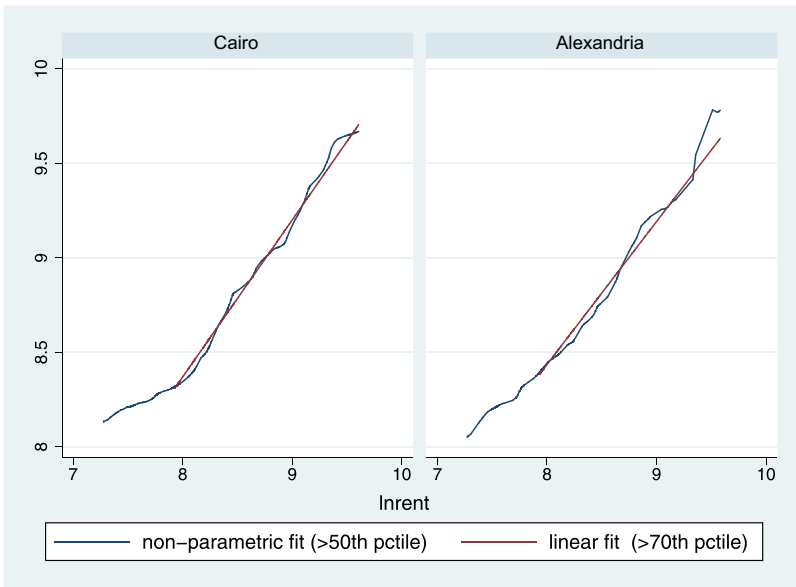


Figure 5. Household Income per capita versus Imputed Rent (Log-Log, Household Survey) [Colour figure can be viewed at [wileyonlinelibrary.com](http://wileyonlinelibrary.com)]

certain degree arbitrary. Toward the end of Section 4.3, we will briefly comment on how the range of  $\alpha$  observed here may translate into a range for  $\theta$  and by implication a range for estimated levels of inequality.

Next, we need estimates of  $\beta_1$ . Here, we fully rely on data from the HIECS. Before we imposed a functional form on  $m(x)$ , which describes the relationship between household income per capita and the value of the household’s house (captured by imputed rent), we first fitted a non-parametric kernel regression to the data (for Cairo and Alexandria separately). The results are presented in Figure 5. It is found that a linear model captures the relationship between log of household income and log of (imputed) rent reasonably well, particularly in the case of Cairo. Alexandria shows a degree of concavity, but also here a linear model arguably provides a good fit for high values of rent and household income; see the fitted linear lines included in the figure.

Estimates of  $\beta_1$  appear to be less sensitive to where we place the cutoff for the data included in the estimation when compared to estimates of  $\alpha$ . See Figure 6, which investigates how  $\hat{\beta}_1$  varies with the number of top observations included in the regression. The gray area indicates the 95 percent confidence interval. Note how  $\hat{\beta}_1$  is reasonably stable across the different cutoffs considered, which is consistent with the degree of linearity observed in Figure 5. The dotted lines denote the estimates that will be used in our analysis (see the values reported the first column of Table 3), which are obtained as the value of  $\hat{\beta}_1$  for the top 10 percent (90th percentile) for Cairo and for the top 13 percent (87th percentile) for Alexandria.<sup>24</sup>

<sup>24</sup>Having more observations for Cairo than for Alexandria allows for a higher income percentile cutoff.

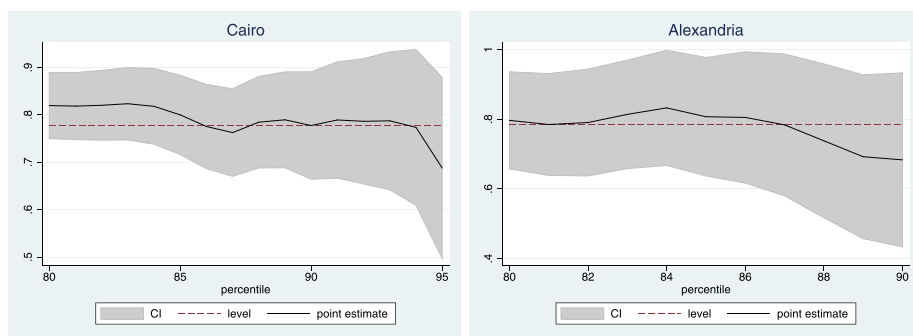


Figure 6. Estimates of  $\beta_1$  Using Increasingly Smaller Numbers of Top Observations (Household Survey) [Colour figure can be viewed at wileyonlinelibrary.com]

TABLE 3  
ESTIMATES OF  $\beta_1$ ,  $\alpha$ ,  $\theta$ , AND  $\gamma$

Subgroup	$\hat{\beta}_1$	$\hat{\alpha}$	$\hat{\theta}_{mix}$	$\hat{\theta}_{svy}$	$\hat{\gamma}_{mix}$	$\hat{\gamma}_{svy}$
Cairo	0.778	1.120	1.441	2.052	3.270	1.951
Alexandria	0.784	1.139	1.453	2.999	3.208	1.500

What does this mean for  $\theta$ ? Our findings are summarized in Table 3, which shows the estimator  $\hat{\theta}_{mix} = \hat{\alpha} / \hat{\beta}_1$  as well as the individual components  $\hat{\alpha}$  and  $\hat{\beta}_1$  that go into the estimator ( $\hat{\gamma}$  denotes the inverted Pareto Lorenz coefficient, defined as  $\hat{\gamma} = \frac{\hat{\theta}}{\hat{\theta}-1}$ ). For comparison, we also include the estimator  $\hat{\theta}_{svy}$  that is obtained using data from the HIECS only (see Section 4.1). The data on house prices give us reason to believe that the top tail of the income distribution is underestimated in Egypt when relying on household survey data only, as is evidenced by the fact that  $\hat{\theta}_{mix}$  is visibly smaller than  $\hat{\theta}_{svy}$ .

#### 4.3. Main Results: Re-estimating Inequality for Egypt

Having new estimates of the Pareto tail indices for the respective income distributions of Cairo and Alexandria is not enough. To see what this means for total inequality for (urban) Egypt, we also need estimates of the share of the population that resides in the respective metropolitan areas and enjoys incomes above  $\tau$ ; that is, estimates of  $\Pr[Y > \tau, \text{district } d]$  for  $d = \text{Cairo}, \text{Alexandria}$ . We estimate these by:  $\Pr[Y > \tau, \text{district } d] = \Pr[Y > \tau | \text{district } d] \Pr[\text{district } d]$ , where  $\Pr[\text{district } d]$  (the share of the urban population residing in district  $d$ ) is obtained from the most recent population census and where  $\Pr[Y > \tau | \text{district } d]$  is estimated using Proposition 7. For comparison, the latter is also estimated using data from the HIECS only. The two different estimators are denoted by  $\hat{\lambda}_{prop7}$  and  $\hat{\lambda}_{svy}$ , respectively.  $\Pr[Y > \tau, \text{district } d]$  and  $\Pr[\text{district } d]$  are denoted by  $P$  and  $\pi$ , respectively, such that  $\hat{P}_{prop7} = \pi \hat{\lambda}_{prop7}$  and  $\hat{P}_{svy} = \pi \hat{\lambda}_{svy}$  (where we have suppressed the subscript  $d$  for ease of notation). The estimates are presented in Table 4.<sup>25</sup>

<sup>25</sup>For other urban,  $\hat{\lambda}_{prop7} = \hat{\lambda}_{svy}$  and  $\hat{P}_{prop7} = \hat{P}_{svy}$  because we do not use any house price data for this group.

TABLE 4  
ESTIMATES OF  $\pi$ ,  $\hat{\lambda}$ , AND  $P$

Subgroup	$\pi$	$\hat{\lambda}_{prop7}$	$\hat{\lambda}_{svy}$	$\hat{P}_{prop7}$	$\hat{P}_{svy}$
Cairo	0.251	0.104	0.084	0.026	0.021
Alexandria	0.130	0.051	0.033	0.007	0.004
Other urban	0.619	0.027	0.027	0.017	0.017

TABLE 5  
ESTIMATES OF  $S$ ,  $GINI$ ,  $MLD$ , AND  $THEIL$  (FOR THE TOP TAIL)

Subgroup	$S_{5,mix}$	$S_{5,svy}$	$Gini_{mix}$	$Gini_{svy}$	$MLD_{mix}$	$MLD_{svy}$	$Theil_{mix}$	$Theil_{svy}$
Cairo	0.194	0.116	0.532	0.226	0.491	0.082	1.085	0.096
Alexandria	0.048	0.024	0.525	0.167	0.477	0.047	1.042	0.053
Other urban	0.071	0.081	0.296	0.296	0.150	0.150	0.210	0.210

Note that our estimate of  $\lambda$  finds that the percentage of households residing in Cairo and Alexandria with incomes exceeding  $\tau$  is larger than what the HIECS alone would have us believe. This, combined with the earlier observation that  $\hat{\theta}_{mix} < \hat{\theta}_{svy}$ , leads us to believe that relying on survey data alone will arguably underestimate both the number of households with high incomes as well as the size of their incomes (either because top income earners are missing in the survey or because they underreport their incomes, or both). Table 5 compares estimates of the income share of the top 5 percent ( $S_5$ ) obtained using the HIECS to those obtained using both the HIECS and the house price data.<sup>26</sup> The additional columns compare estimates of inequality among top income households (i.e. only including households whose income exceeds  $\tau$  and that reside in the respective district) for three different measures of inequality.

Estimates of total inequality for (urban) Egypt are obtained by adding estimates of bottom- and between-inequality to the estimates of top inequality reported in Table 5. Bottom inequality (i.e. inequality among households with income below  $\tau$ ) is estimated using the HIECS only. The between-inequality component is estimated using data from both sources, as it is a function of average income among top earners (which is a function of  $\theta$ ; see equation (17)) as well as a function of  $\lambda$  (in the case of  $MLD$ ) and of the top income share  $S$  (in the case of the Theil index); see equations (6) and (9). In the case of the Gini coefficient, we implement the approximate decomposition that is also used by Alvaredo (2011):  $Gini \approx (1 - \sum_d \lambda_d)(1 - \sum_d s_d)Gini_1 + \sum_d s_d$ .<sup>27</sup>

The total inequality estimates are presented in Table 6.<sup>28</sup> The survey-only estimate of the Gini coefficient for (urban) Egypt in 2008/9 stands at 38.5. This is

<sup>26</sup>For example, 19.4 percent of national (i.e. urban) income is owned by households that are both part of the national top 5 percent and reside in Cairo. The sum of the income shares in Table 5 adds up to the national top income share reported in Table 6.

<sup>27</sup>When the top group is small, the Gini decomposition from equation (3) can be approximated in this way.

<sup>28</sup>The estimates of the top income shares are obtained for different income cutoffs above which the income distribution is assumed to be Pareto (the parameters of which are estimated using house price data). For  $S_{10}$ , the income share of the top 10 percent, the cutoff was placed at the 90th percentile; for  $S_5$  we used the 95th percentile; and so on. The  $Gini$ ,  $MLD$ , and  $Theil$  measures of inequality are obtained using a cutoff at the 95th percentile.

TABLE 6  
ESTIMATES OF INEQUALITY FOR (URBAN) EGYPT IN 2008/9: SURVEY-ONLY VERSUS SURVEY + HOUSE PRICES

	Survey and House Prices	Survey Only
<i>Gini</i>	0.518	0.385
<i>MLD</i>	0.374	0.244
<i>Theil</i>	0.738	0.302
$S_{10}$	0.422	0.321
$S_5$	0.314	0.221
$S_1$	0.151	0.089

relatively low by international standards and hence would suggest that Egypt ranks among lower-inequality countries. Our estimate of the Gini coefficient is 51.8, which is considerably higher than the survey-only estimate. The level of top incomes recorded in the HIECS is found to be at odds with house prices observed toward the top end of the market in Cairo and Alexandria. Our estimates represent an attempt to correct for this. We repeated the analysis for other choices of inequality measures, specifically for the *MLD* and *Theil* measures. Noticeable increases in inequality can be observed for all measures considered. The magnitude of the adjustment is largest for the Theil index, which is consistent with the fact that the Theil index is most sensitive to the top tail of the income distribution when compared to the other two choices of inequality measures.

The precision of our estimate of inequality is largely determined by the precision with which we are able to estimate  $\alpha$  and  $\beta_1$  (provided that the assumptions under which the estimators have been derived reasonably apply to the data at hand). It is instructive to verify what level of inequality would be obtained using rather conservative values for  $\theta$ . Note that a most conservative estimate of  $\theta$  can be obtained by combining a value of  $\alpha$  from the top end of the estimated range with a value of  $\beta_1$  from the low end of the estimated range (but taking  $\beta_1 \geq \frac{1}{2}$ , which rules out housing expenditure shares that are convex increasing functions of household income; see the discussion following Assumption 1). For Cairo, this gives us a value of around 2.0 (1.2/0.60; see Figures 4 and 6). For Alexandria, we obtain a value that is around 2.5 (1.25/0.5; see Figures 4 and 6). Note that these values are still below the respective survey-only estimates of  $\theta$  (see Figure 2). In other words, even with very conservative estimates for  $\hat{\theta}_{mix}$ , we would still obtain estimates of inequality that are higher than the survey-only estimate. The estimate that we consider most reasonable finds a Gini coefficient for (urban) Egypt of 51.8, which is roughly 13 points higher than the survey-only estimate. Of course, by the same token, we may also be underestimating inequality. Working with values of  $\theta$  toward the lower end of our estimated range yields estimates of inequality that are noticeably higher than the mid-range Gini coefficient of 51.8.

Using our estimates for  $\theta$  (and  $\lambda$ ) in a back-of-the-envelope calculation, we find that there are approximately 300 households in Cairo whose household income exceeds 1 million USD per year. Although no other information on the number of millionaires in Egypt is currently available, this estimate seems rather conservative.

## 5. CONCLUDING REMARKS

A growing literature has shown that household surveys provide only limited information about top incomes and therefore underestimate income inequality. This paper presents a method that corrects for this underestimation. We use the household survey for the bottom part of the distribution and combine it with another data source that provides a better coverage of the top tail. The existing literature has restricted itself to the use of tax record data to capture the top tail. Unfortunately, income tax records are unavailable in many countries, including most of the developing world. Our method permits a much larger set of data for the top tail; the only requirements are that the data (i) contain a good predictor of household income and (ii) provide a good coverage of the top tail.

We apply this method to Egypt, where estimates of inequality based on household surveys alone are low by international standards. Using publicly available data from real estate listings to estimate the top tail of the income distribution, we find strong evidence that inequality in Egypt is being underestimated. The Gini index of income for urban Egypt is found to increase from 39 to 52 after correcting for the missing top tail. A natural next step would be to use data on house prices to estimate the top tail of the wealth distribution, and extend the analysis to other countries.

## REFERENCES

- Aguiar, M. and M. Bils, "Has Consumption Inequality Mirrored Income Inequality?" *American Economic Review*, 105(9), 2725–56, 2015.
- Alvaredo, F., "The rich in Argentina over the twentieth century, 1932–2004," in A. B. Atkinson and T. Piketty (eds), *Top Incomes: A Global Perspective*, Oxford University Press, Oxford, 253–98, 2010.
- , "A Note on the Relationship Between Top Income Shares and the Gini Coefficient," *Economics Letters*, 110(3), 274–77, 2011.
- Alvaredo, F., A. B. Atkinson, L. Chancel, T. Piketty, E. Saez and G. Zucman, "World Wealth and Income Database," 2017, <http://wid.world>
- Alvaredo, F. and J. Londoño Vélez, "High Incomes and Personal Taxation in a Developing Economy: Colombia 1993–2013," Working Paper 12, Commitment to Equity—CEQ, 2013.
- Alvaredo, F. and T. Piketty, "Measuring Top Incomes and Inequality in the Middle East: Data Limitations and Illustration with the Case of Egypt," Working Paper 832, ERF, 2014.
- Anand, S. and P. Segal, "The Global Distribution of Income," in A. B. Atkinson and F. Bourguignon (eds), *Handbook of Income Distribution*, volume 2A, Elsevier, Amsterdam, 2015.
- Assouad, L., "Top Incomes and Personal Taxation in Lebanon: An Exploration of Individual Tax Records 2005–2012," Master's thesis, Paris School of Economics, 2015.
- Atkinson, A. B., "Measuring Top Incomes: Methodological Issues," in A. B. Atkinson and T. Piketty (eds), *Top Incomes over the Twentieth Century: A Contrast Between Continental European and English-Speaking Countries*, Oxford University Press, Oxford, 18–42, 2007.
- Atkinson, A. B., T. Piketty, and E. Saez, "Top Incomes in the Long Run of History," *Journal of Economic Literature*, 49(1), 3–71, 2011.
- Burricand, C., "Transition from Survey Data to Registers in the French SILC Survey," in M. Jäntti, V.-M. Törmälehto, and E. Marlier (eds), *The Use of Registers in the Context of EU-SILC: Challenges and Opportunities*, European Union, Luxembourg, 2013.
- Cano, L., "Top Income Shares in a Growing South American Economy: Ecuador 2004–2011," Working Paper, University of Toulouse 1 Capitole—Lereps, 2015.
- CAPMAS (Central Agency for Public Mobilization and Statistics), "Egypt in Figures 2015," CAPMAS, 2015, <http://www.msrinternet.capmas.gov.eg/pdf/EgyptinFigures2015/EgyptinFigures/Tables/PDF/1-> (accessed December 20, 2016).
- Diaz-Bazan, T., "Measuring Inequality from Top to Bottom," Policy Research Working Paper, The World Bank, Washington, D.C., 2014.



- Doudich, M., A. Ezzrari, R. van der Weide, and P. Verme, "Estimating Quarterly Poverty Rates Using Labor Force Surveys: A Primer," *The World Bank Economic Review*, 30(3), 475–500, 2016.
- Himmelberg, C., C. Mayer, and T. Sinai, "Assessing High House Prices: Bubbles, Fundamentals and Misperceptions," *Journal of Economic Perspectives*, 19(4), 67–92, 2005.
- Hlasny, V. and P. Verme, "Top Incomes and the Measurement of Inequality in Egypt," *The World Bank Economic Review*, Advance Access published July 10, 2016.
- Jäntti, M., V.-M. Törmälehto, and E. Marlier, *The Use of Registers in the Context of EU-SILC: Challenges and Opportunities*, European Union, Luxembourg, 2013.
- Jenkins, S. P., "Pareto Models, Top Incomes and Recent Trends in UK Income Inequality," *Economica*, 84(334), 261–89, 2017.
- Kim, N. N. and J. Kim, "Sodok jipyo ui jaegumto" ["Reexamining Income Distribution Indices of Korea"], *Journal of Korean Economic Analysis*, 19(2), 1–57, 2013 (in Korean).
- Korinek, A., J. A. Mistiaen, and M. Ravallion, "Survey Nonresponse and the Distribution of Income," *The Journal of Economic Inequality*, 4(1), 33–55, 2006.
- Lakner, C. and B. Milanovic, "Global Income Distribution: From the Fall of the Berlin Wall to the Great Recession," *The World Bank Economic Review*, 30(2), 203–32, 2016.
- Larsen, E., "The Engel Curve of Owner-Occupied Housing Consumption," *Journal of Applied Economics*, 17(2), 325–52, 2014.
- Morelli, S., T. Smeeding, and J. Thompson, "Post-1970 Trends in Within-Country Inequality and Poverty: Rich and Middle-Income Countries," in A. B. Atkinson and F. Bourguignon (eds), *Handbook of Income Distribution*, volume 2A, Elsevier, Amsterdam, 2015.
- OAMDI, "Harmonized Household Income and Expenditure Surveys (HHIES)," Version 2.0 of Licensed Data Files; HIECS 2008/2009—Central Agency for Public Mobilization and Statistics (CAPMAS), Economic Research Forum (ERF), 2014, <http://erf.org.eg/data-portal/>
- Shorrocks, A. F., "The Class Of Additively Decomposable Inequality Measures," *Econometrica*, 48, 613–25, 1980.
- Székely, M. and M. Hilgert, "What's Behind the Inequality We Measure: An Investigation Using Latin American Data," Research Department Working Paper, Inter-American Development Bank, 1999.
- van der Weide, R., C. Lakner, and E. Ianchovichina, "Is Inequality Underestimated in Egypt? Evidence from House Prices," Policy Research Working Paper Series 7727, The World Bank, Washington, D.C., 2016.
- van Veelen, M., "An Impossibility Theorem Concerning Multilateral International Comparison Of Volumes," *Econometrica*, 70(1), 369–75, 2002.
- van Veelen, M. and R. van der Weide, "A Note on Different Approaches to Index Number Theory," *American Economic Review*, 98(4), 1722–30, 2008.
- Verme, P., B. Milanovic, S. Al-Shawarby, S. El Tawila, M. Gadallah, and E. A. A.El-Majeed, *Inside Inequality in the Arab Republic of Egypt: Facts and Perceptions across People, Time, and Space*, The World Bank, Washington, D.C., 2014.
- World Bank, "Taking Stock: An Update on Vietnam's Recent Economic Development," The World Bank, Washington, D.C., 2014.
- , "Chile: Distributional Effects of the 2014 Tax Reform," The World Bank, Washington, D.C., 2016.

## SUPPORTING INFORMATION

Additional supporting information may be found in the online version of this article at the publisher's web-site:

Appendix

A Small Validation Exercise: Re-estimating Inequality in the Survey after Dropping Top Incomes

**Figure A.1:** Simulated Non-response Probabilities as a Function of Household (Log) Income per capita

**Table A.1:** Estimates of  $\beta_1$ ,  $\alpha$ ,  $\theta$ ,  $\lambda$ , and  $P$

**Table A.2:** Estimates of Inequality for (Urban) Egypt in 2008/9: Survey-Only versus Survey + Imputed Rents