

# Parametric models of income distributions integrating misreporting and non-response mechanisms <sup>\*</sup>

Mathias Silva<sup>†‡</sup>

Aix Marseille Univ, CNRS, AMSE, Marseille, France

This version: November 7, 2023

## Abstract

Several representativeness issues affect the available data sources in studying populations' income distributions. High-income under-reporting and non-response issues have been evidenced to be particularly significant in the literature, due to their consequence in under-estimating income growth and inequality. This paper bridges several past parametric modelling attempts to account for high-income data issues in making parametric inference on income distributions at the population level. A unified parametric framework integrating parametric income distribution models and popular data replacing and reweighting corrections is developed. To exploit this framework for empirical analysis, an Approximate Bayesian Computation approach is developed. This approach updates prior beliefs on the population income distribution and the high-income data issues presumably affecting the available data by attempting to reproduce the observed income distribution under simulations from the parametric model. Applications on simulated and EU-SILC data illustrate the performance of the approach in studying population-level mean incomes and inequality from data potentially affected by these high-income issues.

*Keywords:* 'Missing rich', GB2, Bayesian Inference.

*JEL Code:* D31, C18, C11

---

<sup>\*</sup>I am grateful to Michel Lubrano, Stephen Bazen, Emmanuel Flachaire, Philippe Van Kerm, Markus Jäntti, Jonathan Goupille-Lebret, Duangkamon Chotikapanich, Sylvia Kaufmann, an anonymous referee, participants at the 8th European User Conference for EU-Microdata (GESIS, Eurostat), and seminar participants at the AMSE PhD seminar and ENS Lyon CERGIC Graduate seminar for helpful contributions and comments on earlier versions of this paper.

<sup>†</sup>Aix Marseille Univ, CNRS, AMSE, Marseille, France; and University of Lyon, Ecole Normale Supérieure de Lyon and Center for Economic Research on Governance, Inequality and Conflict.

<sup>‡</sup>Corresponding author: [mathias.silva-vazquez@univ-amu.fr](mailto:mathias.silva-vazquez@univ-amu.fr). Aix Marseille Univ, CNRS, AMSE, Marseille, France. AMSE - Aix-Marseille Université 5-9 Boulevard Bourdet 13001 Marseille, France. The project leading to this publication has received funding from the French government under the "France 2030" investment plan managed by the French National Research Agency (reference: ANR-17-EURE-0020) and from Excellence Initiative of Aix-Marseille University - A\*MIDEX. Declaration of interest: None. All R replication code for the results shown in this paper are available upon request and will soon be made available on the author's Github site.

# 1 Introduction

The recent literature on income inequality has paid increasing attention to the dynamics and the measurement of top incomes (e.g., [Atkinson and Piketty 2007](#), [Leigh 2009](#), [Atkinson et al. 2011](#), [Burkhauser et al. 2017](#)). The slowly-rising availability of tax data for research purposes along with findings concerning the recent rises in the share of incomes accumulated at the very top quantiles of the distribution (e.g., [Lakner and Milanovic 2016](#), [Alvaredo et al. 2018](#)) have jointly brought forward the multiple deficiencies affecting the typical methods and data sources used to study income distributions.

A robustly evidenced shortcoming of these conventional approaches involves the limited quality of publicly-available household survey data, the most commonly used data source on the subject, in capturing the magnitude and trends of the income shares of the highest incomes in their population (e.g., [Deaton 2005](#), [Burdín et al. 2014](#), [Jenkins 2017](#), [Lustig 2019](#), [Flachaire et al. 2022](#)).

Typically these measurement and coverage issues around the upper tail of a population’s income distribution imply non-random missing information in the data (i.e., the errors are more likely or larger in magnitude for higher incomes) and can therefore induce bias into any resulting distributional estimate. When ignored, this can have many clear policy implications when it leads to underestimation of income growth and inequality at the population level, along with a biased reading of their relationship and dynamics. This has motivated a vast literature on correction and estimation methods to overcome this data issue for the study of income distributions.

An important implication of the almost universal problem of missing or misreported high incomes is that any empirical strategy seeking to overcome it requires a decision on the magnitude and distribution of such errors affecting the data. As put forward by [Bourguignon \(2018\)](#), adjusting for measurement errors on high incomes requires a value for some or all of three key parameters: the income level beyond which measurement errors are to be corrected, the true population share of incomes above this level, and the share of under-covered population incomes.

Although external data sources can be instrumentally used to formulate informative choices for these parameters (e.g., [Atkinson and Piketty 2007](#), Chapter 2, [Bustos 2015](#), [Blanchet et al. 2022](#), [Jorda and Niño-Zarazúa 2019](#), [Flachaire et al. 2022](#)), correcting for measurement or coverage errors on high incomes is conditioned by the uncertainty around them. Broadly speaking, the precision with which inference can be made on a population’s income distribution depends on the uncertainty around the form and magnitude of measurement or coverage errors affecting the available data.

This paper proposes a new empirical strategy bridging several previous results in the income inequality literature. Firstly, a parametric modeling approach is developed in the interest of integrating within a single framework all assumptions about the form of the population’s income distribution and the form of the measurement or coverage issues affecting the available data. This parametric framework allows for exploiting several previously explored parametric corrections for high-incomes data issues in making inference at the population level.

Secondly, a Bayesian estimation strategy allows for inference on the population’s income distribution through data presumably affected by representativeness issues on the upper tail. This strategy extends the Approximate Bayesian Computation approach recently explored in [Kobayashi and Kakamu \(2019\)](#) and [Silva \(2023\)](#) in the context of income distributions. In exploiting this approach to estimate income distributions under the proposed parametric framework, the magnitudes and forms of the representativeness issues may be uncertain. Past knowledge on the possible nature of these under similar settings poses information that may be used in dealing with this uncertainty through the use of informative prior beliefs.

Finally, several applications over simulated and household survey data from the European Union’s Statistics on Income and Living Conditions (EU-SILC) illustrate the performance of the proposed approach in controlled and observational settings. These applications evidence the several biases that can hinder making inference on a population’s income distribution if high-income representativeness issues affecting the available data are ignored. Additionally, the presented estimates suggest the presence of both high-income under-reporting and high-income non-response issues in selected EU-SILC samples. This results in population-level estimates of average incomes and inequality that are at higher levels and with higher uncertainty than their sample counterparts.

The following section presents an overview on the common causes and corrections for data errors on high incomes explored in the previous literature. Section 3 develops a parametric framework integrating popular forms of such data errors to parametric income distributions. The fourth section introduces an Approximate Bayesian Computation routine for inference on a population’s income distribution through the proposed parametric framework and under magnitudes and forms for high-income data issues that might be uncertain. Section 5 presents simulated and EU-SILC data applications of the method under typical parametric forms. The sixth and final section of the paper presents concluding remarks with proposals for future work in studying high-income data issues through the proposed approach.

## 2 Dealing with ‘missing rich’ issues

### 2.1 The ‘missing rich’ phenomenon and its sources

In describing the nature of the ‘missing rich’ ( $MR$ ) problem, [Lustig \(2020\)](#) points at the many different issues affecting the upper tail of the observed income distribution in usual data sources. In the context of survey data, the main focus of this paper, one first source of  $MR$  may arise from non-coverage errors in the sampling design itself as a consequence of the sparseness and irreplaceability of high income households. High-income households are generally very few and very dissimilar between themselves such that households on any part of the upper tail of the distribution may have a zero probability of inclusion in the achieved survey sample.

A second possible source for  $MR$  in survey data involves reporting issues either in the form of unit or item non-response (i.e., high-income households refusing to respond to

the survey or particularly to the items concerning their income level, respectively) or in the form of under-reporting of income levels when responding to the survey. Even if the sampling scheme is designed to be representative of the income distribution of the entire population of interest, unit or item non-response may yield an achieved sample which is not and particularly so when this non-response occurs more significantly for households on the upper tail of the income distribution. In a similar way, with income under-reporting the achieved survey sample may yield an income distribution which is not representative of the population's true income distribution if under-reporting is particularly present for high-income households.

Finally, a third possible source can be found within the data provision procedures commonly used by the institutions in charge of distributing publicly-available household survey datasets. In the interest of statistical disclosure control, it is common for such publicly-available datasets to contain a measure of household incomes which is top-coded (i.e., right-censored) meaning that reported incomes above a certain threshold cannot be observed as measured and only an indicator of being above this income threshold is presented.

An additional consideration on the nature of the MR problem concerns the analysis of income concepts which are an aggregate of different income sources. Notably, when working with household survey data it is common to focus the analysis on the distribution of a key income variable which aggregates all income sources surveyed (i.e., labor incomes, capital incomes, social security benefits, etc.). There is, however, vast empirical evidence of differential trends in *MR* issues across the distributions of each income source separately (e.g., see [Moore et al. 2000](#), [Angel et al. 2019](#)). Ultimately, the patterns of *MR* issues that might affect a measure of incomes which aggregates all sources will be the output of the interaction of the different *MR* patterns affecting each income source and the trends in the composition of these aggregated incomes along the distribution.

## 2.2 Income distribution models and treatment of 'missing rich' issues

In modeling income distributions, the use of parametric models is a standard. Some work has fruitfully explored the use of non- or semi-parametric methods for income distribution analysis (e.g., [Jenkins 1995](#)), yet there is vast evidence of parametric models fitting real data on incomes better than these alternatives in many different settings (e.g., [Darvas 2019](#), [Jorda et al. 2021](#)).

The usual modeling step involves assuming that individuals' incomes  $y_i$  are distributed across its population following some three- (i.e.,  $\Theta \subseteq \mathbb{R}^3$ ) or four-parameter (i.e.,  $\Theta \subseteq \mathbb{R}^4$ ) distribution  $y_i \sim f_y(\cdot; \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} \in \Theta$ . Popular choices for  $f_y(\cdot; \boldsymbol{\theta})$  include the Generalized Beta family of distributions (e.g., see [McDonald 1984](#), [Jenkins 2009](#), [citealtgrafnedyalkova2014](#), [Chotikapanich et al. 2018](#), [Jorda and Niño-Zarazúa 2019](#)), in particular the four-parameter Generalized Beta distribution of the second kind (GB2, which is taken as illustrative reference in what follows) and the three-parameter Singh-Maddala (Burr XII) distribution,

and the Double Pareto Log-Normal distribution<sup>1</sup>.

There are several virtues to the parametric approach aside from its generally good fit to real data on incomes. Of particular relevance is its flexibility with respect to the format of data available. Several estimators following parametric expressions for microdata or bracketed/grouped data from incomes following  $y_i \sim f_y(\cdot; \theta)$  are available for most distributions such as Generalized Method-of-Moments (GMM), Maximum Likelihood Estimation (MLE), or Bayesian inference methods.

A central consideration required in analysing data prone to high-income representativeness issues is that the distribution of observed incomes  $y_i^{Obs}$  will very unlikely follow the form of the population's income distribution  $y_i \sim f_y(\cdot; \theta)$ . Within a parametric approach, however, the observed distribution can be derived under assumed parametric forms for the errors affecting the data (e.g., see [Deaton 2005](#)). Jointly modeling the population income distribution component  $f_y(\cdot; \theta)$  and the high-income issues is an attempt at separating which aspects of the data reflect those of the population income distribution and which aspects are due to the high-income problems considered.

Deriving parametric distributions for the sample distribution of incomes  $y_i^{Obs}$  observed under simple forms of measurement errors is the focus of early literature in the field. Models obtained from simple two-parameter distributions 'distorted' through classical measurement errors (i.e., independent of incomes) brought forward implications that would be in strong contrast with recent empirical observations: classical measurement errors can yield sample inequality estimates that overestimate inequality at the population level (e.g., see [Krishnaji 1970](#), [Hartley and Revankar 1974](#), [Hinkley and Revankar 1977](#), [Van Praag et al. 1983](#), [Ransom and Cramer 1983](#), [Chesher and Schluter 2002](#)).

More recent literature, in change, has focused in characterizing under-reporting phenomena affecting income data. The robustly evidenced progressiveness of under-reporting with respect to income levels has resulted in more appropriate *non-classical* parametric expressions for these measurement errors (e.g., see [Gottschalk and Huynh 2010](#), [Bourguignon 2018](#), [Blanchet et al. 2022](#), and [Flachaire et al. 2022](#)) and has consistently found that high-income under-reporting yields sample inequality measures that underestimate inequality at the population level. This recent exploration of high-income under-reporting has given way to what are known as **replacing** corrections: incomes presumed to be under-reported in the data are replaced by imputations from external data sources such as administrative tax data or by imputations from a model for the under-reporting mechanism. The 'corrected' data are then treated as a representative sample of the population's incomes following  $f_y(\cdot; \theta)$ .

---

<sup>1</sup>The GB2 distribution  $GB2(\alpha, \beta, p, q)$  has pdf:

$$y_i \sim f_y^{GB2}(y_i | \alpha, \beta, p, q) = \frac{\alpha y_i^{\alpha p - 1}}{\beta^{\alpha p} B(p, q) \left(1 + \left(\frac{y_i}{\beta}\right)^{\alpha}\right)^{p+q}}, \quad (y_i, \beta, \alpha, p, q) \in \mathbb{R}_+^5$$

with parameters  $\alpha$ ,  $p$ , and  $q$  ruling the shape of the distribution and  $\beta$  ruling the scale and where  $B(p, q)$  denotes the Beta function, defined as  $B(p, q) = \int_0^1 t^{p-1} (1-t)^{q-1} dt$ . See [Chotikapanich et al. \(2018\)](#) for detailed coverage on the use of this model for income distributions.

Another source of *MR* issues, that of missing observations, has also been treated under parametric approaches. The case of item non-response has received significantly more treatment than the more complex case of unit non-response (e.g., see Brunori et al. (2022) for a recent survey). The main aspect determining how to proceed in practice concerns the distinction between observations Missing Completely at Random (MCAR), Missing at Random (MAR), or Missing Not at Random (MNAR), following the works of Rubin (1976), Rubin (1977), and Greenlees et al. (1982). In the MCAR case the probability of a unit/item being missing in the data is independent of any characteristics of the unit and is constant across all units, inducing no particular biases to any distributional estimates from samples affected by this form of missing data. The MAR case allows this probability to change with the characteristics of the unit but requires it to be independent of the unit’s income level. Finally, the more complex MNAR case allows this probability to also change with the unit’s income level and is therefore the only mechanism capable of representing the empirically evidenced negative relationship between response probabilities and income levels in survey data (e.g., Bollinger et al. 2019, Hlasny 2020).

The biases introduced by MAR or MNAR missing data mechanisms in distributional analysis have mostly been treated under the assumption that unit/item missingness is due to non-response. This is, the assumption that conditional on being sampled high-income units are less likely to report their incomes (in the case of item non-response) or any information at all (in the case of unit non-response) than other units. This approach has motivated the use of **reweighting** corrections: the empirical distribution of incomes in the sample is reweighted by the related distribution of (imputed) response probabilities. Like with replacing, the reweighted data is then treated as a representative sample of the population’s incomes following  $f_y(.; \theta)$ .

A particularly lacking aspect of the recent replacing/reweighting approaches in dealing with *MR* is the lack of unified parametric frameworks integrating the modelling assumptions on the income distribution  $f_y(.; \theta)$  and those on the under-reporting and/or non-response mechanisms. This has several consequences on the applicability and generalizability of these methods. As a model for the data directly as it is observed, a unified parametric approach can allow for deriving expressions and estimation strategies suitable for microdata but also for other formats such as bracketed or grouped data. Additionally, this may prove useful in dealing with the challenge that recent approaches in the literature face concerning the choice of correction quantities (i.e., the share of missing and/or under-reported incomes and their distribution). In general these quantities are hand-set by the analyst or are set to match quantities given from more reliable external data.

While setting correction quantities *ad hoc* relies entirely on the analyst’s knowledge about the population’s true income distribution, setting these quantities taking external data as reference poses several issues of its own. Firstly, it is not always the case that more reliable external data sources on incomes are available for research purposes as there may be access restrictions to such data or the data may suffer from *MR* issues of their own such as those induced by tax evasion and tax avoidance on administrative tax data. Secondly, even when external data is available it is generally the case that the population coverage and income components covered differ from those in the primary data available for the analysis and this implies that several harmonizations must be made in transferring quantities from the former to complement the latter. This harmonizations often come at



the cost of forcing different income concepts to represent the same and of a loss in being able to quantify the statistical uncertainty around the resulting distributional estimates and in particular how these are affected by the inherent uncertainty concerning the correction quantities.

The parametric framework proposed in what follows builds on the recent literature exploring replacing and reweighting corrections for *MR* by integrating in a single distribution function both the population income distribution model  $f_y(\cdot; \boldsymbol{\theta})$  and any assumed form for measurement or missing data issues affecting incomes. Several formats of data and inference may be analyzed through the scope of this framework including that of learning about plausible values for the different *MR* correction quantities from the data itself and of integrating the uncertainty around these quantities into distributional estimates such as the Gini coefficient.

### 3 A parametric replacing and reweighting framework

Let individual  $i$ 's *true* income be denoted by  $y_i \sim f_y(\cdot; \boldsymbol{\theta})$ , with probability density function (pdf)  $f_y(\cdot; \boldsymbol{\theta})$  and cumulative distribution function (CDF)  $F_y(\cdot; \boldsymbol{\theta})$  parameterized by the parameter vector  $\boldsymbol{\theta}$ , and consider a sample of *observed* individual incomes from this population  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$ . This sample may be affected by two types of *MR* issues: high-income under-reporting, in which case income  $y_i^{Obs}$  is observed but differs from  $y_i$  following under-reporting of high incomes, and non-response, in which case no income is observed for the individual  $i$ .

To introduce a parametric model for data under both of these possible issues, let  $\varphi(\mathbf{y}, \mathbf{X}; \boldsymbol{\nu})$  a **response probability function** defining the probability for an individual to report her income after being sampled from the population. In its most general formulation this probability, parameterized by the vector  $\boldsymbol{\nu}$ , may depend on the individual's income  $y_i$  and/or other characteristics  $\mathbf{X}_{i\cdot}$ , but also on others' incomes  $\mathbf{y}$  and/or characteristics  $\mathbf{X}$  more generally. This function should allow a representation presenting all typical properties of a univariate probability density function.

Additionally, denote by  $m(\mathbf{y}, \mathbf{X}; \boldsymbol{\eta})$  an **income reporting function**, defining the link between  $i$ 's income  $y_i$ , her characteristics in  $\mathbf{X}$ , and her income reported in the data, if any,  $y_i^{Obs} \equiv m(\mathbf{y}, \mathbf{X}; \boldsymbol{\eta})$  parameterized by the vector  $\boldsymbol{\eta}$ . For any application of empirical relevance defining such a reporting function is only of interest in as much as it allows for an inverse representation taking as input an observed income  $y_i^{Obs}$  from the available data and yielding as output a corresponding income level  $y_i$  or an estimate of this if  $m(\mathbf{y}, \mathbf{X}; \boldsymbol{\eta})$  is non-deterministic. A simple case presenting this property is when the reporting function is deterministic and invertible such that  $y_i \equiv m^{-1}(\mathbf{y}^{Obs}, \mathbf{X}; \boldsymbol{\eta})$  defines a **replacing function**<sup>2</sup>.

---

<sup>2</sup>Formally, this invertibility assumption amounts to assuming the reporting function  $m(y_i, \mathbf{X}; \boldsymbol{\eta})$  to be continuously differentiable with non-zero derivative at all income levels in the population as a sufficient condition of invertibility. Moreover, this assumption also implies:

$$\frac{\partial m^{-1}(y_i^{Obs}, \mathbf{X}; \boldsymbol{\eta})}{\partial y_i^{Obs}} = \left( \frac{\partial m(y_i, \mathbf{X}; \boldsymbol{\eta})}{\partial y_i} \right)^{-1}$$

for all income levels.

In practice, this simplifying assumption may be invalid whenever the income reporting function is non-deterministic such that individuals with equal characteristics and income may randomly report different observed incomes. This assumption may also be invalid whenever individuals with same characteristics but different income levels report a same observed income in a deterministic manner, defining a flat region of incomes where the reporting function is not invertible. These situations could be addressed by devising a vector-output replacing function, yielding several possible income levels as output, in the spirit of Multiple Imputations (e.g., see [Brownstone and Valletta 1996](#)). For simplicity of presentations and derivations, however, the income reporting function  $m(\mathbf{y}, \mathbf{X}; \boldsymbol{\eta})$  is assumed to be a real-output, deterministic, and invertible function in what follows.

Within this framework, we can relate  $i$ 's income to her observed income  $y_i^{Obs}$ , if reported, and to some unobservable income  $y_i^{NObs}$ , in case of non-response, following<sup>3</sup>:

$$y_i = \begin{cases} m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta}) & , \text{ with probability } \varphi(y_i, \mathbf{X}_i; \boldsymbol{\nu}) \\ y_i^{NObs} & , \text{ with probability } 1 - \varphi(y_i, \mathbf{X}_i; \boldsymbol{\nu}) \end{cases} .$$

If no measurement or non-response issues are believed to affect the data, then this amounts to setting  $(m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta}), \varphi(y_i, \mathbf{X}_i; \boldsymbol{\nu})) \equiv (y_i^{Obs}, 1)$  and therefore  $y_i^{Obs} \sim f_{\mathbf{y}}(\cdot; \boldsymbol{\theta})$ .

### 3.1 High-income under-reporting forms.

Whenever some form of measurement error is assumed to affect incomes in the data, then this may be introduced through a specific choice for the replacing function  $m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta})$ . This function serves the purpose of introducing any replacing or imputation step where  $i$ 's income is set as a function of her observed income and characteristics. For simplicity, only forms  $m^{-1}(y_i^{Obs}; \boldsymbol{\eta})$  where measurement errors are exclusively determined by units' income levels are considered in what follows.

Any  $m^{-1}(y_i^{Obs}; \boldsymbol{\eta})$  representing progressive under-reporting of high incomes should imply an increasing and convex quantile ratio  $r(i; \boldsymbol{\eta})$  defined as:

$$r(i; \boldsymbol{\eta}) = \frac{m^{-1}(y_{(i)}^{Obs}; \boldsymbol{\eta})}{y_{(i)}^{Obs}} , \quad \frac{\partial r(i; \boldsymbol{\eta})}{\partial y_{(i)}^{Obs}} \geq 0 , \quad \frac{\partial^2 r(i; \boldsymbol{\eta})}{\partial^2 y_{(i)}^{Obs}} \geq 0 ,$$

with  $y_{(i)}^{Obs}$  denoting the  $i$ -th quantile of  $\mathbf{y}^{Obs}$ . This restricts relative discrepancies between observed  $y_i^{Obs}$  and replaced  $m^{-1}(y_i^{Obs}; \boldsymbol{\eta})$  incomes to be non-decreasing with incomes.

Recently popular replacing approaches can easily be expressed as deterministic forms for  $m^{-1}(y_i^{Obs}; \boldsymbol{\eta})$  including:

---

<sup>3</sup>For simplification reasons, all derivations in what follows are under the assumption that  $i$ 's both response probabilities and reported income depend only on  $i$ 's income and characteristics:  $m^{-1}(\mathbf{y}^{Obs}, \mathbf{X}; \boldsymbol{\nu}) \equiv m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\nu})$ ,  $\varphi(\mathbf{y}, \mathbf{X}; \boldsymbol{\nu}) \equiv \varphi(y_i, \mathbf{X}_i; \boldsymbol{\nu})$



- **Piecewise linear quantile-ratio replacing** (e.g., [Flachaire et al. 2022](#)):

$$m^{-1}(y_{(i)}^{Obs}, \{\bar{p}_k\}_{k=1}^{K-1}, \{\beta_k\}_{k=1}^{K-1}, \{\delta_k\}_{k=1}^{K-1}) \equiv \begin{cases} y_{(i)}^{Obs}, & \text{if } p_{(i)} \leq \bar{p}_1 \\ y_{(i)}^{Obs} \times \sum_{k=1}^{K-1} \mathbf{1}(\bar{p}_k < p_{(i)} \leq \bar{p}_{k+1}) \times \underbrace{(\beta_k + \delta_k p_{(i)})}_{\text{Linear replacing weights for incomes in the } k\text{-th segment.}}, & \end{cases},$$

with  $\delta_j \leq \delta_{j+1} < \infty$ ,  $j = 1, \dots, K-2$ ,  $\bar{p}_K = 1$ , and with  $p_{(i)} = F_{\mathbf{y}}(y_{(i)}; \boldsymbol{\theta})$  denoting ordered-incomes individual  $y_{(i)}^{Obs}$ 's percentile in the population's income distribution<sup>4</sup>. In absence of missing data, the sample percentile  $p_{(i)}^{Obs} \equiv \frac{(i)}{N}$ ,  $(i) = 1, \dots, N$  is also a valid estimate of  $F_{\mathbf{y}}(y_{(i)}; \boldsymbol{\theta})$ . The central assumption under this approach is that the quantile ratio  $r(i; \boldsymbol{\eta}) = \frac{m^{-1}(y_{(i)}^{Obs}; \{\bar{p}_k\}_{k=1}^{K-1}, \{\beta_k\}_{k=1}^{K-1}, \{\delta_k\}_{k=1}^{K-1})}{y_{(i)}^{Obs}}$  can be represented as a continuous piecewise linear function. This piecewise representation allows for progressive under-reporting of high incomes across segments  $(\bar{p}_k, \bar{p}_{k+1}]$  of the income distribution at the cost of introducing 3 additional parameters  $(\bar{p}_k; \beta_k; \delta_k)$  per segment.

- **Linear progressive under-reporting (LPU)**, [Bourguignon 2018](#)):

$$m^{-1}(y_i^{Obs}; \bar{p}, \delta) \equiv y_i^{Obs} + \underbrace{\mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta}))}_{\text{Indiv. with observed incomes above } \bar{p}\text{-th percentile under-report}} \times \underbrace{\left( \frac{\delta(y_i^{Obs} - F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta}))}{1 - \delta} \right)}_{\text{Under-reported amount linearly increases with true incomes with slope } \delta},$$

with  $\delta \in [0, 1)$  and with  $F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta})$  denoting the  $\bar{p}$ -th population income quantile. This replacing scheme assumes that all individuals with incomes above the  $\bar{p}$ -th percentile under-report their incomes in the observed sample and do so in a linearly progressive manner with under-reporting increasing by  $\delta$  with every additional unit of income. The incomes quantile ratio implied under LPU  $r(i; \boldsymbol{\eta}) = \frac{m^{-1}(y_{(i)}^{Obs}; \bar{p}, \delta)}{y_{(i)}^{Obs}}$  is strictly convex for income levels above  $F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta})$ .

- **Generalized Pareto replacing:** (e.g., [Atkinson and Piketty 2007](#), Chapter 2,, [Jenkins 2017](#), [Hlasny and Verme 2022](#), [Charpentier and Flachaire 2022](#)):

$$m^{-1}(y_i^{Obs}; \mu, \sigma, \zeta) \equiv y_i^{Obs} + \underbrace{\mathbf{1}(y_i^{Obs} > \mu)}_{\text{Indiv. with observed incomes above } \mu \text{ under-report}} \times \underbrace{\left[ \left( \frac{\left( 1 - \left( \frac{p_i - \bar{p}}{1 - \bar{p}} \right) \right)^{-\zeta} - 1}{\zeta} \right) \times \sigma - (y_i^{Obs} - \mu) \right]}_{\text{Observed incomes are replaced by corresponding percentile under a Generalized Pareto dist.}},$$

where  $(\mu, \sigma, \zeta)$  are respectively the location, scale, and shape parameters of a Generalized Pareto distribution  $GPD(\mu, \sigma, \zeta)$  with CDF given by ([Pickands, 1975](#)):

$$F_{\mathbf{y}}(y_i; \mu, \sigma, \zeta) = \begin{cases} 1 - \left( 1 + \frac{\zeta(y_i - \mu)}{\sigma} \right)^{-\frac{1}{\zeta}}, & \text{if } \zeta \neq 0 \\ 1 - e^{-\left( \frac{y_i - \mu}{\sigma} \right)}, & \text{if } \zeta = 0 \end{cases}, \quad y_i > \mu.$$

<sup>4</sup>In what follows  $\mathbf{1}(\cdot)$  represents the indicator function, taking value 1 whenever the condition it takes as argument holds true and 0 otherwise.

Under this replacing scheme, any income above a level  $\mu$  is assumed to be under-reported. True incomes are assumed to follow a Generalized Pareto distribution or some specific case such as the Pareto I, corresponding to  $GPD(\frac{\sigma}{\zeta}, \sigma, \zeta)$ ,  $\zeta > 0$ , or Pareto II, corresponding to  $GPD(\mu, \sigma, \zeta)$ ,  $\zeta > 0$  above this income level. In absence of missing data issues an individual in the  $p_i$ -th sample percentile with  $p_i > \bar{p}$  (equivalently,  $y_i^{Obs} > \mu$ ) has true income corresponding to the  $\left(\frac{p_i - \bar{p}}{1 - \bar{p}}\right)$ -th quantile on this Pareto distribution. Similarly to LPU, the incomes quantile ratio implied under Generalized Pareto replacing  $r(i; \boldsymbol{\eta}) = \frac{m^{-1}(y_i^{Obs}; \mu, \sigma, \zeta)}{y_i^{Obs}}$  is strictly convex for income levels above  $\mu$ , representing progressiveness of the under-reporting.

These common replacing schemes all exploit the assumption that under-reporting is a deterministic function of individual incomes (or their sample percentile/rank equivalently), and that individuals have the same rank in the population's income distribution as in the observed sample. It's also important to note that each specific replacing scheme implies within it specific assumptions on under-reporting behavior at the individual level.

Figure 1 below illustrates a comparative example of the quantile ratios  $r(i, \boldsymbol{\eta})$  under these three common forms for  $m^{-1}(\cdot; \boldsymbol{\eta})$ . The respective parameter values  $\boldsymbol{\eta}$  are set to represent a same progressive under-reporting pattern: LPU affecting observed incomes from the .75-th percentile of the population income distribution upwards and with slope of  $\delta = .67$ . A first observation illustrated in this figure is that a piecewise linear approximation to the considered LPU under-reporting pattern, introducing six parameters in  $\boldsymbol{\eta}$  in total (i.e., a linear approximation with two segments), is not flexible enough to correctly represent it. Secondly, replacing under a Generalized Pareto tail all incomes above the .75-th percentile can represent the reference LPU pattern accurately except for the top of the income distribution. For the highest incomes the differences between observed and replaced incomes under this form can be particularly large as a consequence of the heavy Pareto tail used for the purpose of replacing. Finally, this similarity across GPD and LPU schemes suggests the latter as the more stable and parsimonious alternative of the two.

Further choices for  $m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta})$  may exploit other individual characteristics  $\mathbf{X}_i$  to represent larger heterogeneities in under-reporting patterns. A popular choice when data on individual consumption is available without reporting errors of its own is to define the income reporting function as an Engel curve (e.g., see [Pissarides and Weber 1989](#), [Lyssiotou et al. 2004](#), [Hurst et al. 2014](#)). Additionally, a stochastic component may be introduced in the definition of  $m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta})$  to allow for heterogeneity in under-reporting behavior across individuals with same level of incomes (e.g., [Flachaire et al. 2022](#)).

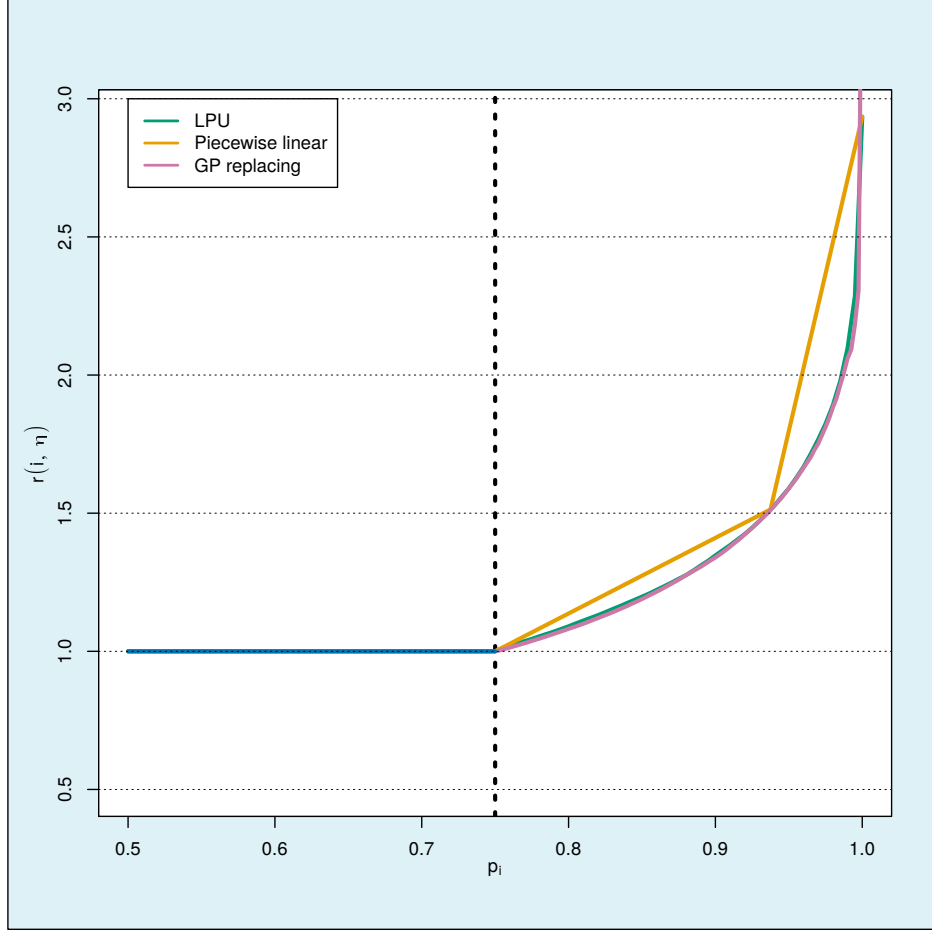


Figure 1: Quantile ratios under common replacing schemes

**Note:** Three common replacing schemes: LPU, piecewise linear quantile ratio, and Generalized Pareto replacing (GP in legend), as applied to a same income distribution following  $y_i \sim GB2(2.257, 17393, 1, 1.033)$  and affected by LPU with  $\bar{p} = .75$  (represented by the dashed vertical line) and  $\delta = .67$ . The piecewise linear approximation was calibrated to fit this LPU pattern at the  $\bar{p}_1 = .75$  and  $\bar{p}_2 = .9375$  sample percentiles. The GPD coefficients  $\zeta$  and  $\sigma$  were estimated conditional on  $\mu$  being the .995-th sample quantile as a typical empirical practice (e.g., see [Jenkins 2017](#)) and imposing finite variance ( $\zeta < \frac{1}{2}$ ).

### 3.2 High-income non-response forms.

Concerning the modeling assumptions for the response probabilities  $\varphi(\mathbf{y}_i, \mathbf{X}_i; \boldsymbol{\nu})$ , the several possible types of missing data mechanisms may be considered, following [Rubin \(1976\)](#). If income non-response follows a random process which is unrelated to incomes  $\mathbf{y}$  and other characteristics  $\mathbf{X}$ , then the mechanism corresponds to a MCAR process. In the MCAR case, we observe incomes  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$  which are a random sample of the population's incomes and therefore no particular bias is induced by the missing data. A simple MCAR mechanism is such that  $\varphi(\mathbf{y}_i, \mathbf{X}_i; \boldsymbol{\nu}) \equiv \varphi(\mathbf{y}_i, \mathbf{X}_i; p) \equiv p$ ,  $p \in (0, 1]$ , where all individuals are just as likely to report incomes after they have been sampled.

A second potential mechanism concerns the case where non-response in incomes is not completely at random but where the missingness can be fully explained by other non-missing characteristics of the individuals and/or by the observed incomes, i.e.,  $\varphi(\mathbf{y}_i, \mathbf{X}_i; \boldsymbol{\nu}) \equiv \varphi(\mathbf{y}_i^{Obs}, \mathbf{X}_i; \boldsymbol{\nu})$ . This mechanism represents an MAR process and is an appropriate representation for scenarios of item non-response, where sampled

individuals report information about their characteristics  $\mathbf{X}_i$  but not about their income, as long as their unobserved income  $y_i^{Nobs}$  is unnecessary to account for the non-random non-response probabilities. MAR mechanisms are usually dealt with in analysis through multiple imputations of incomes for those individuals in the data with missing incomes but observed characteristics.

Finally, it may be the case that response probabilities may not be fully accounted for from observed information. For instance, it may be the case that the reason why individuals do not report their incomes in the data has everything to do with their unobserved level of incomes  $y_i^{Nobs}$ . This corresponds to the MNAR scenario and is particularly complex to deal with, as it may include non-random unit non-response mechanisms, where sampled individuals do not report neither incomes nor characteristics and where their unobserved incomes  $y_i^{Nobs}$  are a determinant of this.

Forms for  $\varphi(\mathbf{y}_i, \mathbf{X}_i; \boldsymbol{\eta})$  suitable for MAR mechanisms have been the focus of the recent survey in Brunori et al. (2022). Recently popular reweighting approaches allowing for dealing also with MNAR mechanisms can easily be expressed as deterministic forms for  $\varphi(\mathbf{y}_i, \mathbf{X}_i; \boldsymbol{\eta})$  including:

- **Right-truncation** (e.g., Alvaredo 2011, Jorda and Niño-Zarazúa 2019):

$$\varphi(y_i; t, \alpha) \equiv \begin{cases} \alpha, & \text{if } p_{(i)} \leq t \\ 0, & \text{if } p_{(i)} > t \end{cases},$$

which amounts to assuming that any and all individuals above the  $t$ -th percentile on the population income distribution will not report incomes in the data, while anyone below this threshold will report an income with probability  $\alpha$ . The limiting case  $\alpha \rightarrow 1$  corresponds to assuming that any unit with income below the  $t$ -th percentiles will always report an income when sampled.

- **Regional non-response reweighting** (e.g., Korinek et al. 2007, Hlasny and Verme 2018):

$$\varphi(y_i, \mathbf{X}_i; \boldsymbol{\beta}) \equiv \frac{e^{g(y_i, \mathbf{X}_i; \boldsymbol{\beta})}}{1 + e^{g(y_i, \mathbf{X}_i; \boldsymbol{\beta})}},$$

with  $g(y_i, \mathbf{X}_i; \boldsymbol{\beta})$  being a twice continuously differentiable function of observed unit  $i$ 's characteristics parameterized by the vector  $\boldsymbol{\beta}$ . The comparative analysis in Hlasny and Verme (2015) suggests a simple logarithmic form for  $g$  taking as input a linear combination of income  $y_i$  and region indicator variables to be equally efficient as more complex specifications in many scenarios. This approach infers response probabilities for units from modeling the relationship between non-response rates and units' characteristics at aggregate (i.e., regional) levels, when this information is available. The key assumption is that individual characteristics relate to individual response probabilities in the same way that they do at the aggregate level (i.e., that ecological inference is feasible), such that individual response probabilities  $\varphi(y_i, \mathbf{X}_i; \boldsymbol{\beta})$  may be properly estimated and used for the purpose of reweighting the observed data.

- **Income-proportional reweighting** (e.g., Blanchet et al. 2022):

$$\varphi(y_i; \gamma_0, \gamma_1, t, \alpha) = \begin{cases} e^{\gamma_0(y_i)^{-\gamma_1}}, & \gamma_1 > 0, \text{ if } p_{(i)} > t \\ \alpha, & \text{if } p_{(i)} \leq t \end{cases}.$$

This scheme corresponds to a non-response mechanism where individuals with true incomes above the  $t$ -th percentile have increasingly lower response probabilities, with the parameter  $\gamma_1$  representing the income elasticity of non-response (i.e., how much response probabilities decrease with an increase in incomes of 1%). For a given value of such elasticity,  $\gamma_0$  serves as an intercept to assure the continuity of  $\varphi(y_i; \gamma_0, \gamma_1, t, \alpha)$  at  $t$ . Similarly to right-truncation, the parameter  $\alpha$  represents the (constant) response probability for units with incomes below the  $t$ -th percentile. Moreover, this reweighting scheme includes the right-truncation  $\varphi(y_i; t, \alpha)$  as the limiting case  $\gamma_1 \rightarrow \infty$ . Figure 2 provides an illustrative example of how these two cases relate and their resulting contrasts with the population's income distribution.

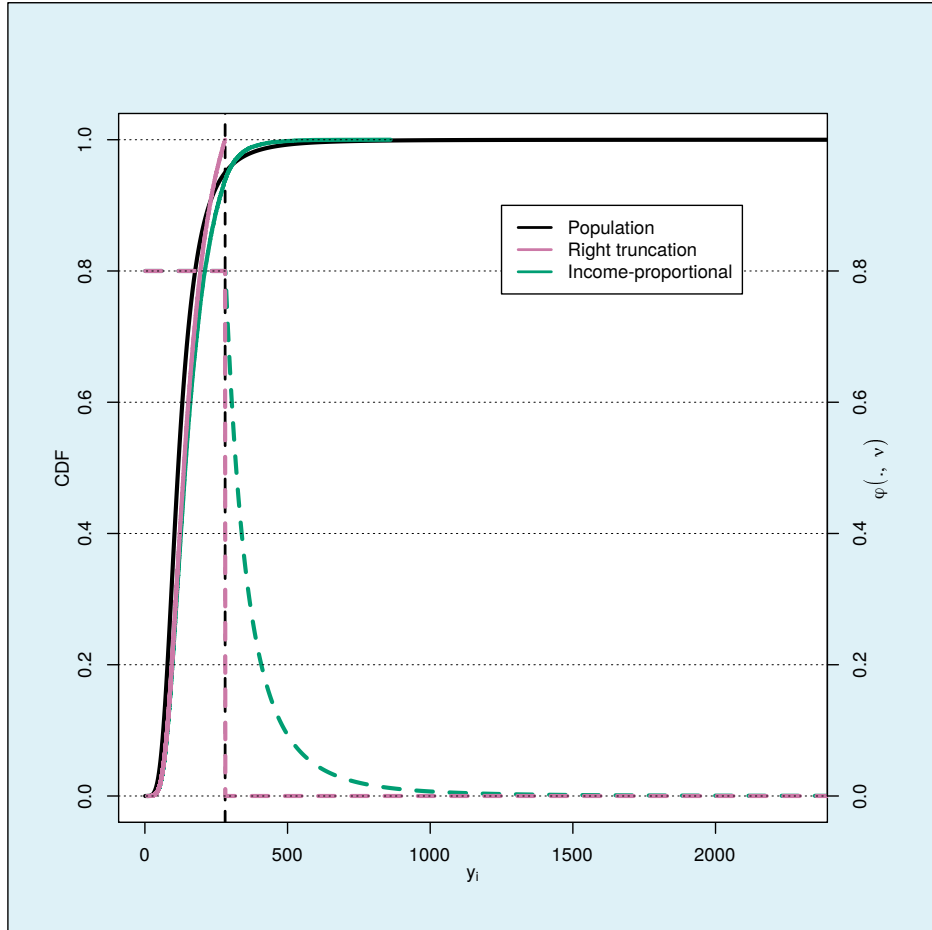


Figure 2: Income-proportional reweighting schemes

**Note:** A population income distribution following  $y_i \sim GB2(2.75, 100, 1.75, 1.25)$  and its corresponding CDF under two cases of income-proportional non-response schemes: Right truncation, with parameter values set at  $(\alpha, t) = (.8, .95)$ , and income-proportional with parameter values  $(\gamma_1, \alpha, t) = (3.75, .8, .95)$  which requires  $\gamma_0 = 20.92$  for continuity. Solid lines represent respective CDFs, on the left axis, and dashed lines represent response probabilities, on the right axis. The dashed vertical line represents the  $t$ -th population percentile.

### 3.3 Jointly accounting for high-income under-reporting and non-response.

Both replacing and reweighting corrections may interact within this framework. Response probabilities are modeled as a function of true incomes  $\mathbf{y}$ , yet these may differ from

observed incomes  $\mathbf{y}^{Obs}$  following an assumed form for  $m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta})$ . The interaction allowing for modeling MNAR non-response mechanisms through observed incomes directly simply amounts to the composite function  $\varphi(m^{-1}(y_i^{Obs}, \mathbf{X}_i; \boldsymbol{\eta}), \mathbf{X}_i; \boldsymbol{\nu})$ .

A crucial question this parametric framework allows to answer is: given i) an assumed form for the population income distribution  $f_{\mathbf{y}}(.; \boldsymbol{\theta})$ , ii) an assumed income reporting form  $m(y_i, \mathbf{X}_i; \boldsymbol{\eta})$ , and iii) an assumed response probability function  $\varphi(y_i, \mathbf{X}_i; \boldsymbol{\nu})$ , then what distribution  $f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  will observed incomes under this framework follow? This distribution can be derived applying the deterministic transformation  $m^{-1}(.; \boldsymbol{\eta})$  to  $f_{\mathbf{y}}(.; \boldsymbol{\theta})$  and reweighting the resulting density by the response probabilities  $\varphi(.; \boldsymbol{\nu})$ , yielding the relationship<sup>5</sup>:

$$f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}) = \frac{\overbrace{f_{\mathbf{y}}(m^{-1}(y_i^{Obs}; \boldsymbol{\eta}); \boldsymbol{\theta}) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \boldsymbol{\eta})}{\partial y_i^{Obs}} \right)}^{\text{Reporting function: Replacing transformation of } \mathbf{y}} \times \overbrace{\varphi(m^{-1}(y_i^{Obs}; \boldsymbol{\eta}); \boldsymbol{\nu})}^{\text{Non-response: Reweighting of } f_{\mathbf{y}}}}{\underbrace{\int f_{\mathbf{y}}(m^{-1}(y_i^{Obs}; \boldsymbol{\eta}); \boldsymbol{\theta}) \times \varphi(m^{-1}(y_i^{Obs}; \boldsymbol{\eta}); \boldsymbol{\nu}) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \boldsymbol{\eta})}{\partial y_i^{Obs}} \right) dy^{Obs}}_{\text{Normalizing constant}}} \quad (1)$$

The main application of the result in (1) is that of parametrically integrating all assumptions about the population income distribution and the *MR* issues affecting the data in a model for the observed data itself. In doing so, this model for the data constitutes a case of continuous model expansion to accommodate for non-response and/or measurement errors (e.g., see [Nandram and Choi 2002](#), [Gustafson 2005](#), [Gelman et al. 2013](#), Chapter 7). In particular, the parametric income distribution model  $f_{\mathbf{y}}(.; \boldsymbol{\theta})$  is expanded by the inclusion of parametric *MR* forms defining a continuum of possible corrections for under-reporting or non-response and including the specific case where no such issues affect the observed data. Fitting such a model to data is an attempt at separately identifying characteristics of the population income distribution, captured by the  $\boldsymbol{\theta}$  vector, and features representing the *MR* forms presumed to affect the data, captured by the  $\boldsymbol{\eta}$  and  $\boldsymbol{\nu}$  vectors.

There are several virtues to integrating the replacing and/or reweighting corrections considered relevant into a model to be taken to the data *as-is*. Firstly, because the correction quantities are completely defined through  $\boldsymbol{\eta}$  and/or  $\boldsymbol{\nu}$  this approach guarantees that all corrections are done on the income concept and population being analyzed. This avoids the issue of manipulating these concepts to be compatible with correction quantities defined in terms of different income concepts or population. On a related note, if external data informative on the forms and magnitudes of *MR* are available then these should be introduced by specifying adequate representations  $m^{-1}(y_i^{Obs}; \boldsymbol{\eta})$  and  $\varphi(y_i^{Obs}; \boldsymbol{\nu})$  and setting  $\boldsymbol{\eta}$  and  $\boldsymbol{\nu}$  to quantify these magnitudes.

A second virtue of this integrated approach is that uncertainty on the parameter values or estimates of these may be translated into uncertainty on distributional statistics like the Gini coefficient. Importantly, this can produce estimates of the population's

---

<sup>5</sup>For simplicity and without loss of generality, only non-response forms of the type  $\varphi(.; \boldsymbol{\eta}) \equiv \varphi(y_i; \boldsymbol{\eta})$  are considered in what follows.



income distribution which account for the uncertainty surrounding the possible  $MR$  issues affecting the available data. Additionally, from a same set of parameter estimates both corrected and non-corrected distributional estimates may be computed in the interest of gaining insight on the impact of the corrections considered in terms of the estimated population's income distribution<sup>6</sup>.

Thirdly, the 'building blocks' nature of the framework allows for exploring several candidate forms for replacing and/or reweighting corrections leaving other components unchanged in a straightforward manner. In particular, this allows for studying the robustness of the estimated  $\theta$  to different assumptions on the form of  $MR$  affecting the data.

Finally, stating the model as a properly defined parametric distribution implies that all observed units are re-weighted under any assumed form for  $\phi(.;\nu)$ , either directly through the reweighting of units prone to non-response in the numerator of (1) (i.e., through down-weighting the density at their respective level of income with respect to the population density) or indirectly through the correction for missing observations in the normalization constant of (1). This 'indirect' reweighting accommodates for the fact that if some units are under-represented in the data due to higher non-response probabilities then necessarily the rest of units are over-represented and therefore need to be reweighted under any correction for these non-response probabilities.

Making inference about the features of the population income distribution and the  $MR$  aspects of the data simultaneously poses several challenges. Issues of identifiability, in particular, require attention as a given model specified following (1) might fit equally well a sample of observed incomes for very different values of the  $\theta$ ,  $\eta$ , and  $\nu$  parameters, making inference on them invalid. The type of continuous model expansion underlying (1) to introduce uncertainty about the specific form and magnitudes of the  $MR$  issues affecting the data falls in line with previous empirical strategies within Bayesian inference (e.g., see Nandram and Choi 2002). The use of prior probabilities on parameter values under a Bayesian approach can overcome some identifiability issues. The following section details a Bayesian inference approach for this purpose which can exploit external information on the  $MR$  correction quantities in dealing with this.

## 4 Parameter inference under 'missing rich'

### 4.1 A Bayesian inference approach

Under the framework developed in the previous section, inference on the population's income distribution  $f_y(.;\theta)$  is made through inference on the values of the  $\theta$  vector given the sample of observed incomes  $y^{Obs}$ . This task is considerably less complex whenever the correction quantities  $\eta$  and  $\nu$  are given fixed values. However, it is rarely the case

---

<sup>6</sup>Corrected distributional estimates correspond to distributional estimates computed from estimates for  $\theta$  only and under  $\nu$  and/or  $\eta$  fixed to correspond to the scenario in absence of  $MR$  issues. These correspond to the estimates of the distribution of true incomes at the population level. Non-corrected distributional estimates correspond to those computed from estimates for all the model's parameters in  $\theta$ ,  $\nu$ , and/or  $\eta$ . These correspond to the fitted income distribution representing the observed data on incomes.

that sound candidate values for these quantities are available. The central challenge in learning about the population's income distribution through  $\theta$  is therefore to exploit the framework under an empirical strategy that can properly make inference on these parameters but also on  $\eta$  and  $\nu$  at the same time.

To make such inference, the task is to learn about which values for the parameters  $\theta \in \Theta_\theta \subseteq \mathbb{R}^{\dim(\theta)}$ ,  $\eta \in \Theta_\eta \subseteq \mathbb{R}^{\dim(\eta)}$ , and  $\nu \in \Theta_\nu \subseteq \mathbb{R}^{\dim(\nu)}$  are more likely to have generated the observed data  $\mathbf{y}^{Obs}$  than others within some region of possible values  $\Theta \equiv \Theta_\theta \times \Theta_\eta \times \Theta_\nu$ . In the Bayesian framework, this information takes the form of a posterior probability distribution  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  defined by two main components under Bayes' theorem. Firstly, all prior beliefs about the values of the  $(\theta, \eta, \nu)$  parameters must be elicited through a prior probability distribution  $p(\theta, \eta, \nu)$  over  $\Theta$ . Secondly, for any fixed value for the parameters  $(\tilde{\theta}, \tilde{\eta}, \tilde{\nu})$  the model's likelihood  $L(\mathbf{y}^{Obs} | \tilde{\theta}, \tilde{\eta}, \tilde{\nu})$  quantifies how likely the observed data  $\mathbf{y}^{Obs}$  is to have been generated from  $f_{\mathbf{y}^{Obs}}(\cdot; \tilde{\theta}, \tilde{\eta}, \tilde{\nu})$ <sup>7</sup>.  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  is then a probability distribution proportional to the prior probability distribution *updated* (or reweighted, equivalently) by the likelihood function:

$$\pi(\theta, \eta, \nu | \mathbf{y}^{Obs}) \propto L(\mathbf{y}^{Obs} | \theta, \eta, \nu) \times p(\theta, \eta, \nu). \quad (2)$$

As an evidence-weighted conversion of prior beliefs, the information contained in the  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  posterior distribution can be interpreted as all remaining uncertainty on the values of  $(\theta, \eta, \nu)$  after having 'learnt' from the data through the likelihood  $L(\mathbf{y}^{Obs} | \theta, \eta, \nu)$ . Whenever the data are informative about these parameters, the posterior distribution reflects less uncertainty around their values than that in  $p(\theta, \eta, \nu)$ .

Estimating a posterior distribution  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  for the model parameters faces several complexities. As is usual in most Bayesian inference settings, it is rarely the case that  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  admits a known form given a model  $L(\mathbf{y}^{Obs} | \theta, \eta, \nu)$  and a prior  $p(\theta, \eta, \nu)$ . This is typically circumvented by studying the posterior distribution through samples generated to converge to  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  under the Monte Carlo principle<sup>8</sup> or the Markov Chain Monte Carlo (MCMC) extension of this principle (e.g., see Gelman et al. 2013, Chapter 11).

A second complexity in estimating  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  concerns the possible 'non-identifiability' of at least some of the parameters in  $(\theta, \eta, \nu)$ . As an illustrative example of this issue, consider a model specified following (1) with a parameter  $\lambda_\theta \in \theta$  ruling the right tail of the  $f_{\mathbf{y}}(\cdot; \theta)$  income distribution and a replacing correction  $m^{-1}(\cdot; \eta)$  with parameter  $\lambda_\eta \in \eta$  also affecting only the right tail. It can be the case that a same sample of incomes  $\mathbf{y}^{Obs}$  may be equally well fit under two different parameter values  $(\tilde{\theta}, \tilde{\eta}, \tilde{\nu}) \in \Theta$  and  $(\tilde{\theta}', \tilde{\eta}', \tilde{\nu}') \in \Theta$  including  $(\tilde{\lambda}_\theta, \tilde{\lambda}_\eta)$  and  $(\tilde{\lambda}'_\theta, \tilde{\lambda}'_\eta)$  respectively. This can render the model incapable of separately identifying variations in high incomes in  $\mathbf{y}^{Obs}$  that would occur with changes in  $\lambda_\theta$  and those due to  $\lambda_\eta$ .

<sup>7</sup>For example, in the case of  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$  being  $N$  independent observations their joint likelihood follows  $L(\mathbf{y}^{Obs} | \theta, \eta, \nu) = \prod_{i=1}^N f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \theta, \eta, \nu)$

<sup>8</sup>The Monte Carlo principle states that any quantity of  $\pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$  which can be expressed as an expectation can be studied through a sufficiently large sample of  $J$  independent draws  $\{(\tilde{\theta}^{(j)}, \tilde{\eta}^{(j)}, \tilde{\nu}^{(j)})\}_{j=1}^J$  from this distribution  $(\tilde{\theta}^{(j)}, \tilde{\eta}^{(j)}, \tilde{\nu}^{(j)}) \sim \pi(\theta, \eta, \nu | \mathbf{y}^{Obs})$

If  $\mathbf{y}^{Obs}$  is not informative about differences in the respective likelihoods  $L(\mathbf{y}^{Obs}|\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})$  and  $L(\mathbf{y}^{Obs}|\tilde{\boldsymbol{\theta}}', \tilde{\boldsymbol{\eta}}', \tilde{\boldsymbol{\nu}}')$ , then prior beliefs on these values will not be updated. The respective posterior probabilities  $\pi(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}|\mathbf{y}^{Obs})$  and  $\pi(\tilde{\boldsymbol{\theta}}', \tilde{\boldsymbol{\eta}}', \tilde{\boldsymbol{\nu}}'|\mathbf{y}^{Obs})$  will therefore be dominated entirely by differences in prior beliefs  $p(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})$  and  $p(\tilde{\boldsymbol{\theta}}', \tilde{\boldsymbol{\eta}}', \tilde{\boldsymbol{\nu}}')$ . If available, external information about the plausibility of  $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})$  and  $(\tilde{\boldsymbol{\theta}}', \tilde{\boldsymbol{\eta}}', \tilde{\boldsymbol{\nu}}')$  may be exploited to set an informative prior distribution  $p(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  giving a lower prior probability to the one set of parameter values less compatible with this external data among the two. Informative priors are a way of exploiting prior knowledge to justify differences in posterior densities for parameter values where  $\mathbf{y}^{Obs}$  is uninformative through the model  $L(\mathbf{y}^{Obs}|\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$ .

Returning to the illustrative example above, consider an application of (1) as a model for a survey's sample on incomes  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$  integrating a GB2 income distribution  $f_y^{GB2}(\cdot; \boldsymbol{\theta})$ ,  $\boldsymbol{\theta} = (\alpha, \beta, p, q)$  with a LPU form for  $m^{-1}(\cdot; \boldsymbol{\eta})$ ,  $\boldsymbol{\eta} = (\bar{p}, \delta)$ . LPU affects only the tail above the  $\bar{p}$ -th percentile of the income distribution, while the  $p$  and  $q$  parameters of the GB2 distribution rule its right tail. This allows for identifiability issues as described above, as there might be configurations 'trading' values of  $p$  and  $q$  with values of  $\bar{p}$  and  $\delta$  while representing two observably identical income distributions. In this example, external information might be introduced in the form of prior probabilities by setting the marginal prior distributions for  $\bar{p}$  and  $\delta$  around previous empirical findings on *MR* issues in similar settings<sup>9</sup>.

Several sampling algorithms can be devised to obtain samples  $\{(\tilde{\boldsymbol{\theta}}^{(j)}, \tilde{\boldsymbol{\eta}}^{(j)}, \tilde{\boldsymbol{\nu}}^{(j)})\}_{j=1}^J$  from  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}|\mathbf{y}^{Obs})$  under a model following (1) and an informative prior  $p(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$ . The Metropolis-Hastings (MH) algorithm defines a type of MCMC sampler suitable for estimating parametric income distribution models in several contexts (e.g., see [Chotikapanich and Griffiths 2000](#), [Peters and Sisson 2006](#), [Chotikapanich and Griffiths 2008](#)). A standard MH sampler for the joint parameter vector  $\boldsymbol{\phi} = (\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  is possible following algorithm 1 below. Such an MH algorithm yields as output a sample  $\{\tilde{\boldsymbol{\theta}}^{(j)}, \tilde{\boldsymbol{\eta}}^{(j)}, \tilde{\boldsymbol{\nu}}^{(j)}\}_{j=1}^J$  resulting from a global exploration of the support of  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}|\mathbf{y}^{Obs})$  through local accept-reject steps. Any  $j$ -th,  $j = 1, \dots, J$ , local accept-reject step is defined by the MH acceptance probability:

$$\rho^{(j)} = \min \left\{ 1, \frac{\pi(\tilde{\boldsymbol{\phi}}^{(j)}|\mathbf{y}^{Obs}) \times g(\tilde{\boldsymbol{\phi}}^{(j-1)}, \tilde{\boldsymbol{\phi}}^{(j)})}{\pi(\tilde{\boldsymbol{\phi}}^{(j-1)}|\mathbf{y}^{Obs}) \times g(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\boldsymbol{\phi}}^{(j-1)})} \right\}, \quad \boldsymbol{\phi} = (\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}),$$

with  $g(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\boldsymbol{\phi}}^{(j-1)})$  denoting a candidate function from which the  $j$ -th candidate value  $\tilde{\boldsymbol{\phi}}^{(j)}$  is sampled, given the previously retained value  $\tilde{\boldsymbol{\phi}}^{(j-1)}$ .

---

<sup>9</sup>For example, in their study comparing household survey incomes to linked tax return data for Uruguay [Flachaire et al. \(2022\)](#) find evidence of progressive under-reporting potentially affecting the survey data above  $\bar{p} = .50$ . In studying similar linked data for the Austrian case, [Angel et al. \(2019\)](#) find evidence of progressive under-reporting of wages potentially affecting their survey above the  $\bar{p} = .50$  percentile. The degree of progresiveness of under-reporting can be quantified in terms of  $\delta$  under a linear approximation to the observed under-reporting patterns.

---

**Algorithm 1:** A Metropolis-Hastings algorithm (*MH*).

---

**Initialization:**

**Until**  $L(\mathbf{y}^{Obs}|\boldsymbol{\phi}^{(0)}) > 0$ :

- 1: Sample  $\tilde{\boldsymbol{\phi}}^{(0)}$  from  $p(\boldsymbol{\phi})$

**Sampling:**

**for**  $j = 1, \dots, J$  **do**

- 2: Sample  $\tilde{\boldsymbol{\phi}}^{(j)} \sim g(\boldsymbol{\phi}, \tilde{\boldsymbol{\phi}}^{(j-1)})$  from the candidate  $g$
- 3: Accept and store  $\tilde{\boldsymbol{\phi}}^{(j)}$  with probability:

$$\rho^{(j)} = \min \left\{ 1, \frac{\overbrace{L(\mathbf{y}^{Obs}|\tilde{\boldsymbol{\phi}}^{(j)}) \times p(\tilde{\boldsymbol{\phi}}^{(j)}) \times g(\tilde{\boldsymbol{\phi}}^{(j-1)}, \tilde{\boldsymbol{\phi}}^{(j)})}^{\propto \pi(\tilde{\boldsymbol{\phi}}^{(j)}|\mathbf{y}^{Obs}) \text{ under (2)}}}{L(\mathbf{y}^{Obs}|\tilde{\boldsymbol{\phi}}^{(j-1)}) \times p(\tilde{\boldsymbol{\phi}}^{(j-1)}) \times g(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\boldsymbol{\phi}}^{(j-1)})} \right\}$$

▷ e.g., if  $u^{(j)} \leq \rho^{(j)}$  where  $u^{(j)}$  is a draw from a  $Uniform(0, 1)$  distribution  
otherwise store  $\tilde{\boldsymbol{\phi}}^{(j)} = \tilde{\boldsymbol{\phi}}^{(j-1)}$   
**end**

---

A common choice of candidate function is that of the Adaptive Random-Walk Metropolis (*AM*) algorithm (Haario et al., 2001). In this case the proposal  $g \equiv g_\Sigma$  is defined by the following adaptive random walk process:

$$\begin{aligned} (\boldsymbol{\theta}^{(j)}, \boldsymbol{\eta}^{(j)}, \boldsymbol{\nu}^{(j)}) &\equiv \boldsymbol{\phi}^{(j)} \sim g_{\Sigma^{(j-1)}}(\boldsymbol{\phi}, \tilde{\boldsymbol{\phi}}^{(j-1)}) \Rightarrow \tilde{\boldsymbol{\phi}}^{(j)} = \tilde{\boldsymbol{\phi}}^{(j-1)} + \boldsymbol{\epsilon}^{(j)} \\ \boldsymbol{\epsilon}^{(j)} &\sim N_d(0, \Sigma^{(j-1)}) \\ \Sigma^{(j-1)} &= \begin{cases} \Sigma^{(0)}, & \text{if } j \leq J_0 \\ s_d \times \frac{1}{(j-1)} \left( \sum_{i=1}^{(j-1)} \tilde{\boldsymbol{\phi}}^{(i)} \tilde{\boldsymbol{\phi}}^{(i)'} - i \times \bar{\boldsymbol{\phi}} \bar{\boldsymbol{\phi}}' \right) + s_d \times \chi \times I_d, & \text{if } j > J_0, \quad 0 < \chi \ll 1 \end{cases} \end{aligned}$$

with  $\bar{\boldsymbol{\phi}}$  denoting the mean value of all draws up to and including the  $(j-1)$ -th and with  $s_d$  suggested, following Gelman et al. (1996), to be set to  $s_d = \frac{2.4^2}{d}$  where  $d$  is the number of parameters in  $\boldsymbol{\phi}$ <sup>10</sup>.

Under this proposal distribution the  $j$ -th candidate value  $\tilde{\boldsymbol{\phi}}^{(j)}$  is obtained by sampling from a multivariate Gaussian distribution centered at the previously retained draw  $\tilde{\boldsymbol{\phi}}^{(j-1)}$  and with covariance matrix  $\Sigma^{(j-1)}$ . Being initially set to a given matrix  $\Sigma^{(0)}$ , this covariance matrix starts adapting exploiting all past draws after a sufficiently large initial period  $J_0$  following the sample covariance matrix. An (*AM*) algorithm can thus focus on sampling more densely in regions near values  $\tilde{\boldsymbol{\phi}}$  with high posterior density and less densely in regions of low posterior density. It is also possible to extend the scope of the local accept-reject exploration by sampling  $M$  candidates at once from  $g_{\Sigma^{(j-1)}}(\boldsymbol{\phi}, \tilde{\boldsymbol{\phi}}^{(j-1)})$  in the spirit of the multiple-try Metropolis sampler of Liu et al. (2000).

---

<sup>10</sup>As discussed in Haario et al. (2001), the addition of the diagonal matrix  $\chi \times I_d$  is needed with an insignificantly small but non-zero  $\chi$  to assure the non-singularity of  $\Sigma^{(j-1)}$  and assure the convergence of the MCMC sampling distribution to  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}|\mathbf{y}^{Obs})$ .

Output  $\{(\tilde{\boldsymbol{\theta}}^{(j)}, \tilde{\boldsymbol{\eta}}^{(j)}, \tilde{\boldsymbol{\nu}}^{(j)})\}_{j=1}^J$  from an MCMC algorithm may be used to diagnose sampler issues in terms of convergence. One possible expression of convergence issues, where the retained samples cannot be taken as representative of the posterior distribution they target, is that of non-stationarity with respect to the sample average. Several traceplots computed on the sampler output can be useful in detecting this issue, such as the cusum path plots of [Yu and Mykland \(1998\)](#).

## 4.2 Approximate Bayesian inference through the Generalized Lorenz curve

Implementing the (**AM**) algorithm requires being able to compute the likelihood function  $L(\mathbf{y}^{Obs}|\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$ . For a model following (1), the joint likelihood for a sample of independent microdata  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$  can be computed as  $L(\mathbf{y}^{Obs}|\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}) = \prod_{i=1}^N f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$ . However, joint likelihood functions for microdata prone to uncertain measurement error and/or truncation points are often challenging to compute or study numerically. In particular, deterministic under-reporting or non-response schemes like LPU or right-truncation that only affect observations above a fixed threshold income will introduce jumps into  $f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  that in consequence introduce discontinuities in the joint data likelihood  $L(\mathbf{y}^{Obs}|\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  (e.g., see [Chernozhukov and Hong 2004](#)). These discontinuities are a function of the observed data  $\mathbf{y}^{Obs}$  and the model parameters  $(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$ . For example, a model including a right-truncation form for non-response can jump to a joint likelihood value of zero when evaluated at parameter values  $(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  which imply a truncation below the highest observed income(s) in  $\mathbf{y}^{Obs}$  and this case is likely to happen over several different parameter values. Additionally, a computable likelihood function may not be available in several contexts, such as when data on incomes is only available at group level (e.g., see [Kobayashi and Kakamu 2019](#), [Eckernkemper and Gribisch 2021](#)).

In devising a more flexible implementation of the (**AM**) sampling idea in light of the potentially high complexity of the likelihood function one possibility is to focus on a different representation of the income distribution. Instead of opting for its observed density  $f_{\mathbf{y}^{Obs}}(\cdot; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  the model and data can be represented by their corresponding Generalized Lorenz curve (*GLC*). The *GLC* ([Shorrocks 1983](#), [Kakwani 1984](#)) can be defined as the cumulative of the quantile function below the  $u$ -th percentile:

$$GLC(u; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}) \equiv \int_0^u Q_{\mathbf{y}^{Obs}}(x; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}) dx,$$

where  $Q_{\mathbf{y}^{Obs}}$  denotes the quantile function of the parametric model  $f_{\mathbf{y}^{Obs}}$  and where  $GLC(1; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}) \equiv E[y_i^{Obs}|\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}]$  defines the average observed income. In addition to allowing for comparisons between income distributions in terms of relative inequality as any Lorenz curve, the *GLC* also allows for absolute comparisons in terms of mean incomes.

Opting for the *GLC* as a representation for the income distribution has as main advantage a more regular functional form than the underlying density  $f_{\mathbf{y}^{Obs}}$  due to its cumulative nature. Additionally, the *GLC* is a representation compatible with both microdata or grouped data. The main challenge when working under this representation is that a tractable likelihood function for the *GLC* is rarely available. However, this

likelihood function can be approximated in practice whenever a mechanism for simulating samples from  $f_{y^{Obs}}(\cdot; \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  is available.

An empirical  $GLC$  may be computed from a sample of observed incomes'  $\mathbf{y}^{Obs}$  order statistics as<sup>11</sup>:

$$GLC_k^{Obs} = \underbrace{\frac{\sum_{i=1}^k y_{(i)}^{Obs}}{\sum_{i=1}^N y_{(i)}^{Obs}}}_{s_k^{Obs}} \times \underbrace{\frac{1}{N} \sum_{i=1}^N y_{(i)}^{Obs}}_{\mu^{Obs}} = \frac{\sum_{i=1}^k y_{(i)}^{Obs}}{N}, \quad k = 1, \dots, N, \quad GLC_0^{Obs} = 0,$$

with  $s_k^{Obs}$  denoting the cumulative income share up to the  $k$ -th observation in the ordered sample and  $\mu^{Obs}$  denoting the sample average income. As most common formats of grouped data on incomes allow for computing  $K < N$  cumulative income shares and a sample mean, a grouped-data  $GLC$  is simply a subset of the  $N$  points  $\{GLC_k^{Obs}\}_{k=1}^N$  from the empirical  $GLC$  of its underlying microdata.

Being able to simulate samples of observed incomes  $\tilde{\mathbf{y}}^{Obs}$  for given parameter values following  $\tilde{\mathbf{y}}^{Obs} \sim f_{y^{Obs}}(\cdot; \tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})$  is generally feasible for any typical income distribution model and non-response or under-reporting forms<sup>12</sup>. This can be exploited for the purpose of Bayesian inference on  $(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  from a sample of observed incomes  $\mathbf{y}^{Obs}$  represented through their empirical  $GLC$  (denoted  $\{GLC_k^{Obs}\}_{k=1}^N$  in what follows) through a class of simulation-based inference methods known as Approximate Bayesian Computation (ABC) (e.g., see [Kobayashi and Kakamu 2019](#), [Silva 2023](#)).

ABC can approximate the unavailable likelihood function of the  $GLC$   $L(\{GLC_k^{Obs}\}_{k=1}^N | \tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})$  in a non-parametric way through comparing the empirical  $GLC$  from simulated data samples  $\tilde{\mathbf{y}}^{Obs}$  (denoted  $\{GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})\}_{k=1}^N$ ) and the observed  $\{GLC_k^{Obs}\}_{k=1}^N$ . This approximation requires a way of assessing how closely the empirical income distribution  $\{GLC_k^{Obs}\}_{k=1}^N$  resembles data simulated from the model for any given parameter values  $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}) \in \boldsymbol{\Theta}$ .

The overall degree of discrepancy between the observed and simulated-data empirical income distributions may be summarized by the following unidimensional metric:

$$d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})\}_{k=1}^N) = \sum_{k=1}^N |(GLC_k^{Obs} - GLC_{k-1}^{Obs}) - (GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}) - GLC_{k-1}^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}))|,$$

which corresponds to the empirical Wasserstein-1 distance ([Kantorovich, 1939](#)) in the case of microdata<sup>13</sup>. Explored in the context of ABC by [Bernton et al. \(2019\)](#), this distance

<sup>11</sup>In the case of microdata samples from surveys with income-ordered weights  $\{w_{(i)}\}_{i=1}^N$  the weighted empirical  $GLC$  may be computed similarly as:

$$GLC_k^{Obs} = \frac{\sum_{i=1}^k y_{(i)}^{Obs} w_{(i)}}{\sum_{i=1}^N y_{(i)}^{Obs} w_{(i)}} \times \frac{1}{\sum_{i=1}^N w_{(i)}} \sum_{i=1}^N y_{(i)}^{Obs} w_{(i)} = \frac{\sum_{i=1}^k y_{(i)}^{Obs} w_{(i)}}{\sum_{i=1}^N w_{(i)}}, \quad k = 1, \dots, N, \quad GLC_0^{Obs} = 0.$$

<sup>12</sup>For simplicity, it is assumed in what follows that simulated data are in the form of independent microdata  $\tilde{\mathbf{y}} = \{\tilde{y}_i^{Obs}\}_{i=1}^N$ .

<sup>13</sup>See the derivations in [Appendix A](#) for a detailed description of this distance and a grouped-data implementation of this discrepancy.



summarizes the absolute discrepancies between all order statistics across observed and simulated data  $|y_{(i)}^{Obs} - \tilde{y}_{(i)}^{Obs}|$ ,  $i = 1, \dots, N$ .

In approximating  $L(\{GLC_k^{Obs}\}_{k=1}^N | \boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  for an ABC implementation of the (**AM**) algorithm, parameter values  $(\boldsymbol{\theta}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}) \in \boldsymbol{\Theta}$  yielding simulated data  $\{GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})\}_{k=1}^N$  resembling  $\{GLC_k^{Obs}\}_{k=1}^N$  more closely under  $d(\cdot, \cdot)$  than others should be given larger importance. This is commonly introduced exploiting a kernel function  $K_\tau$  giving increasingly larger weight to parameter values with a lower discrepancy  $\varepsilon \equiv d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})\}_{k=1}^N)$ . A common 'smooth' kernel for this purpose is the Gaussian kernel (e.g., see [Ratmann 2010](#)):

$$K_\tau^{gauss}(\varepsilon) = \frac{1}{\tau} \times \frac{1}{\sqrt{2\pi}} \times \exp \left\{ -\frac{1}{2} \left( \frac{\varepsilon}{\tau} \right)^2 \right\}, \quad \varepsilon \equiv d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}})\}_{k=1}^N).$$

Under this kernel the ABC discrepancies are weighted following a Normal distribution centered at zero (i.e., highest weight is given to values  $(\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}, \tilde{\boldsymbol{\nu}}) \in \boldsymbol{\Theta}$  exactly reproducing  $\{GLC_k^{Obs}\}_{k=1}^N$ ), and with a standard deviation given by the bandwidth parameter  $\tau$ .

By opting for the  $GLC$  as representation of the income distribution the target posterior distribution for Bayesian inference is no longer  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu} | \mathbf{y}^{Obs})$  but  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu} | \{GLC_k^{Obs}\}_{k=1}^N)$ . By approximating the likelihood function for the  $GLC$  under the ABC approach this latter target posterior distribution is also approximated. The ABC target posterior distribution for the parameter vector  $\boldsymbol{\phi} = (\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  can be stated as ([Drovandi and Frazier, 2022](#)):

$$\pi_\tau(\boldsymbol{\phi} | \{GLC_k^{Obs}\}_{k=1}^N) \propto \underbrace{\int_{\mathbb{R}^N} K_\tau(\epsilon) \times L(\{GLC_k^{Obs}\}_{k=1}^N | \boldsymbol{\phi}) d\tilde{GLC} \times p(\boldsymbol{\phi})}_{L_\tau(\{GLC_k^{Obs}\}_{k=1}^N | \boldsymbol{\phi})}, \quad \tilde{GLC} = \{GLC_k^{Obs}(\boldsymbol{\phi})\}_{k=1}^n, \quad (3)$$

with  $\epsilon \equiv d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\boldsymbol{\phi})\}_{k=1}^N)$ . The intractable integral defining  $L_\tau(\{GLC_k^{Obs}\}_{k=1}^N | \boldsymbol{\phi})$  may be unbiasedly estimated in practice for any given point  $\tilde{\boldsymbol{\phi}} \in \boldsymbol{\Theta}$  using  $Z$  simulated income distributions from the parametric model following:

$$\hat{L}_\tau(\{GLC_k^{Obs}\}_{k=1}^N | \tilde{\boldsymbol{\phi}}) = \frac{1}{Z} \sum_{z=1}^Z K_\tau(\epsilon^{(z)}), \quad \epsilon^{(z)} \equiv d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs;(z)}(\tilde{\boldsymbol{\phi}})\}_{k=1}^N), \quad z = 1, \dots, Z,$$

and where  $Z$  is commonly set to  $Z = 1$ .

As a bandwidth parameter,  $\tau$  rules the strictness of the ABC non-parametric approximation to  $L(\{GLC_k^{Obs}\}_{k=1}^N | \tilde{\boldsymbol{\phi}})$  by defining how the weights  $K_\tau^{gauss}(\varepsilon)$  decrease with an increase in the discrepancy  $\varepsilon$ . The approximation is exact when  $\tau \rightarrow 0$ , as only parameter values which exactly reproduce the observed income distribution are given a non-zero weight, and  $\tau \rightarrow \infty$  amounts to considering any and all parameter values in  $\boldsymbol{\Theta}$  equally likely to have generated the observed data (i.e., the likelihood is approximated as a flat function).

The ABC posterior distribution  $\pi_\tau(\boldsymbol{\phi} | \{GLC_k^{Obs}\}_{k=1}^N)$  might differ from that in (2) for several reasons. One first source of differences lies on the quality of the approximation to the exact posterior distribution  $\pi(\boldsymbol{\phi} | \{GLC_k^{Obs}\}_{k=1}^N)$ . The main determinant of this is

the choice for the bandwidth parameter  $\tau$ . ABC implementations of sampling algorithms targeting  $\pi_\tau(\boldsymbol{\phi}|\{GLC_k^{Obs}\}_{k=1}^N)$  in the spirit of the (**AM**) algorithm pay an increasing computational cost for a stricter approximation through a lower  $\tau$ . In practice, the choice for this bandwidth results from calibrating the sampling algorithm through several initial runs balancing strictness of the approximation and computational cost.

The second main source for differences between  $\pi_\tau(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}|\{GLC_k^{Obs}\}_{k=1}^N)$  and  $\pi(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu}|\mathbf{y}^{Obs})$  concerns the possible loss of information due to summarizing the data through the GLC and not through the microdata directly. If what can be learnt about  $(\boldsymbol{\theta}, \boldsymbol{\eta}, \boldsymbol{\nu})$  from the data represented through the *GLC* is less than what can be learnt from microdata then their respective estimated posterior distributions will differ even when the ABC approximation to the likelihood is exact (i.e., when  $\tau \rightarrow 0$ ).

Following Silva (2023), an ABC (**AM**) algorithm with these settings can be devised extending the Marjoram et al. (2003) ABC implementation of the (**MH**) algorithm. Algorithm 2 below presents a possible implementation, denoted (**ABC-AM**) in what follows.

---

**Algorithm 2:** An AM ABC (**ABC-AM**) algorithm.

---

**Initialization:**

- 1: Set  $\Sigma^{(0)}$ ,  $J_0$ ,  $M$ ,  $\tau$   
**Until**  $K_\tau^{gauss}(\tilde{\varepsilon}^{(0)}) > 0$ :
- 2: Sample  $(\tilde{\boldsymbol{\theta}}^{(0)}, \tilde{\boldsymbol{\eta}}^{(0)}, \tilde{\boldsymbol{\nu}}^{(0)}) \equiv \tilde{\boldsymbol{\phi}}^{(0)}$  from  $p(\tilde{\boldsymbol{\phi}})$
- 3: Generate  $\{GLC_k^{Obs}(\tilde{\boldsymbol{\phi}}^{(0)})\}_{k=1}^N$  by simulating from  $f_y^{Obs}(\cdot; \tilde{\boldsymbol{\phi}}^{(0)})$
- 4: Generate  $\tilde{\varepsilon}^{(0)} = d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\phi}}^{(0)})\}_{k=1}^N)$

**Sampling:**

**for**  $j = 1, \dots, J$  **do**

- 5: Sample  $\{\tilde{\boldsymbol{\phi}}^{(m)}\}_{m=1}^M \sim g_{\Sigma^{(j-1)}}(\boldsymbol{\phi}, \tilde{\boldsymbol{\phi}}^{(j-1)})$  from the candidate  $g_{\Sigma^{(j-1)}}$
- 6: Generate  $\{GLC_k^{Obs}(\tilde{\boldsymbol{\phi}}^{(m)})\}_{k=1}^N$  by simulating from  $f_y^{Obs}(\cdot; \tilde{\boldsymbol{\phi}}^{(m)})$ ,  $m = 1, \dots, M$
- 7: Generate  $\tilde{\varepsilon}^{(j)} = \min_{m \in \{1, \dots, M\}} d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\phi}}^{(m)})\}_{k=1}^N)$  and candidate  $\tilde{\boldsymbol{\phi}}^{(j)} = \arg \min_{\tilde{\boldsymbol{\phi}}^{(m)}} d(\{GLC_k^{Obs}\}_{k=1}^N, \{GLC_k^{Obs}(\tilde{\boldsymbol{\phi}}^{(m)})\}_{k=1}^N)$
- 8: Accept and store  $(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\varepsilon}^{(j)})$  with probability:

$$\rho^{(j)} = \min \left\{ 1, \frac{K_\tau^{gauss}(\tilde{\varepsilon}^{(j)}) \times p(\tilde{\boldsymbol{\phi}}^{(j)}) \times g_{\Sigma^{(j-1)}}(\tilde{\boldsymbol{\phi}}^{(j-1)}, \tilde{\boldsymbol{\phi}}^{(j)})}{K_\tau^{gauss}(\tilde{\varepsilon}^{(j-1)}) \times p(\tilde{\boldsymbol{\phi}}^{(j-1)}) \times g_{\Sigma^{(j-1)}}(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\boldsymbol{\phi}}^{(j-1)})} \right\}$$

▷ e.g., if  $u^{(j)} \leq \rho^{(j)}$  where  $u^{(j)}$  is a draw from a *Uniform*(0, 1) distribution

otherwise store  $(\tilde{\boldsymbol{\phi}}^{(j)}, \tilde{\varepsilon}^{(j)}) = (\tilde{\boldsymbol{\phi}}^{(j-1)}, \tilde{\varepsilon}^{(j-1)})$

**if**  $j > J_0$  **then**

- 9: Update  $\Sigma^{(j)}$

**end**

**end**

---

At any  $j$ -th step, the (**ABC-AM**) algorithm draws  $M$  candidate parameter values  $\{\tilde{\boldsymbol{\phi}}^{(m)}\}_{m=1}^M$  from the adaptive proposal  $g_{\Sigma^{(j-1)}}$ , simulates a single income distribution from

the model for each such candidate, and computes their respective discrepancies with respect to the observed income distribution. The candidate with the lowest discrepancy is then taken as the  $j$ -th candidate  $\tilde{\phi}^{(j)}$ , along with its associated ABC discrepancy  $\tilde{\varepsilon}^{(j)}$ , in the same spirit as [Clarté et al. \(2021\)](#). Finally, the MH accept-reject rule is computed with respect to the ABC approximation of the likelihood through  $K_\tau(\tilde{\varepsilon}^{(j)})$ .

Together, the parametric framework for income distributions under  $MR$  issues developed in the previous section along with the Bayesian empirical strategy presented in this section allow for a broad range of applications. The following section illustrates some of the main income distribution analysis possible under this approach.

## 5 Applications and examples

### 5.1 Applications on simulated data

Simulated-data applications can give insight on the performance of the ABC approach in making inference on  $(\theta, \eta, \nu)$  in a controlled setting exploiting a model following (1). Consider a hypothetical population's income distribution following a GB2 distribution  $y_i \sim f_{\mathbf{y}}^{GB2}(\cdot; \theta) \equiv GB2(\alpha, \beta, p, q)$ , with parameters  $\alpha$ ,  $p$ , and  $q$  ruling the shape of the distribution and  $\beta$  ruling the scale. Typically, these parameters are the focus of the analysis of the income distribution. However, if the available data  $\mathbf{y}^{Obs}$  is presumably affected by any of the  $MR$  forms considered in the previous sections, additional parameters ruling assumed parametric forms for these issues must also be introduced into the analysis.

Assume that microdata samples from this population's income distribution may be jointly affected by high-income under-reporting following an LPU scheme with parameters  $(\bar{p}, \delta)$  and high-income non-response following a right-truncation scheme with  $\alpha$  fixed to  $\alpha = 1$  and parameter  $t$  where  $t \gg \bar{p}$ . Under this joint scheme a model for the observable data  $\mathbf{y}^{Obs}$  can be obtained applying (1)<sup>14</sup>:

$$\begin{aligned}
f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \theta, \bar{p}, \delta, t) &= \frac{f_{\mathbf{y}}^{GB2}(m^{-1}(y_i^{Obs}; \bar{p}, \delta); \theta) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} \right) \times \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t)}{\int f_{\mathbf{y}}^{GB2}(m^{-1}(y_i^{Obs}; \bar{p}, \delta); \theta) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} \right) \times \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t) dy^{Obs}} \\
&= \frac{f_{\mathbf{y}}^{GB2} \left( y_i^{Obs} + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1; GB2}(\bar{p}; \theta)) \times \left( \frac{\delta(y_i^{Obs} - F_{\mathbf{y}}^{-1; GB2}(\bar{p}; \theta))}{1 - \delta} \right); \theta \right)}{t \times (1 - \delta \times \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1; GB2}(\bar{p}; \theta)))} \\
&\times \frac{\mathbf{1}(y_i^{Obs} \leq (1 - \delta)F_{\mathbf{y}}^{-1; GB2}(t; \theta) + \delta F_{\mathbf{y}}^{-1; GB2}(\bar{p}; \theta))}{t \times (1 - \delta \times \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1; GB2}(\bar{p}; \theta)))}. \tag{4}
\end{aligned}$$

[Appendix A](#) presents details on the derivations required for this expression.

<sup>14</sup>In what follows,  $F^{-1; GB2}(\cdot; \alpha, \beta, p, q)$  denotes the quantile function of the GB2 distribution. For simplicity of notation,  $\theta = (\alpha, \beta, p, q)$  also holds in what follows.

Equation (4) expands the GB2 distribution to allow for LPU (whenever  $\bar{p} \ll t$  and  $\delta > 0$ ) and for non-response in the form of a right-truncation (whenever  $\bar{p} \ll t < 1$ ). For illustrative purposes, a first experiment of interest consists in estimating the posterior distribution  $\pi(\boldsymbol{\theta}, \bar{p}, \delta, t | \mathbf{y}^{Obs})$  using the (**ABC-AM**) algorithm through this model over a sample of  $N$  simulated incomes  $\mathbf{y}^{Obs} = \{y_i^{Obs}\}_{i=1}^N$ . In particular, this exercise is most interesting when the simulated data is effectively affected by LPU and right-truncation forms of  $MR$  jointly.

Benchmark parameter values can be set to  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\bar{p}, \delta, t) = (.5, .15, .99)$  in this interest<sup>15</sup>. These correspond to a population income distribution with an average income of 15054 and a Gini coefficient of 0.348. Data simulated under this setting corresponds to a sample from a GB2 distribution which starts being affected by LPU above the median with a slope of  $\delta = .15$  and which contains no observations for units above the .99-th population's income distribution percentile.

Data can be simulated in this specific case by sampling  $\frac{N}{t}$  incomes from the GB2 distribution and applying the LPU and right-truncation transformations under the benchmark values. This yields a single random sample of  $N$  observed incomes. The samples used in this exercise were generated in this way, for a hypothetical population of 10000 units (i.e.,  $N = 9900$ ). Figure 3 below illustrates how a sample generated under this setting relates to the theoretical observed incomes' distribution  $f_{\mathbf{y}^{Obs}}$  under (4) and to the respective complete population's  $f_{\mathbf{y}}^{GB2}$  income distribution.

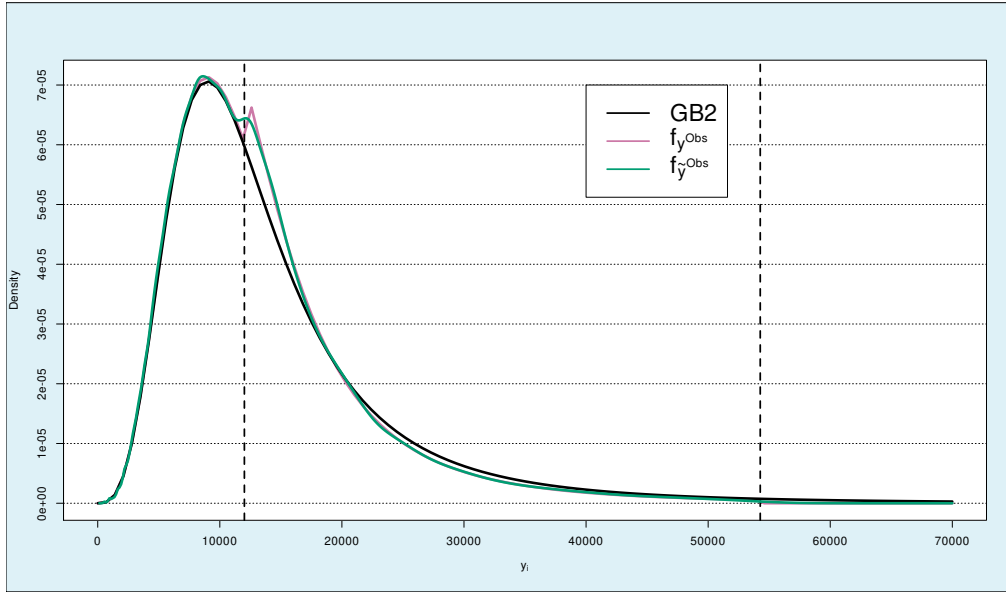


Figure 3: Population, theoretical, and sample densities for model (4)

**Note:** Three density curves describing a population's GB2 income distribution (GB2 in legend), the observed incomes density function following (4) ( $f_{\mathbf{y}^{Obs}}$  in legend), and kernel density estimate from an estimating sample generated from this model ( $N = 9900$ ) ( $f_{\tilde{\mathbf{y}}^{Obs}}$  in legend). Benchmark parameter values taken as  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\bar{p}, \delta, t) = (.5, .15, .99)$ , with the left-most vertical dashed line representing the population's  $\bar{p}$ -th income percentile and the right-most vertical dashed line representing the right-truncation point.

<sup>15</sup>The benchmark value for  $\beta$  being 10000, it is here scaled by 1000 in the interest of numerical stability when applying (**ABC-AM**).

Several conditions must be considered in eliciting a joint prior probability distribution for the model parameters  $p(\boldsymbol{\theta}, \bar{p}, \delta, t)$  in practice. Firstly, this joint prior distribution can be set as the product of several marginal prior distributions:

$$p(\boldsymbol{\theta}, \bar{p}, \delta, t) = p(\alpha) \times p\left(\frac{\beta}{1000}\right) \times p(p) \times p(q) \times p(\bar{p}) \times p(\delta) \times p(t).$$

Secondly, given the high flexibility of the GB2 distribution it is possible to represent virtually any specific case of this distribution in a constrained range of parameter values. In this sense, the marginal prior distributions for the GB2 parameters may be set as follows:

$$\begin{aligned}\alpha &\sim p(\alpha) \equiv \text{Gamma}(1, 1) \\ \frac{\beta}{1000} &\sim p\left(\frac{\beta}{1000}\right) \equiv \text{Gamma}(5, 2) \\ p &\sim p(p) \equiv \text{Gamma}(1, 1) \\ q &\sim p(q) \equiv \text{Gamma}(1, 1).\end{aligned}$$

This amounts to prior beliefs on the shape parameters  $\alpha$ ,  $p$ , and  $q$  following a right-skewed Gamma distribution with mode at the value 1 and to prior beliefs on  $\frac{\beta}{1000}$  following another right-skewed Gamma distribution with mode approximately at the value 8.

Thirdly, reflecting a strong prior belief on the presence *MR* issues in the data, the  $(\bar{p}, \delta, t)$  parameters may be given the following prior distributions:

$$\begin{aligned}\bar{p} &\sim p(\bar{p}) \equiv \text{Beta}(8, 5) \\ \delta &\sim p(\delta) \equiv \text{Beta}(1, 5) \\ (1 - t) &\sim p(1 - t) \equiv \text{Beta}(1, 25).\end{aligned}$$

These reflect empirically-relevant values for the literature using right-truncation forms for non-response (e.g., see [Jorda and Niño-Zarazúa 2019](#)) and that exploring high-income under-reporting in survey data (e.g., see [Flachaire et al. 2022](#)). Importantly, these prior beliefs also give considerable probability to the ‘complete data’ scenario where no under-reporting or non-response issues affect the sample. This is, it is also made likely *a priori* that the observed income distribution may be correctly represented by a single GB2 distribution without introducing *MR* phenomena.

Finally, several constraints may be imposed on the elicited joint prior distribution to further constrain the parameter space. Imposing restrictions for finite variance on the GB2 income distribution amounts to giving 0 prior probability to parameter values with  $\alpha < \frac{2}{q}$  and  $\alpha < -\frac{1}{p}$ . Additionally, because under-reported incomes have no relevance if they correspond to a true income above the truncation point, the restriction  $t > \bar{p}$  is imposed<sup>16</sup>. Figure 4 below summarizes these elicited prior distributions for each of the

---

<sup>16</sup>Formally, these restrictions impose the following joint prior distribution:

$$p(\boldsymbol{\theta}, \bar{p}, \delta, t) = p(\alpha) \times p\left(\frac{\beta}{1000}\right) \times p(p) \times p(q) \times p(\bar{p}) \times p(\delta) \times p(t) \times \prod_{i=1}^3 C_{(i)}(\boldsymbol{\theta}, \bar{p}, \delta, t),$$

model's parameters.

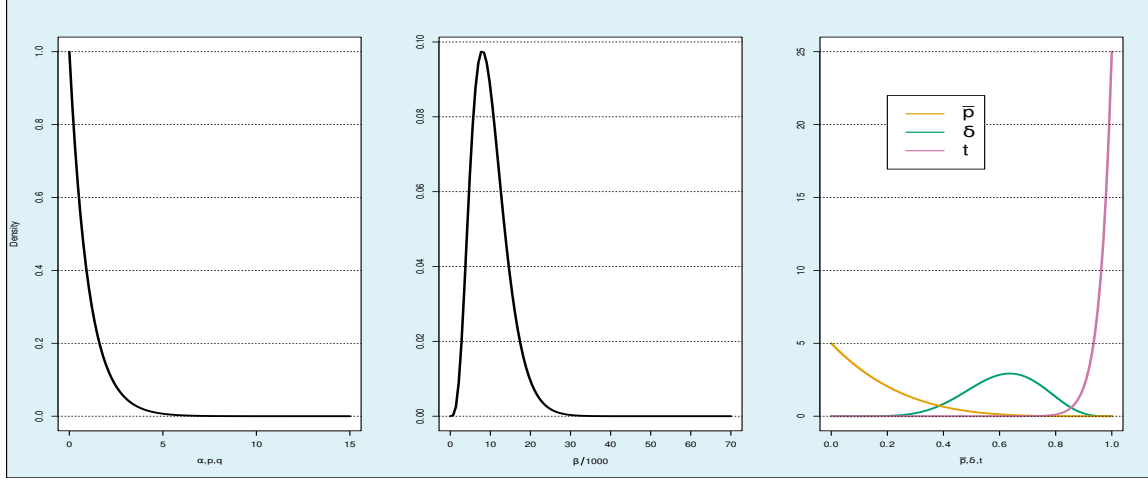


Figure 4: Prior distributions elicited for the parameters of model (4)

**Note:** *Left:* Prior distributions for  $\alpha$ ,  $p$ , and  $q$  parameters. *Center:* Prior distribution for  $\frac{\beta}{1000}$ . *Right:* Prior distributions for  $\bar{p}$ ,  $\delta$ , and  $t$  parameters.

## Applications on a single sample

As a first exercise, three central scenarios are explored applying the (**ABC-AM**) algorithm to a single simulated-data sample for clarity of illustration. Firstly, to evidence the possible biases that these forms of  $MR$  may induce if the issue is not taken into consideration, a simple GB2 distribution is fit to the data. A second scenario consists of estimating the income distribution parameters  $(\alpha, \beta, p, q)$  under (4) conditional on fixing the correction quantities  $(\boldsymbol{\eta}, \boldsymbol{\nu}) = (\bar{p}, \delta, t)$  at their true values. Finally, a third scenario consists of estimating all parameters in (4) eliciting prior uncertainty in  $(\boldsymbol{\eta}, \boldsymbol{\nu})$ . In all cases, the algorithm is set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples are obtained, taking the initial  $J_0 = 50000$  draws as the burn-in period where the algorithm's adaptive terms are calibrated. For computational ease, the simulated data taken as estimating sample was summarized by its GLC computed at sample centiles  $\{GLC_k^{Obs}\}_{k=1}^{100}$  which provides a highly detailed summary of the overall sample  $\{GLC_k^{Obs}\}_{k=1}^N$ .

Figure 5 below illustrates for all three scenarios explored the goodness-of-fit of the resulting posterior distribution estimates in terms of their fit to the true pdf for this simulated sample of observed incomes computed under (4)<sup>17</sup>. These estimated posterior

with

$$\begin{cases} C_{(1)}(\boldsymbol{\theta}, \bar{p}, \delta, t) = \mathbf{1}\left(\alpha > \frac{2}{q}\right) \\ C_{(2)}(\boldsymbol{\theta}, \bar{p}, \delta, t) = \mathbf{1}\left(\alpha > -\frac{1}{p}\right) \\ C_{(3)}(\boldsymbol{\theta}, \bar{p}, \delta, t) = \mathbf{1}(t > \bar{p}) \end{cases} \quad .$$

<sup>17</sup>The posterior distribution for the pdf coordinates  $\pi_{\tau}^{f_{\mathbf{y}^{Obs}}} (f_{\mathbf{y}^{Obs}}(.; \boldsymbol{\theta}, \bar{p}, \delta, t) | \mathbf{y}^{Obs})$  is computed in what follows from the retained samples from the parameters' posterior distribution  $\{(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^J$  as  $\{f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^J$  with  $f_{\mathbf{y}^{Obs}}$  following (4) and with  $\bar{p}^{(j)}, \delta^{(j)}, t^{(j)}$  taking a same fixed value when no uncertainty is introduced on these parameters.



distributions are summarized in each case by their corresponding 95% Highest Posterior Density Interval (HPDI)<sup>18</sup> which constitutes an interval estimate of the pdf at each level of incomes in the range of the estimating sample  $\mathbf{y}^{Obs}$ .

As a first remark, the figure illustrates that a simple GB2 model, corresponding to fixing  $(\bar{p}, \delta, t) = (1, 0, 1)$ , is insufficient in this case to accommodate for the relatively high complexity of the observed-incomes distribution, producing results which under-estimate the mass of top incomes and which over-estimate the mass of incomes at the mode of the distribution as well as at the middle-high incomes region. Secondly, both scenarios which allow for the possibility of *MR* issues affecting the sample achieve a good fit to the true pdf with a higher degree of uncertainty at the mode of the distribution in the case where prior uncertainty is elicited for the *MR* parameters  $(\bar{p}, \delta, t)$ . Only in these two latter cases the true pdf is contained within the computed *HPDI*(.95) credibility intervals all along the income distribution.

Estimated ABC marginal posterior distributions for each of the income distribution parameters in (4) are summarized in figure 8 in [Appendix A](#). As a first observation, all scenarios yield posterior distribution estimates which significantly differ from the elicited prior distributions, effectively updating these prior beliefs. A most relevant result is the strong bias of parameter estimates obtained without considering *MR* issues. In this scenario the  $\alpha$  shape parameter is under-estimated while the scale parameter  $\beta$  and the shape parameters  $p$  and  $q$  are all over-estimated. In contrast, both scenarios correcting for these issues yield estimated posterior distributions centered at their true value. Additionally, introducing uncertainty on the  $(\bar{p}, \delta, t)$  parameters yields posterior distributions for the income distribution parameters which reflect only slightly higher uncertainty than the respective estimates obtained under the known true values for these. Finally, insight on the performance of the (**ABC-AM**) sampler underlying these marginal density estimates can be obtained from traceplots illustrating the sequence of draws from the algorithm. For the more complex case where *MR* parameters are uncertain a priori, the traceplots presented in figure 8 suggest a stable behavior of the MCMC sampler around the respective parameter's true value after the initial burn-in period, not providing any evidence of convergence issues.

---

<sup>18</sup>A 95% highest posterior density interval (*HPDI*(.95)) may be estimated summarizing a region of values for the posterior predictive distribution of the pdf at any fixed income level  $y^{Obs}$  with estimated posterior density  $\pi_{\tau}^{f_{\mathbf{y}^{Obs}}}(f_{\mathbf{y}^{Obs}}(y^{Obs}; \boldsymbol{\theta}, \bar{p}, \delta, t) | \mathbf{y}^{Obs})$  above a threshold  $c \in (0, 1)$ . In the case of a unimodal posterior distribution, this  $c$  defines the narrowest continuous interval of values accumulating an estimated posterior mass of .95 on the posterior distribution represented by  $\{f_{\mathbf{y}^{Obs}}(y^{Obs}; \boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^J$ . This is equivalent to identifying a threshold value  $c$  defining the interval of values with estimated posterior densities above it:

$$HPDI(.95) = \left\{ f_{\mathbf{y}^{Obs}}(y^{Obs}; \boldsymbol{\theta}, \bar{p}, \delta, t) : \int (\pi_{\tau}^{f_{\mathbf{y}^{Obs}}}(x | \mathbf{y}^{Obs}) \geq c) \times \pi_{\tau}^{f_{\mathbf{y}^{Obs}}}(x | \mathbf{y}^{Obs}) dx = .95 \right\}.$$

*HPDI*(.95) then provides an interval estimate for  $f_{\mathbf{y}^{Obs}}(y^{Obs}; \boldsymbol{\theta}, \bar{p}, \delta, t)$  from integrating information on the estimated joint parameter posterior distribution  $\pi_{\tau}(\boldsymbol{\theta}, \bar{p}, \delta, t | \mathbf{y}^{Obs})$ .

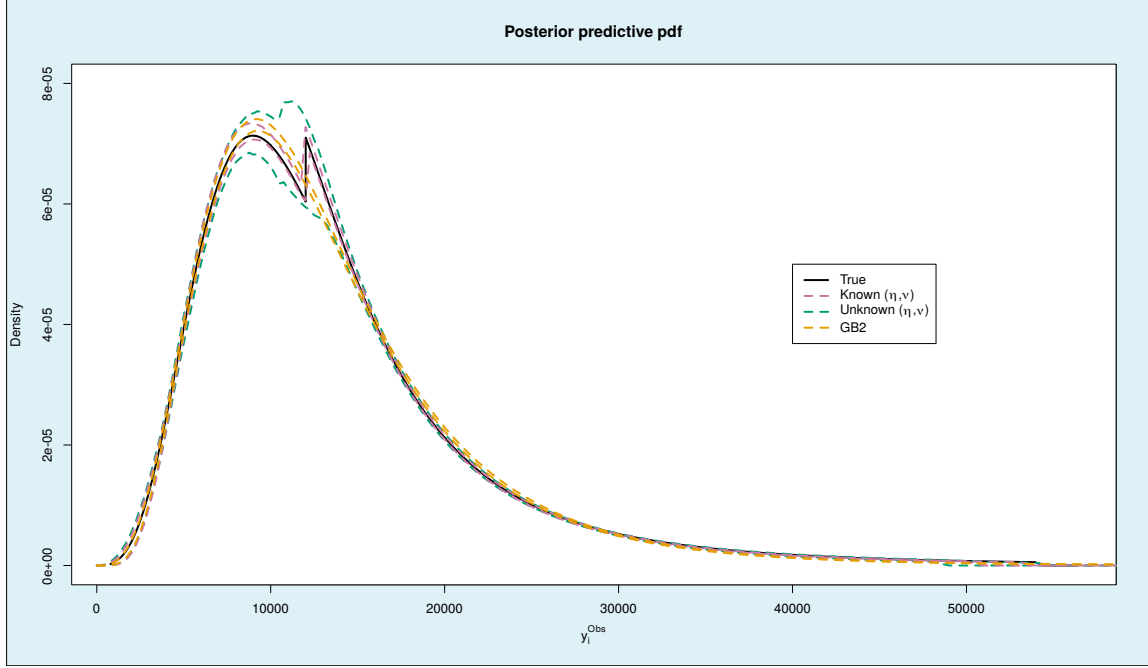


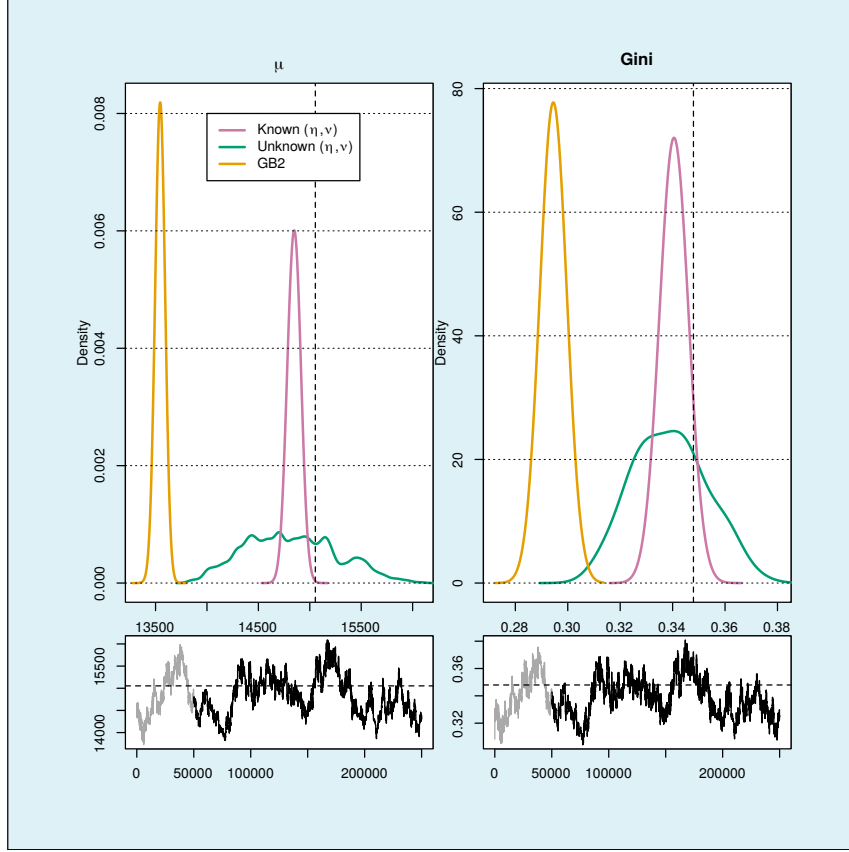
Figure 5: Sample and ABC posterior predictive pdf estimates for model (4) on a simulated sample.

**Note:** Observed incomes density function following (4) for a simulated sample of  $N = 9900$  incomes under parameter values  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$  in solid black. ABC HPDI(.95) interval-estimates for this pdf plotted in dashed curves. Estimates obtained applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a GB2), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases, the algorithm is set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained ( $\{\theta^{(j)}\}_{j=1}^{250000}$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period.

For the scenario where  $(\bar{p}, \delta, t)$  are uncertain a priori and are therefore also to be inferred from the data, figure 9 in Appendix A summarizes the estimated posterior distributions for these parameters. The estimates showcase a significant update of the elicited prior distributions, with posterior distributions centered at the true values for these correcting quantities. This result provides support of the ABC approach as a fruitful empirical strategy for parametric inference on income distributions through data affected by MR issues of uncertain magnitude.

One additional illustration that this exercise provides concerns the impact of accounting for MR issues in the available data when making inference on income growth and inequality at their respective population levels. Figure 6 below presents the posterior predictive distributions of the population's mean income, which is often used to track aggregate income growth across time, and Gini coefficient of inequality. Both statistics are determined by the  $(\alpha, \beta, p, q)$  coefficients alone under a GB2 distribution, following expressions denoted  $G^{GB2}(\alpha^{(j)}, p^{(j)}, q^{(j)})$  and  $\mu^{GB2}(\alpha^{(j)}, \beta^{(j)}, p^{(j)}, q^{(j)})$  in what follows, as derived in McDonald and Ransom (2008) and implemented in Graf and Nedyalkova. (2015).

Figure 6: Posterior predictive estimates of population mean income and Gini coefficient



**Note:** Kernel density estimates for ABC predictive posterior distribution estimates of population mean income and Gini coefficient computed on a single simulated sample of  $N = 9900$  observed incomes following (4) under parameter values  $\left(\alpha, \frac{\beta}{1000}, p, q\right) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Estimates obtained applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a GB2), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases, the algorithm is set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained ( $\{\theta^{(j)}\}_{j=1}^{250000}$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period. Traceplots of the underlying MCMC samples for estimates with uncertainty on  $(\eta, \nu)$  below, with burn-in period in gray. True parameter values in dashed black lines. Both statistics computed as in Graf and Nedyalkova. (2015).

A main implication of these results is illustrated in this figure. These posterior predictive distributions evidence a significant under-estimation of both mean income and income inequality when neglecting MR issues is inappropriate. Only the scenarios estimated under a margin for corrections through  $(\bar{p}, \delta, t)$  achieve estimates for these statistics closely reproducing their true values at the population level. This provides further evidence of the biased reading that might be made of a population's income distribution when the possibility of MR issues affecting the available data is not integrated into the analysis.

## Monte Carlo experiments

Having detailed the performance of the proposed method on a single simulated-data sample limited by MR issues, a next possible exercise is that of analyzing this performance across several samples and setups. For this purpose, the same analysis as in the previous

sub-section was performed over 50 microdata samples of size  $N/t = 10000$  generated under each of three setups defined as follows. Firstly, in the interest of studying the performance of the ABC approach under a simpler parametric model affected by  $MR$  and nested in (4), setup *I* sets the true parameter values to  $(\alpha, \beta, p, q) = (2.3, 10000, 1, 1.25)$ , corresponding to a three-parameter Singh-Maddala population income distribution with parameters  $(\alpha, \beta, q) = (2.3, 10000, 1.25)$ , and  $(\bar{p}, \delta, t) = (.5, .15, .99)$  which corresponds to observed income distributions affected by the same  $MR$  scenario studied in the previous sub-section.

Setup *II* also focuses on the specific case of (4) when the population income distribution follows a Singh-Maddala distribution (i.e., when  $p$  is fixed to  $p = 1$ ). In this setup, however, no  $MR$  issues affect the simulated samples (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$ ). This serves to address the question of whether allowing for the possibility of  $MR$  with prior uncertainty on the corresponding parameters can nonetheless produce unbiased estimates properly reflecting the lack of such issues in the observed data<sup>19</sup>.

Finally, setup *III* reproduces the same situation as explored in the previous sub-section but over 50 simulated income distributions. This setup serves to illustrate how the more complex GB2 population income distribution affects the performance of the ABC method in contrast with the simpler Singh-Maddala setups.

Results obtained under setup *I* are graphically summarized by figures 10 and 11 in [Appendix A](#). A first interesting remark concerns the marked differences between the interval estimates obtained for each of the model parameters under all three estimation scenarios and the prior distributions elicited for them which evidences that the data effectively provides information to learn about these parameters and update significantly the prior beliefs. This observation is particularly relevant for the estimates obtained under prior uncertainty on the  $MR$  parameters  $(\bar{p}, \delta, t)$  which illustrate that learning jointly about the population's income distribution parameters and these  $MR$  parameters affecting the estimating sample occurs across all samples. Secondly, these results also illustrate the biases affecting those estimates obtained neglecting the possibility of  $MR$  issues in the data. These estimates provide an over-estimated value for the scale of the population's income distribution and for the shape parameter  $q$  ruling its right tail. Finally, comparing parameter estimates obtained conditional on the true values for the high-income LPU and right-truncation affecting the data with those obtained under prior uncertainty on these illustrates that the latter reflect only slightly larger uncertainty in their estimated posterior distributions.

Attending to the possibility of model misspecification issues introduced by considering high-income under-reporting and/or non-response when these do not affect the sample of observed incomes, results obtained under setup *II* are illustrated in figures 12 and 13 in [Appendix A](#). These figures compare estimates produced under a simple Singh-Maddala distribution, obtained as a special case of the GB2 distribution by fixing the parameter  $p$  to  $p = 1$ , to those obtained under a Singh-Maddala distribution affected by LPU and right-truncation  $MR$  forms, obtained as a special case of (4) also by fixing  $p = 1$ , on simulated data samples that have not been affected by  $MR$  issues.

---

<sup>19</sup>Note that in both setups exploiting the Singh-Maddala distribution the initial covariance matrix of the proposal function for the (**ABC-AM**) algorithm is reduced to  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01)$ .

A main observation that these figures provide is that the ABC method provides unbiased and precise estimates of the population income distribution’s parameters in both situations across all samples, with only slightly lower precision when estimates are produced under prior uncertainty on the  $MR$  parameters. This is further illustrated by the interval-estimates obtained for the  $(\bar{p}, \delta, t)$  parameters which correctly estimate a posterior distribution for  $t$  heavily concentrated at its true value of  $t = 1$  and for the under-reporting rate  $\delta$  heavily concentrated at its true value of  $\delta = 0$ . Interestingly, these latter estimates also illustrate that in absence of under-reporting the ABC method accommodates this by inferring one of the two parameters ruling the LPU form to represent no under-reporting and is unable to learn about the other parameter. In this case, this can be seen by posterior estimates for  $\bar{p}$  which closely follow the prior distribution elicited for this parameter while the posterior distribution for  $\delta$  significantly updates the prior beliefs on this parameter to yield a posterior distribution reflecting negligible levels of high-income LPU, if any.

Figures 14 and 15 in [Appendix A](#) illustrate the performance of the proposed ABC method under (4) for the GB2 population income distribution case. Similarly to what is illustrated on the case of a single simulated-data sample in the previous subsection, the resulting estimates only appropriately infer the income distribution parameters when the possibility of  $MR$  issues affecting the observed incomes is considered. Fitting a GB2 distribution on these simulated samples produces estimates which are biased and, for some parameters, very imprecise. This is particularly illustrated by the obtained interval-estimates for the  $p$  and  $q$  shape parameters which often are incompatible with the true parameter values underlying the data and which reflect much larger uncertainty than the respective estimates obtained considering  $MR$  issues. When compared to the Singh-Maddala setups, all resulting posterior distribution estimates show significantly larger uncertainty and variability across samples which can reflect a loss in precision due to the introduction of prior uncertainty on the  $p$  parameter. Interestingly, the interval-estimates obtained in this setup for the  $(\bar{p}, \delta, t)$  parameters ruling the forms of  $MR$  issues explored when these are considered uncertain a priori are almost identical in terms of precision to those obtained in the analogous scenario under setup  $I$ .

## 5.2 Real data applications:

As an illustrative example on real data, the European Union’s Statistics on Income and Living Conditions (EU-SILC) provide an interesting household survey setting. EU-SILC data provides information on people and households within the EU representative at the country level, covering most countries in the EU yearly since 2005, and under a common framework defining the exact definitions of incomes and populations to be surveyed.

A key income variable in income distribution analysis on EU-SILC data is household disposable income under the OECD-modified equivalence scale<sup>20</sup> (HX090). This considers

---

<sup>20</sup>The OECD-modified equivalence scale computes a household’s size HX050 as:

$$\begin{aligned} \text{HX050} = & 1 + 0.5 \times (\text{number of household members aged 14 and over} - 1) \\ & + 0.3 \times (\text{number of household members aged 13 or less}). \end{aligned}$$

all gross household incomes in the data net of regular taxes on wealth, regular inter-household transfers paid, and regular taxes on income and social insurance contributions.

Although this aggregate variable includes definitions of income variables that are common to all countries covered by EU-SILC, the sources from which the data are obtained differ across cross-sectional waves of data and countries. In particular, while some countries rely entirely on survey responses to measure these income variables, other countries source these variable entirely or partially from administrative registers. These differences in sources across waves and countries introduce large heterogeneities in the quality of the data in terms of under-reporting as registers are considered more reliable than survey responses. Additionally, the rising use of register sources determines that for some countries some of the waves of data are sourced from surveys and other waves are sourced from register. This can produce trends in the observed income distributions across waves without necessarily reflecting trends of the true population's income distribution.

Although several calibrations are done over EU-SILC samples and sample weights for enforcing population representativeness on several dimensions, these are not done on income variables (with the exception of The Netherlands). Recent analysis have explored *MR* phenomena on EU-SILC data (e.g., see [Hlasny and Verme 2018](#), [Bartels et al. 2019](#), [Angel et al. 2019](#), [Carranza et al. 2023](#), [Ederer et al. 2022](#)), suggesting this issue to be present with different magnitudes in all countries and periods studied. This makes it such that exploiting the provided survey weights for inference on a country-year's population income distribution is not exempt of representativeness issues arising from *MR* phenomena. In this application, and in the interest on making inference on a population's income distribution, EU-SILC data is exploited using the provided cross-sectional household weights (DB090) alongside the possibility of *MR* issues in the available sample of incomes.

Information on non-response rates for each country and wave of EU-SILC is publicly available through the corresponding quality reports published by the European Commission. Household non-response rates, in particular, can be informative about the overall degree of non-response affecting an observed distribution of household incomes. These rates are computed from a country-level household response rate, which is the product of address contact rates (i.e., the share of households in the sampling frame that were successfully contacted) and household response rates (i.e., the share of households in the sampling frame that completed their survey after being successfully contacted).

Table 1 below summarizes EU-SILC samples for five selected countries (Austria, Germany, France, Spain, and Italy) and for the 2005, 2007, 2011, and 2016 waves. With the exception of Germany, all other selected four countries are known to exploit register data sources on incomes in complementing EU-SILC survey responses, although the timing and extent of this practice is poorly documented in general (e.g., see information on EU-SILC income data sources in [Jäntti et al. 2013](#), [Carranza et al. 2023](#), and [Wirth and Pforr 2022](#)). The mean and Gini coefficient for household disposable income distributions computed using their respective survey weights summarize the observable trends in growth and inequality across countries and waves. Because these distributions

---



are presumably affected by *MR* issues, these values may provide a biased estimate of the corresponding population’s mean and Gini coefficient of incomes. With many heterogeneities, all countries experienced mean income growth and, with the exception of Italy, income inequality increased from 2005 to 2016 based on observed incomes alone. Finally, overall household non-response rates show large disparities across countries and years in terms of levels and trends, illustrating possible heterogeneities in the incidence of this issue on the respective observed income distributions.

Table 1: EU-SILC sample descriptives for selected countries

Country	Wave	$N$	Household non-response rate	$\mu^{Obs}$	Gini
Austria (AT)	2005	5146	0.38	20212.24	0.27
	2007	6805	0.22	20405.17	0.28
	2011	6182	0.23	23948.16	0.29
	2016	5992	0.27	26274.72	0.28
Germany (DE)	2005	13078	0.35	18078.73	0.27
	2007	14047	—	20084.84	0.31
	2011	13473	0.21	21047.33	0.30
	2016	13260	0.23	23424.24	0.31
France (FR)	2005	9745	0.16	18237.49	0.29
	2007	10485	0.14	18423.25	0.27
	2011	11348	0.18	23934.22	0.31
	2016	11446	0.17	25788.05	0.30
Spain (ES)	2005	12865	0.28	12289.05	0.33
	2007	12234	0.23	13520.56	0.32
	2011	12993	0.22	16535.78	0.33
	2016	14168	0.20	16151.14	0.34
Italy (IT)	2005	21874	0.15	16648.63	0.33
	2007	20809	0.14	17422.81	0.32
	2011	19234	0.25	18491.99	0.32
	2016	20966	0.21	18839.71	0.32

**Source:** Own calculations from EU-SILC.

**Note:** EU-SILC samples for Austria (AT), Germany (DE), France (FR), Spain (ES), and Italy (IT) from 2005, 2007, 2011, and 2016 waves. Only considers households with reported household disposable income (HX090) of at least 1 euro. Household non-response rates as reported in the publicly-available quality reports for each wave. Weighted-sample estimates of mean incomes ( $\mu^{Obs}$ ) and Gini coefficients computed using cross-sectional household weights (DB090).

Under the same settings explored in the simulated data applications, model (4) can be fit to the EU-SILC samples through the (**ABC-AM**) algorithm. In summarizing the distribution of observed incomes in any given EU-SILC sample through its empirical GLC survey weights were exploited to compute sample GLC coordinates following the details in footnote 11 and then only those coordinates corresponding to sample centiles were retained as estimating data  $\{GLC_k^{Obs}\}_{k=1}^{100}$ . These GLC coordinates summarize the corresponding country-years’ EU-SILC observed incomes distribution at the population level and corresponds therefore to the distribution that (4) is a model for.

To obtain estimates of the posterior distributions of interest, the (**ABC-AM**)

algorithm was implemented for each and all selected EU-SILC samples. In all cases, the algorithm was set similarly to the applications explored in simulated data, with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples were obtained, taking the initial  $J_0 = 50000$  draws as the burn-in period.

In the interest of comparing the resulting estimates with the observed income distributions in EU-SILC data, posterior predictive distributions of the mean observed income and the Gini coefficient of observed incomes may be computed for the population and compared with their corresponding weighted-sample values. Without analytical expressions for these statistics corresponding to model (4) they may be computed from the simulated observed income distributions corresponding to each MCMC sample  $\{\{GLC_k^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{k=1}^{100}\}_{j=1}^{250000}$ . The posterior predictive distribution of the mean observed income  $\{\mu^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{J=250000}$  can be computed simply as  $\{GLC_{100}^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  while that of the Gini coefficient  $\{G^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  can be computed from the associated Lorenz curve coordinates  $\{LC_k^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{k=1}^{100}$  with:

$$\begin{cases} LC_k^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)}) &= \frac{GLC_k^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})}{GLC_{100}^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})} \\ G^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)}) &= \frac{2}{100} \times \sum_{k=1}^{100} \frac{k}{100} - LC_k^{Obs}(\boldsymbol{\theta}^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)}). \end{cases}$$

Figure 7 below summarizes the resulting estimates in terms of goodness-of-fit to the weighted-sample mean and Gini coefficient of observed incomes, as well as the estimated values for these corresponding to the population's GB2 income distribution. In what follows, the fitted income distribution refers to estimates of the population-level distribution of observed incomes, which may be affected by *MR* phenomena, while the population income distribution refers to estimates of the distribution of 'true' incomes as modeled by the GB2 distribution component in (4). As a first observation, the obtained results for the fitted income distribution reproduce very accurately the levels and trends of mean observed incomes and Gini coefficients for all waves and countries considered.

Concerning the estimates of each population's income distribution a first important result is that these do not match their corresponding estimates fitting the observed incomes' distribution in any case. This suggests that *MR* corrections are indeed required under this model in order to fit the observed distributions accurately under (4). Consequently, the estimated GB2 population income distribution parameters imply levels of mean incomes and inequality above their weighted-sample values. These population level estimates reproduce similar changes in mean incomes across waves as those in the fitted distributions, with some increase in the uncertainty around these quantities for the latter years in the case of France and Italy. Population estimates that reproduce the same dynamic of mean incomes as those reflected in the fitted distribution of incomes can be indicative of the total mass of incomes missing in this latter distribution due to *MR* issues changing very little across EU-SILC waves.

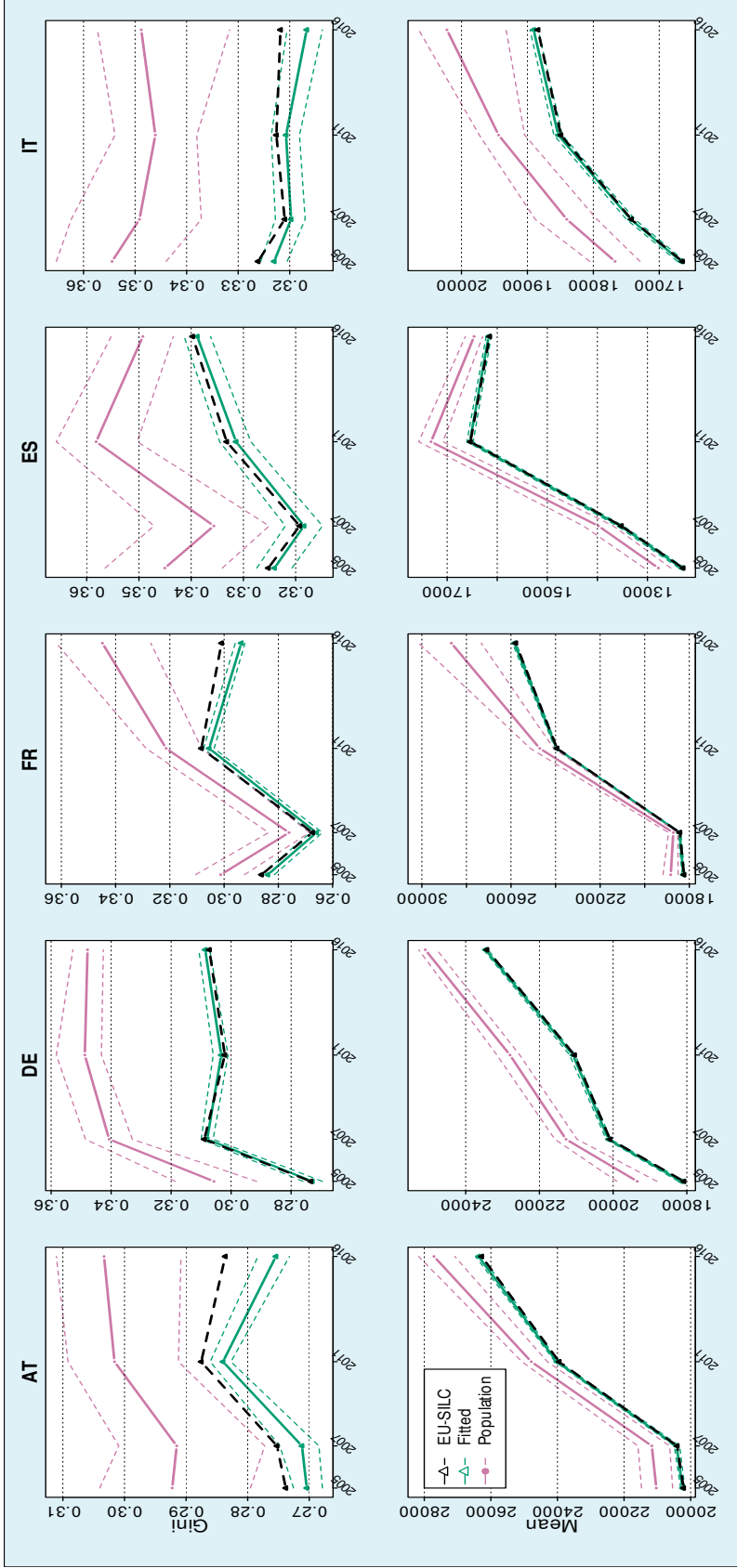


Figure 7: Posterior predictive estimates of mean income and Gini coefficients - EU-SILC countries.

**Source:** Own calculations from EU-SILC.  
**Note:** ABC predictive posterior mean estimates of mean household disposable income (HX090) and Gini coefficients for Austria (AT), Germany (DE), France (FR), Spain (ES), and Italy (IT) from 2005, 2007, 2011, and 2016 EU-SILC waves. Only considers households with reported household disposable income (HX090) of at least 1 euro. Respective 95% HPDI in dashed lines. In black, weighted-sample estimates of mean incomes and Gini coefficients computed using cross-sectional household weights (DB090). Estimates obtained following (4) applying the (ABC-AM) algorithm. In all cases, the algorithm is set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(0.1, 1, 0.1, 0.1, 0.1, 0.1)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained, taking the initial  $J_0 = 50000$  draws as the burn-in period. Fitted values computed from the posterior predictive distribution of observed incomes at the population-level as summarized by the simulated income distributions corresponding to each MCMC sample  $\{GLC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{k=1}^{100}, \{GLC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{k=1}^{100}\}$  while that of the Gini coefficient is computed from the associated Lorenz curve coordinates  $\{LC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{k=1}^{100}$  with  $LC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)}) = GLC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)}) / GLC_{100}^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})$  as  $\{(2/100) \times \sum_{k=1}^{100} (k/100) - LC_k^{Obs}(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$ . Population-level statistics' predictive posterior distribution computed parametrically from the retained  $\{\theta^{(j)}\}_{j=1}^{250000}$  for the assumed GB2 distribution using the implementation of Graf and Nedjalkova. (2015).

As summarized by the Gini coefficient, the income inequality dynamics implied by these estimates at the population level pose some contrasts with respect to their observed incomes' counterpart. In all cases inequality is estimated to be higher at the population level than what estimates on observed incomes alone suggest. The non-overlap of the computed credible intervals across population and fitted observed incomes' estimates of the Gini provide strong evidence of this. Another observation is that the uncertainty around the population Gini is relatively stable across EU-SILC waves for all countries, with strong heterogeneities across countries. This uncertainty is in some cases high enough to make inequality increases and decreases across time equally likely, as is the case for the Austrian Gini between the 2011 and 2016 waves. In other cases, however, the estimates provide clearer evidence of significant increases in income inequality across periods, as can be seen for France and Germany across 2005 and 2016 waves.

Conditional on the assumed forms for  $MR$  issues in these applications, the estimated parameters can suggest margins where the representativeness of the data changes across waves. The detailed estimates in table 2 in [Appendix B](#) suggest that the right-truncation parameter  $t$  introducing high-income non-response in the model (4) is estimated to be within the top .1% of the population income distribution in all waves and countries. This can suggest that significant non-response issues are mostly concentrated on households within this income group.

Concerning under-reporting issues, these estimates suggest strong heterogeneities in the share of the population affected by progressive under-reporting and the progressiveness of under-reporting across waves and countries. As quantified by the estimates for the  $\bar{p}$  and  $\delta$  parameters these are estimated to range from .5401 to .8525, and from .0779 to .3997 respectively across the selected EU-SILC samples. Taken together with the estimates for  $t$ , these results illustrate that fitting the observed income distributions accurately under model (4) always requires jointly correcting for progressive under-reporting of incomes above the median and for right-truncation non-responses somewhere within the top .1% of the income distribution.

## 6 Concluding remarks

Building on previous  $MR$  correction methods explored in the literature, a new framework for parametric inference on income distributions has been explored in this paper. This framework directly integrates this type of corrections into the process of making inference on a population's income distribution by expanding conventional parametric income distribution models with parametric functional forms for both reporting and non-response mechanisms. As a model for observed data on incomes presumably affected by  $MR$  this parametric approach then allows for devising empirical strategies to infer jointly features of the associated population's income distribution and features of the  $MR$  issues in the data.

In dealing with the several constraints that must be faced in devising an empirical strategy for this purpose, the ABC approach has been implemented as a suitable method. This Bayesian estimator allows for updating prior uncertainty on the 'true' population income distribution and the often uncertain  $MR$  quantities affecting the observed data.

The main criteria driving inference under this approach is attempting to reproduce the observed incomes distribution, summarized through its empirical GLC, with simulated GLCs from the parametric model for this distribution.

The illustrative applications presented in this paper evidence several virtues of the ABC approach for inference on the parameters of GB2 population income distributions under data affected by *MR*. In the Monte Carlo simulated-data setting, the analysis illustrates the several possible biases affecting inference on a population's income distribution when *MR* issues affecting the observed data are neglected. This experiment also suggests the ABC approach to be fruitful in learning about uncertain *MR* quantities given informative prior beliefs about these. Finally, no particular risks are evidenced in scenarios allowing for the possibility of *MR* corrections when these are not affecting the estimating data and, in particular, the proposed method is able to update prior beliefs on the *MR* quantities to correctly reflect that no significant *MR* corrections are required in reproducing the observed income distribution.

Applications to cross-sectional EU-SILC data on household disposable income distributions of Austria, Germany, France, Spain, and Italy between 2005 and 2016 give further insight on the suitability of the framework in a typical household survey microdata setting. The resulting estimates imply that reproducing the observed income distributions accurately requires considering some amount of both high-income under-reporting and non-response phenomena in all settings analysed. These applications also illustrate how inference on population income distributions can be made under a priori uncertain *MR* quantities, uncovering contrasts between most of the observed incomes' distribution trends and those inferred for the respective population distributions.

As a first implementation of the framework developed in this paper, however, several aspects both theoretical and empirical have been left unattended and can propose venues for future research. Future work could build on these developments firstly by exploring the empirical strategy implemented in this paper in a linked-data setting where individual observed incomes from a survey prone to be affected by *MR* phenomena may be matched to a corresponding register income from sources presumably less prone to these issues like fiscal data on incomes. If comparing survey-sourced incomes and register-sourced incomes evidences some form of progressive under-reporting and high-income non-response, then estimates obtained under this framework using the survey data alone could be validated in terms of reproducing *MR* patterns consistent with these.

In the specific case of the EU-SILC, an additional direction for future work concerns integrating available external information on the representativeness of the observed income distribution into the prior beliefs for the *MR* quantities. In particular, household non-response rates and information about the specific sources for the observed incomes for a given wave of data and country can be exploited in setting up informative prior probabilities. This could help in accounting for possible artificial trends arising from changes in the sampling or income sources and not from actual changes in the population's income distribution.

Further work seeking to provide better understanding of the possible pitfalls of the proposed approach could explore further setups in terms of the specified *MR* parametric

forms beyond the LPU and right-truncation forms explored in the applications of this paper. In principle, if several alternative forms are capable of representing similar patterns of under-reporting and non-response, then estimates obtained under each of these forms should yield very similar results on the population’s income distribution. Assessing the robustness of the ABC approach in situations of model misspecification due to invalid assumptions on the forms of under-reporting or non-response in a Monte Carlo experimental setting can also provide further insight on the properties of this framework.

Better understanding of the properties of the ABC estimator explored in this paper could also be achieved by studying calibration schemes for the several parameters involved in applying the (**ABC-AM**) algorithm. In particular, calibrating the bandwidth parameter  $\tau$  is of key importance as it ultimately rules the strictness of the approximation to the posterior distribution that inference is made on while at the same time can also heavily determine the computational cost of achieving samples representative of this approximation. While a stricter approximation is likely to reduce the uncertainty reflected in the approximated posterior distribution and so may allow for finer inference on the income distribution, it is also likely to increase the rejection rate of the MCMC sampling algorithm as it forces a stronger constraint on which simulated samples are considered of sufficiently close resemblance to the income distribution observed in the data. The convergence and efficiency properties of the MCMC sampler, which configure the overall computational cost of applying the ABC estimator, are not only determined by the choice for  $\tau$  but also by the parameters ruling the adaptive proposal function and so all of them must be taken into account jointly for the purpose of sampler calibration.

Implemented as it is in this paper, the ABC estimator is compatible with applications exploiting grouped data. Whenever the available data allows for computing coordinates of the empirical GLC of observed incomes (e.g., when the data is presented in the widespread form of mean observed incomes for different groups along the income distribution) then these may be exploited for inference on the income distribution through comparisons with simulated observed income distributions from the parametric model corresponding to the same percentiles of the distribution. Studying the possible losses of quality on the inference that is made on the population’s income distribution under grouped data of different sizes is an exercise for further extensions of this empirical strategy.

Another possible future development involves extending the framework to accommodate for other non-response or measurement error issues presumably affecting the distribution of observed incomes. In particular, several studies have evidenced the existence of differential non-response and income misreporting phenomena in the lower end of the distribution of observed incomes in household survey data (e.g., see [Pedace and Bates 2000](#), [Mathiowetz et al. 2002](#), [Meyer and Mittag 2019](#), [Angel et al. 2019](#), [Hlasny et al. 2022](#), [Flachaire et al. 2022](#)). Although the influence of these issues on estimates of population income inequality might be meager they can introduce important biases to measures of poverty or income growth of groups at the lower end of the distribution. Considering these issues jointly with the *MR* phenomena in a single parametric framework may provide an interesting approach for making inference on the population’s income distribution through a yet more comprehensive modelling of observed incomes.



Finally, a possible extension of this framework involves making inference on income distributions of populations defined at aggregate levels such as regions or the globe. Taking the mixture of all countries' income distributions, estimates obtained accounting for *MR* issues at the country level can be used to study patterns of income growth and distribution on aggregate levels which also take into account the heterogeneities that these *MR* issues can present across countries and years.

## References

- Alvaredo, F. (2011). A note on the relationship between top income shares and the Gini coefficient. *Economics Letters*, 110(3):274–277.
- Alvaredo, F., Chancel, L., Piketty, T., Saez, E., and Zucman, G. (2018). The elephant curve of global inequality and growth. In *AEA Papers and Proceedings*, volume 108, pages 103–08.
- Angel, S., Disslbacher, F., Humer, S., and Schnetzer, M. (2019). What did you really earn last year?: explaining measurement error in survey income data. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 182(4):1411–1437.
- Atkinson, A. B. and Piketty, T. (2007). *Top incomes over the twentieth century: a contrast between continental European and english-speaking countries*. Oxford University Press.
- Atkinson, A. B., Piketty, T., and Saez, E. (2011). Top incomes in the long run of history. *Journal of economic literature*, 49(1):3–71.
- Bartels, C., Metzing, M., et al. (2019). An integrated approach for a top-corrected income distribution. *Journal of Economic Inequality*, 17(2):125–143.
- Bernton, E., Jacob, P. E., Gerber, M., and Robert, C. P. (2019). Approximate Bayesian computation with the Wasserstein distance. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 81(2):235–269.
- Berthet, P., Fort, J.-C., and Klein, T. (2020). A central limit theorem for Wasserstein type distances between two distinct univariate distributions. In *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, volume 56, pages 954–982. Institut Henri Poincaré.
- Blanchet, T., Flores, I., and Morgan, M. (2022). The weight of the rich: improving surveys using tax data. *The Journal of Economic Inequality*, 20(1):119–150.
- Bollinger, C. R., Hirsch, B. T., Hokayem, C. M., and Ziliak, J. P. (2019). Trouble in the tails? what we know about earnings nonresponse 30 years after Lillard, Smith, and Welch. *Journal of Political Economy*, 127(5):2143–2185.
- Bourguignon, F. (2018). Simple adjustments of observed distributions for missing income and missing people. *The Journal of Economic Inequality*, 16(2):171–188.
- Brownstone, D. and Valletta, R. G. (1996). Modeling earnings measurement error: A multiple imputation approach. *The Review of Economics and Statistics*, pages 705–717.
- Brunori, P., Salas-Rojo, P., and Verme, P. (2022). Estimating inequality with missing incomes. *ECINEQ Working Paper*, 616.
- Burdín, G., Esponda, F., and Vigorito, A. (2014). Inequality and top incomes in Uruguay: a comparison between household surveys and income tax micro-data. *World Top Incomes Database Working Paper*, 1.

- Burkhauser, R. V., Hérault, N., Jenkins, S. P., and Wilkins, R. (2017). Top incomes and inequality in the UK: reconciling estimates from household survey and tax return data. *Oxford Economic Papers*, 70(2):301–326.
- Bustos, A. (2015). Estimation of the distribution of income from survey data, adjusting for compatibility with other sources. *Statistical Journal of the IAOS*, 31(4):565–577.
- Carranza, R., Morgan, M., and Nolan, B. (2023). Top income adjustments and inequality: An investigation of the EU-SILC. *Review of Income and Wealth*, 69(3):725–754.
- Charpentier, A. and Flachaire, E. (2022). Pareto models for top incomes and wealth. *The Journal of Economic Inequality*, 20(1):1–25.
- Chernozhukov, V. and Hong, H. (2004). Likelihood estimation and inference in a class of nonregular econometric models. *Econometrica*, 72(5):1445–1480.
- Chesher, A. and Schluter, C. (2002). Welfare measurement and measurement error. *The Review of Economic Studies*, 69(2):357–378.
- Chotikapanich, D., Griffiths, W., Hajargasht, G., Karunaratne, W., and Rao, D. (2018). Using the GB2 income distribution. *Econometrics*, 6(2):21.
- Chotikapanich, D. and Griffiths, W. E. (2000). Applications: posterior distributions for the Gini coefficient using grouped data. *Australian & New Zealand Journal of Statistics*, 42(4):383–392.
- Chotikapanich, D. and Griffiths, W. E. (2008). Estimating income distributions using a mixture of Gamma densities. In Chotikapanich, D., editor, *Modeling income distributions and Lorenz curves*, pages 285–302. Springer.
- Clarté, G., Robert, C. P., Ryder, R. J., and Stoeck, J. (2021). Componentwise approximate Bayesian computation via Gibbs-like steps. *Biometrika*, 108(3):591–607.
- Darvas, Z. (2019). Global interpersonal income inequality decline: The role of China and India. *World Development*, 121:16–32.
- Deaton, A. (2005). Measuring poverty in a growing world (or measuring growth in a poor world). *Review of Economics and statistics*, 87(1):1–19.
- Drovandi, C. and Frazier, D. T. (2022). A comparison of likelihood-free methods with and without summary statistics. *Statistics and Computing*, 32(3):42.
- Eckernkemper, T. and Gribisch, B. (2021). Classical and Bayesian inference for income distributions using grouped data. *Oxford Bulletin of Economics and Statistics*, 83(1):32–65.
- Ederer, S., Četković, P., Humer, S., Jestl, S., and List, E. (2022). Distributional national accounts (DINA) with household survey data: Methodology and results for European countries. *Review of Income and Wealth*, 68(3):667–688.
- Flachaire, E., Lustig, N., and Vigorito, A. (2022). Underreporting of top incomes and inequality: A comparison of correction methods using simulations and linked survey and tax data. *Review of Income and Wealth*.

- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2013). *Bayesian data analysis*. CRC press.
- Gelman, A., Roberts, G., and Gilks, W. (1996). Efficient Metropolis jumping rules. In Bernardo, J. M., Berger, J. O., Dawid, A. P., and Smith, A. F. M., editors, *Bayesian Statistics V. Proceedings of the Fifth Valencia International Meeting*. Oxford University Press.
- Gottschalk, P. and Huynh, M. (2010). Are earnings inequality and mobility overstated? the impact of nonclassical measurement error. *The Review of Economics and Statistics*, 92(2):302–315.
- Graf, M. and Nedyalkova, D. (2015). *GB2: Generalized Beta Distribution of the Second Kind: Properties, Likelihood, Estimation*. R package version 2.1.
- Greenlees, J. S., Reece, W. S., and Zieschang, K. D. (1982). Imputation of missing values when the probability of response depends on the variable being imputed. *Journal of the American Statistical Association*, 77(378):251–261.
- Gustafson, P. (2005). On model expansion, model contraction, identifiability and prior information: Two illustrative scenarios involving mismeasured variables. *Statistical Science*, 20(2):111–140.
- Haario, H., Saksman, E., Tamminen, J., et al. (2001). An adaptive Metropolis algorithm. *Bernoulli*, 7(2):223–242.
- Hartley, M. J. and Revankar, N. S. (1974). On the estimation of the Pareto law from under-reported data. *Journal of Econometrics*, 2(4):327–341.
- Hinkley, D. V. and Revankar, N. S. (1977). Estimation of the Pareto law from underreported data: A further analysis. *Journal of Econometrics*, 5(1):1–11.
- Hlasny, V. (2020). Nonresponse bias in inequality measurement: Cross-country analysis using Luxembourg Income Study surveys. *Social Science Quarterly*, 101(2):712–731.
- Hlasny, V., Ceriani, L., and Verme, P. (2022). Bottom incomes and the measurement of poverty and inequality. *Review of Income and Wealth*, 68(4):970–1006.
- Hlasny, V. and Verme, P. (2015). Top incomes and the measurement of inequality: A comparative analysis of correction methods using Egyptian, EU, and US survey data. In *ECINEQ Conference Paper*, volume 145.
- Hlasny, V. and Verme, P. (2018). Top incomes and inequality measurement: A comparative analysis of correction methods using the EU-SILC data. *Econometrics*, 6(2):1–21.
- Hlasny, V. and Verme, P. (2022). The impact of top incomes biases on the measurement of inequality in the United States. *Oxford Bulletin of Economics and Statistics*, 84(4):749–788.
- Hurst, E., Li, G., and Pugsley, B. (2014). Are household surveys like tax forms? evidence from income underreporting of the self-employed. *Review of economics and statistics*, 96(1):19–33.

- Jäntti, M., Tormalehto, V.-M., and Marlier, E. (2013). The use of registers in the context of EU-SILC: challenges and opportunities. Statistical working papers, Eurostat.
- Jenkins, S. P. (1995). Did the middle class shrink during the 1980s? UK evidence from kernel density estimates. *Economics letters*, 49(4):407–413.
- Jenkins, S. P. (2009). Distributionally-sensitive inequality indices and the GB2 income distribution. *Review of Income and Wealth*, 55(2):392–398.
- Jenkins, S. P. (2017). Pareto models, top incomes and recent trends in UK income inequality. *Economica*, 84(334):261–289.
- Jorda, V. and Niño-Zarazúa, M. (2019). Global inequality: How large is the effect of top incomes? *World Development*, 123:104593.
- Jorda, V., Sarabia, J. M., and Jäntti, M. (2021). Inequality measurement with grouped data: parametric and non-parametric methods. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 184(3):964–984.
- Kakwani, N. (1984). Welfare ranking of income distributions. *Advances in econometrics*, 3:191–213.
- Kantorovich, L. V. (1939). The mathematical method of production planning and organization. *Management Science*, 6(4):363–422.
- Kobayashi, G. and Kakamu, K. (2019). Approximate Bayesian computation for Lorenz curves from grouped data. *Computational Statistics*, 34(1):253–279.
- Korinek, A., Mistiaen, J. A., and Ravallion, M. (2007). An econometric method of correcting for unit nonresponse bias in surveys. *Journal of Econometrics*, 136(1):213–235.
- Krishnaji, N. (1970). Characterization of the pareto distribution through a model of underreported incomes. *Econometrica*, 38(2):251–55.
- Lakner, C. and Milanovic, B. (2016). Global income distribution: From the fall of the Berlin wall to the great recession. *The World Bank Economic Review*, 30(2):203–232.
- Leigh, A. (2009). Top incomes. In Salverda, W., Nolan, B., and Smeeding, T. M., editors, *The Oxford handbook of economic inequality*. New York: Oxford University Press.
- Liu, J. S., Liang, F., and Wong, W. H. (2000). The multiple-try method and local optimization in Metropolis sampling. *Journal of the American Statistical Association*, 95(449):121–134.
- Lustig, N. (2019). Measuring the distribution of household income, consumption and wealth: State of play and measurement challenges. In Stiglitz, J. E., Fitoussi, J.-P., and Durand, M., editors, *For Good Measure: An Agenda for Moving Beyond GDP*. The New Press.
- Lustig, N. (2020). The missing rich in household surveys: Causes and correction approaches. *ECINEQ Working Paper No. 2020-520*.
- Lyssiotou, P., Pashardes, P., and Stengos, T. (2004). Estimates of the black economy based on consumer demand approaches. *The Economic Journal*, 114(497):622–640.

- Marjoram, P., Molitor, J., Plagnol, V., and Tavar. (2003). Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328.
- Mathiowetz, N., Brown, C., and Bound, J. (2002). Measurement error in surveys of the low-income population. *Studies of welfare populations: Data collection and research issues*, pages 157–194.
- McDonald, J. B. (1984). Some generalized functions for the size distribution of income. *Econometrica*, 52(3):647–663.
- McDonald, J. B. and Ransom, M. (2008). The generalized beta distribution as a model for the distribution of income: estimation of related measures of inequality. In Chotikapanich, D., editor, *Modeling Income Distributions and Lorenz Curves*, pages 147–166. Springer.
- Meyer, B. D. and Mittag, N. (2019). Using linked survey and administrative data to better measure income: Implications for poverty, program effectiveness, and holes in the safety net. *American Economic Journal: Applied Economics*, 11(2):176–204.
- Moore, J. C., Stinson, L. L., and Welniak, E. J. (2000). Income measurement error in surveys: A review. *Journal of Official Statistics-Stockholm-*, 16(4):331–362.
- Nandram, B. and Choi, J. W. (2002). Hierarchical Bayesian nonresponse models for binary data from small areas with uncertainty about ignorability. *Journal of the American Statistical Association*, 97(458):381–388.
- Pedace, R. and Bates, N. (2000). Using administrative records to assess earnings reporting error in the survey of income and program participation. *Journal of Economic and Social Measurement*, 26(3-4):173–192.
- Peters, G. and Sisson, S. A. (2006). Bayesian inference, Monte Carlo sampling and operational risk. *Journal of Operational Risk*, 1(3):27–50.
- Pickands, J. (1975). Statistical inference using extreme order statistics. *The Annals of Statistics*, 3:119–131.
- Pissarides, C. A. and Weber, G. (1989). An expenditure-based estimate of Britain’s black economy. *Journal of public economics*, 39(1):17–32.
- Ransom, M. R. and Cramer, J. S. (1983). Income distribution functions with disturbances. *European Economic Review*, 22(3):363–372.
- Ratmann, O. R. (2010). *Approximate Bayesian Computation under model uncertainty*. PhD thesis, Imperial College London (University of London).
- Rubin, D. B. (1976). Inference and missing data. *Biometrika*, 63(3):581–592.
- Rubin, D. B. (1977). Formalizing subjective notions about the effect of nonrespondents in sample surveys. *Journal of the American Statistical Association*, 72(359):538–543.
- Shorrocks, A. F. (1983). Ranking income distributions. *Economica*, 50(197):3–17.



- Silva, M. (2023). Parametric estimation of income distributions using grouped data: an Approximate Bayesian computation approach. *Aix-Marseille School of Economics (AMSE) Working Paper*, 2023(10).
- Van Praag, B., Hagenaaars, A., and van Eck, W. (1983). The influence of classification and observation errors on the measurement of income inequality. *Econometrica*, 51(4):1093–1108.
- Wirth, H. and Pforr, K. (2022). The European Union statistics on income and living conditions after 15 years. *European Sociological Review*, 38(5):832–848.
- Yu, B. and Mykland, P. (1998). Looking at Markov samplers through cusum path plots: a simple diagnostic idea. *Statistics and Computing*, 8:275–286.

## Appendix A

### Comparing income distributions through the Wasserstein distance.

When microdata on a sample of incomes  $\{y_{(i)}\}_{i=1}^n$  is available, we can estimate empirical quantiles of the income distribution and make inference on a parametric model's parameters for the population income distribution by fitting the model at these quantiles. This is so because quantiles are informative on both shape and scale of the distribution. This allows for an ABC approach without the need of summarizing the data through a small set of summary statistics.

The Wasserstein distance, originally developed in [Kantorovich \(1939\)](#), was recently explored for the purpose of summary-free ABC inference in [Bernton et al. \(2019\)](#) and [Drovandi and Frazier \(2022\)](#). The distance between an income distribution  $f_y$  with quantile function  $F_y^{-1}$  and a parametric distribution model for this distribution  $f_y(\cdot; \theta)$  with quantile function  $F_y^{-1}(\cdot; \theta)$  follows in this case:

$$W_p(f_y, f_y(\cdot; \theta)) = \left( \int_0^1 d\{F_y^{-1}(\lambda), F_y^{-1}(\lambda; \theta)\}^p d\lambda \right)^{\frac{1}{p}}.$$

In the case of  $p = 1$  and  $d(x, y) = |x - y|$  this can be consistently estimated from the sample of incomes  $\{y_{(i)}\}_{i=1}^n$  with empirical distribution  $\hat{f}_y$  and a simulated sample of equal size  $\{\tilde{y}_{(i)}\}_{i=1}^n$  from the model with empirical distribution  $\hat{f}_y(\cdot; \theta)$  as (e.g., [Berthet et al. 2020](#)):

$$W_1(\hat{f}_y, \hat{f}_y(\cdot; \theta)) = \frac{1}{n} \sum_{i=1}^n |y_{(i)} - \tilde{y}_{(i)}|.$$

This latter formulation can be interpreted as a metric comparing all sample order statistics (i.e., all sample quantiles). In essence, this metric estimates the average absolute difference between quantiles of the two distributions.

A metric  $\rho(\hat{f}_y, \hat{f}_y(\cdot; \theta))$  may be specified under a similar logic comparing the first-order differences of all empirical *GLC* coordinates (i.e., estimates of the quantiles by definition of the GLC) between these microdata samples instead of order statistics directly. This amounts simply to taking the Wasserstein distance as defined above:

$$\begin{aligned} \rho(\hat{f}_y, \hat{f}_y(\cdot; \theta)) &= \sum_{i=1}^n |(GLC(\hat{F}_y(y_{(i)})) - GLC(\hat{F}_y(y_{(i-1)}))) - (GLC(\hat{F}_y(\tilde{y}_{(i)}; \theta)) - GLC(\hat{F}_y(\tilde{y}_{(i-1)}; \theta)))| \\ &= \sum_{i=1}^n \left| \left( \frac{y_{(i)}}{\sum_{i=1}^n y_{(i)}} \right) \times \left( \frac{\sum_{i=1}^n y_{(i)}}{n} \right) - \left( \frac{\tilde{y}_{(i)}}{\sum_{i=1}^n \tilde{y}_{(i)}} \right) \times \left( \frac{\sum_{i=1}^n \tilde{y}_{(i)}}{n} \right) \right| \\ &= \sum_{i=1}^n \frac{|y_{(i)} - \tilde{y}_{(i)}|}{n} \\ &= W_1(\hat{f}_y, \hat{f}_y(\cdot; \theta)), \end{aligned}$$

with  $\hat{F}_y(\cdot)$  and  $\hat{F}_y(\cdot; \theta)$  denoting the corresponding empirical CDFs for observed and simulated samples respectively. This result supports the use of the Wasserstein-1

distance as a common unidimensional discrepancy allowing for ABC inference either with microdata or grouped data summarized through the *GLC*.

Under grouped data in the form of  $K$  groups' mean incomes, where observed incomes are split into  $K$  segments of sizes  $n_k$  with bounds  $z_{(k)}$ ,  $k = 0, \dots, K+1$ ,  $z_{(0)} \equiv 0$ ,  $z_{(K+1)} \equiv \infty$  and a mean income  $\bar{y}_k$  is provided for each group, we can compute a grouped-data Wasserstein-1 distance from the corresponding empirical GLC coordinates for each group in the observed data  $GLC_k$  and in simulated samples  $GLC_k(\theta)$  as:

$$\begin{aligned}
\rho(\hat{f}_y, \hat{f}_y(\cdot; \theta)) &= \sum_{k=1}^K |(GLC_k - GLC_{k-1}) - (GLC_k(\theta) - GLC_{k-1}(\theta))| \\
&= \sum_{k=1}^K \left| \left( \frac{\bar{y}_{(k)} \times n_{(k)}}{\sum_{k=1}^K \bar{y}_{(k)} \times n_{(k)}} \right) \times \left( \frac{\sum_{k=1}^K \bar{y}_{(k)} \times n_{(k)}}{K} \right) - \left( \frac{\tilde{\bar{y}}_{(k)} \times n_{(k)}}{\sum_{k=1}^K \tilde{\bar{y}}_{(k)} \times n_{(k)}} \right) \times \left( \frac{\sum_{k=1}^K \tilde{\bar{y}}_{(k)} \times n_{(k)}}{K} \right) \right| \\
&= \frac{1}{K} \sum_{k=1}^K |\bar{y}_{(k)} \times n_{(k)} - \tilde{\bar{y}}_{(k)} \times n_{(k)}| \\
&= \frac{1}{K} \sum_{k=1}^K |(\bar{y}_{(k)} - \tilde{\bar{y}}_{(k)}) \times n_{(k)}| \\
&= \frac{1}{K} \sum_{k=1}^K \left| \sum_{i=1}^n (y_{(i)} - \tilde{y}_{(i)}) \times I(z_{(k)} \geq y_{(i)} \geq z_{(k-1)}) \times I(z_{(k)} \geq \tilde{y}_{(i)} \geq z_{(k-1)}) \right|,
\end{aligned}$$

where  $\tilde{\bar{y}}_{(k)}$  denotes the  $k$ -th group's mean income in the simulated sample of incomes and which, in the trivial case of having  $K = n$  groups (i.e., one observation per group) corresponds to the expression for this distance on microdata. These results suggest that in the case of grouped data we can exploit the discrepancies between *GLC* curves through their first-order difference (i.e., through the approximation to the Wasserstein-1 distance) and proceed analogously to the microdata case.

Geometrically, the Wasserstein-1 distance computes the average absolute difference between the quantile functions of two distributions. When only grouped data is available, this average distance is approximated in the expressions above by first computing the area between both empirical quantile curves within each interval of the grouped data, summing these areas across all intervals and dividing by the number of intervals. The approximation comes the fact that in computing these areas the curves might cross within an interval and so we would have no way of accounting for those differences which counteract within the interval (i.e., the absolute value is applied at the interval level in the grouped data expression above). For a same population size  $n$ , however, the quality of the approximation always increases with  $K$ .

Having access to microdata allows a computationally-cheap alternative in which user-specified groups or bins can be defined for exploiting the grouped-data approximation to the Wasserstein-1 distance. For instance, instead of grouping the data on sample deciles, one could define broader groups for lower incomes and finer groups for higher incomes allowing for a particularly stricter fit on the upper tail of the distribution.

## Deriving parametric distribution functions with high-income under-reporting and non-response components

Assuming that microdata samples from a population's GB2 income distribution may be jointly affected by high-income under-reporting following an LPU scheme with parameters  $(\bar{p}, \delta)$  and high-income non-response following a right-truncation scheme with  $\alpha$  fixed to  $\alpha = 1$  and parameter  $t$  where  $t \gg \bar{p}$ , then under this joint scheme:

$$\begin{aligned} \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} &= \frac{\partial [y_i^{Obs} \times [1 + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})) \times \left(\frac{\delta}{1-\delta}\right)]]}{\partial y_i^{Obs}} \\ &= \frac{1}{1 - \delta \times \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \alpha, \beta, p, q))} , \end{aligned}$$

$$\begin{aligned} \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t) &= \begin{cases} 1 , & \text{if } F_{\mathbf{y}}^{GB2} \left( y_i^{Obs} + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})) \times \left( \frac{\delta \times (y_i^{Obs} - F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta}))}{1-\delta} \right) \right) \leq t \\ 0 , & \text{otherwise} \end{cases} \\ &= \begin{cases} 1 , & \text{if } y_i^{Obs} + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})) \times \left( \frac{\delta \times (y_i^{Obs} - F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta}))}{1-\delta} \right) \leq F_{\mathbf{y}}^{-1;GB2}(t; \boldsymbol{\theta}) \\ 0 , & \text{otherwise} \end{cases} \\ &= \begin{cases} 1 , & \text{if } y_i^{Obs} \leq (1 - \delta) F_{\mathbf{y}}^{-1;GB2}(t; \boldsymbol{\theta}) + \delta F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta}) \\ 0 , & \text{otherwise} \end{cases} , \end{aligned}$$

and

$$\int f_{\mathbf{y}}^{GB2}(m^{-1}(y_i^{Obs}; \bar{p}, \delta); \alpha, \beta, p, q) \times \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} \right) dy_i^{Obs} = t ,$$

allow for stating a model for the observable data  $\mathbf{y}^{Obs}$  applying (1):

$$\begin{aligned} f_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \bar{p}, \delta, t) &= \frac{f_{\mathbf{y}}^{GB2}(m^{-1}(y_i^{Obs}; \bar{p}, \delta); \boldsymbol{\theta}) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} \right) \times \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t)}{\int f_{\mathbf{y}}^{GB2}(m^{-1}(y_i^{Obs}; \bar{p}, \delta); \boldsymbol{\theta}) \times \left( \frac{\partial m^{-1}(y_i^{Obs}; \bar{p}, \delta)}{\partial y_i^{Obs}} \right) \times \varphi(m^{-1}(y_i^{Obs}; \bar{p}, \delta); t) dy_i^{Obs}} \\ &= \frac{f_{\mathbf{y}}^{GB2} \left( y_i^{Obs} + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})) \times \left( \frac{\delta(y_i^{Obs} - F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta}))}{1-\delta} \right); \boldsymbol{\theta} \right)}{t \times (1 - \delta \times \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})))} \\ &\quad \times \frac{\mathbf{1}(y_i^{Obs} \leq (1 - \delta) F_{\mathbf{y}}^{-1;GB2}(t; \boldsymbol{\theta}) + \delta F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta}))}{t \times (1 - \delta \times \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1;GB2}(\bar{p}; \boldsymbol{\theta})))} . \end{aligned}$$

Additionally, in this setting the population-level CDF can be stated as:

$$F_{\mathbf{y}^{Obs}}(y_i^{Obs}; \boldsymbol{\theta}, \bar{p}, \delta, t) = F_{\mathbf{y}} \left( y_i^{Obs} + \mathbf{1}(y_i^{Obs} > F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta})) \times \left( \frac{\delta(y_i^{Obs} - F_{\mathbf{y}}^{-1}(\bar{p}; \boldsymbol{\theta}))}{1 - \delta} \right) \right) \times \frac{1}{t} .$$

## Simulated data applications: further results

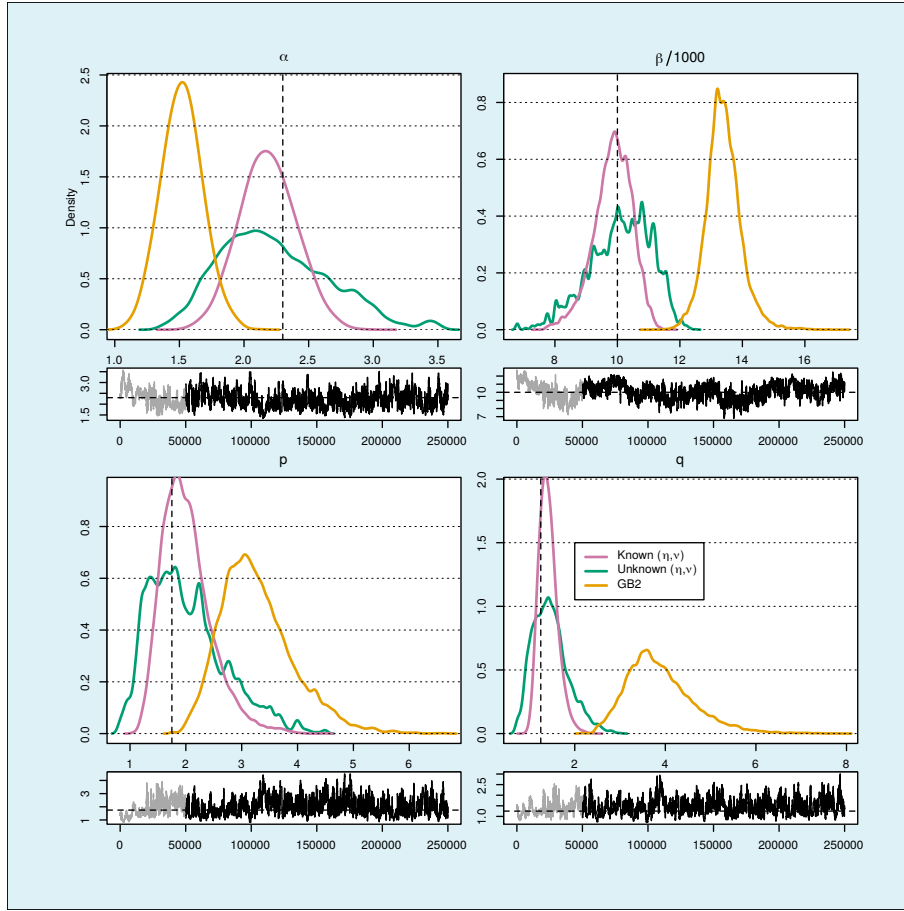


Figure 8: Estimated ABC marginal posterior distributions for income distribution parameters

**Note:** Kernel density estimates for ABC marginal posterior distribution estimates computed on a single simulated sample of  $N = 9900$  observed incomes following (4) under parameter values  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Estimates obtained applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a GB2), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases, the algorithm is set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained ( $\{\theta^{(j)}\}_{j=1}^{250000}$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period. Traceplots of the underlying MCMC samples for estimates with uncertainty on  $(\eta, \nu)$  below, with burn-in period in gray. True parameter values in dashed black lines.

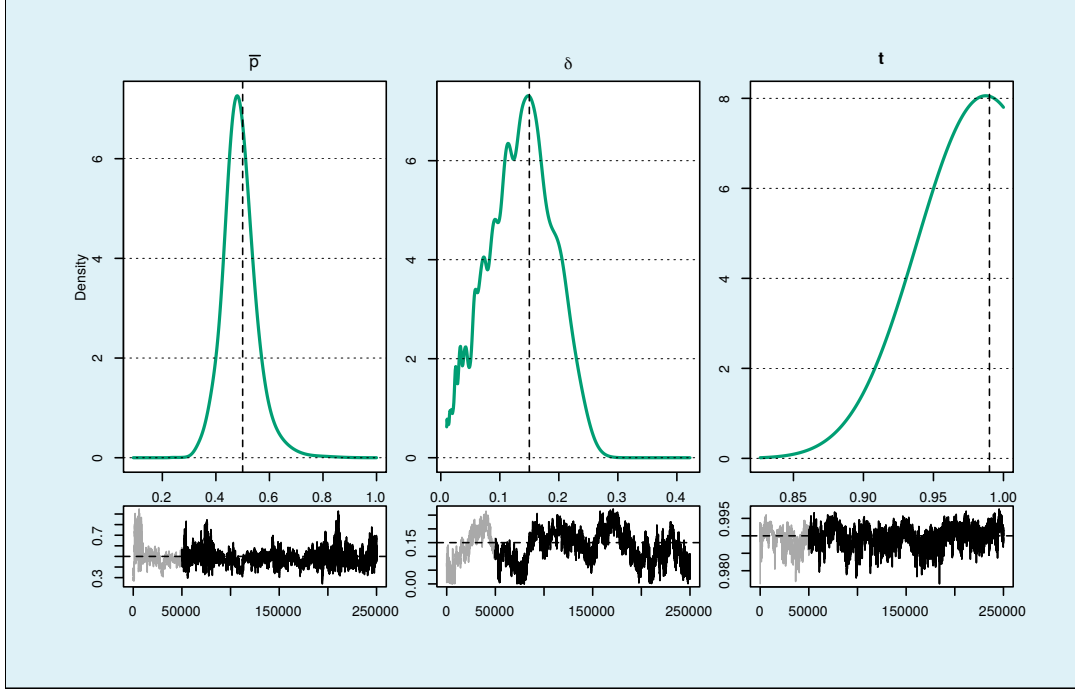


Figure 9: Estimated ABC marginal posterior distributions for  $MR$  parameters **Note:** Kernel density estimates for ABC marginal posterior distribution estimates computed on a single simulated sample of  $N = 9900$  observed incomes following (4) under parameter values  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Estimates obtained applying the (ABC-AM) algorithm with prior uncertainty on the MR parameters  $(\eta, \nu) = (\bar{p}, \delta, t)$ . The algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained, taking the initial  $J_0 = 50000$  draws as the burn-in period. Traceplots of the underlying MCMC samples below, with burn-in period in gray. True parameter values in dashed black lines.



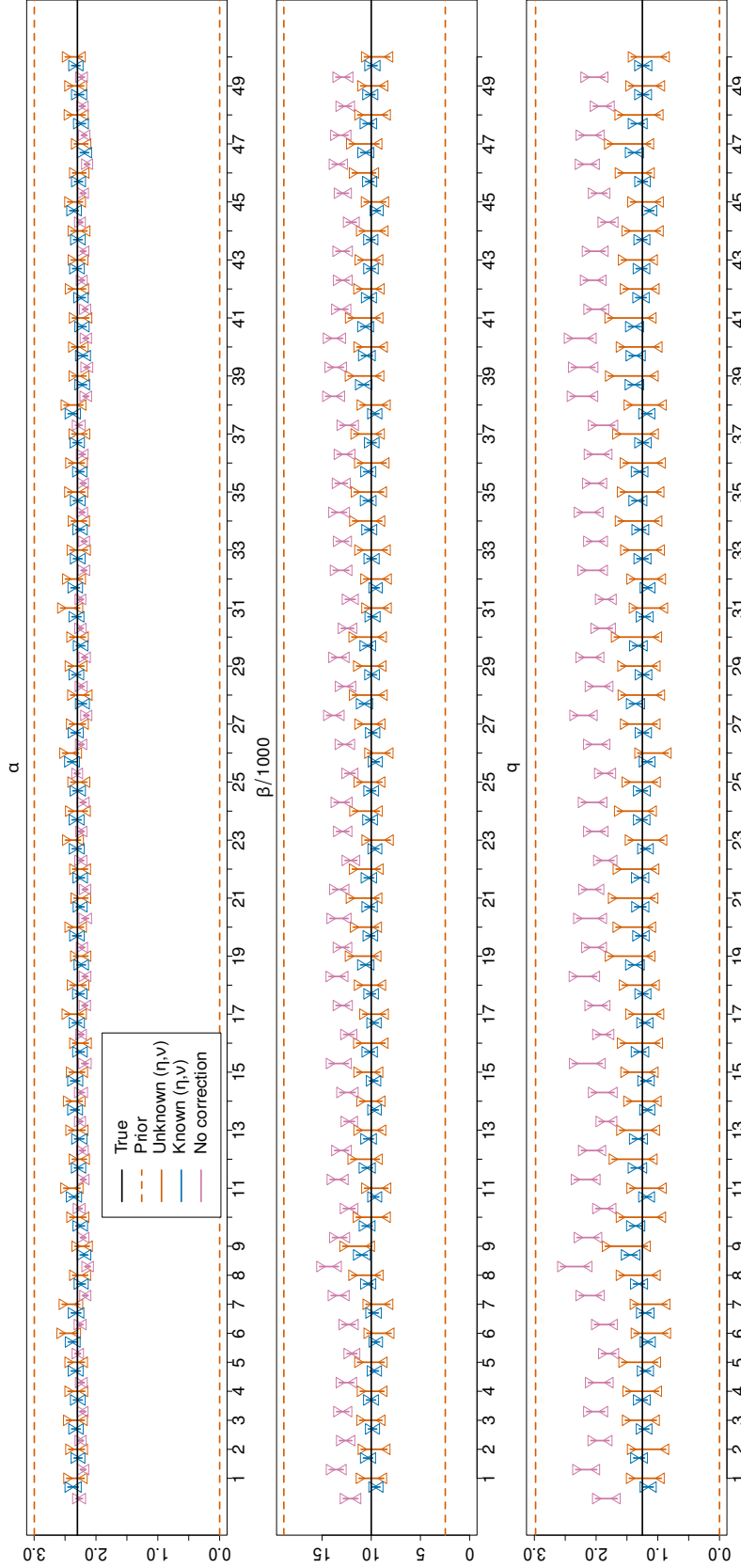


Figure 10: ABC parameter estimates under simulated data setup I: Population income distribution parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with  $p$  fixed at  $p = 1$ , corresponding to a Singh-Maddala population income distribution, with parameter values  $(\alpha, \frac{\beta}{1000}, q) = (2.3, 10, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Posterior distribution estimates obtained under model (4) with the constraint  $p = 1$  applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a simple Singh-Maddala), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases the algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total  $\{(\theta^{(j)})\}_{j=1}^{250000}$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.

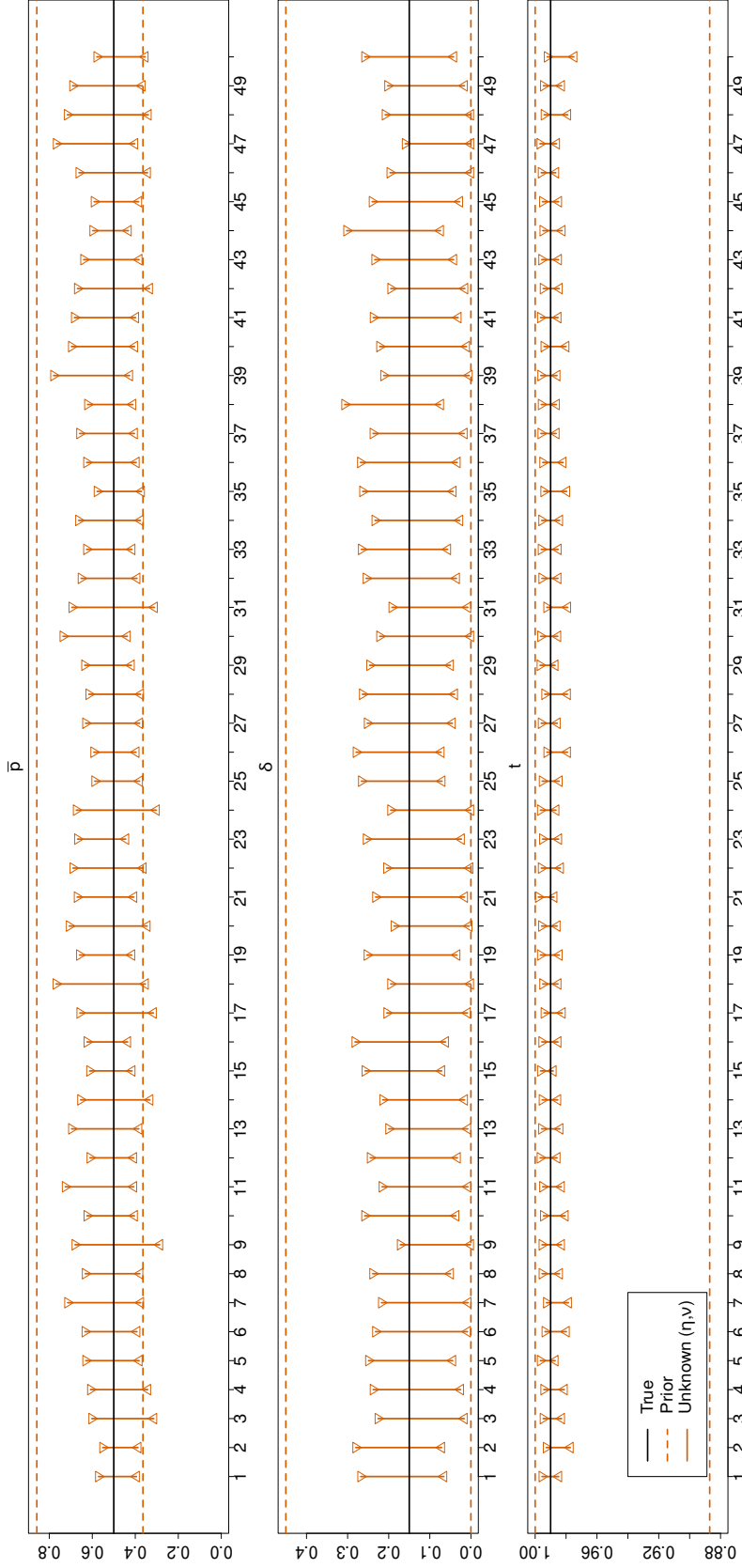


Figure 11: ABC parameter estimates under simulated data setup I: MR parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with  $p$  fixed at  $p = 1$ , corresponding to a Singh-Maddala population income distribution, with parameter values  $(\alpha, \frac{\beta}{1000}, q) = (2.3, 10, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Posterior distribution estimates obtained under model (4) with the constraint  $p = 1$  applying the (ABC-AM) algorithm with prior uncertainty on the MR parameters  $(\bar{p}, \delta, t)$ . The algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total, taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.

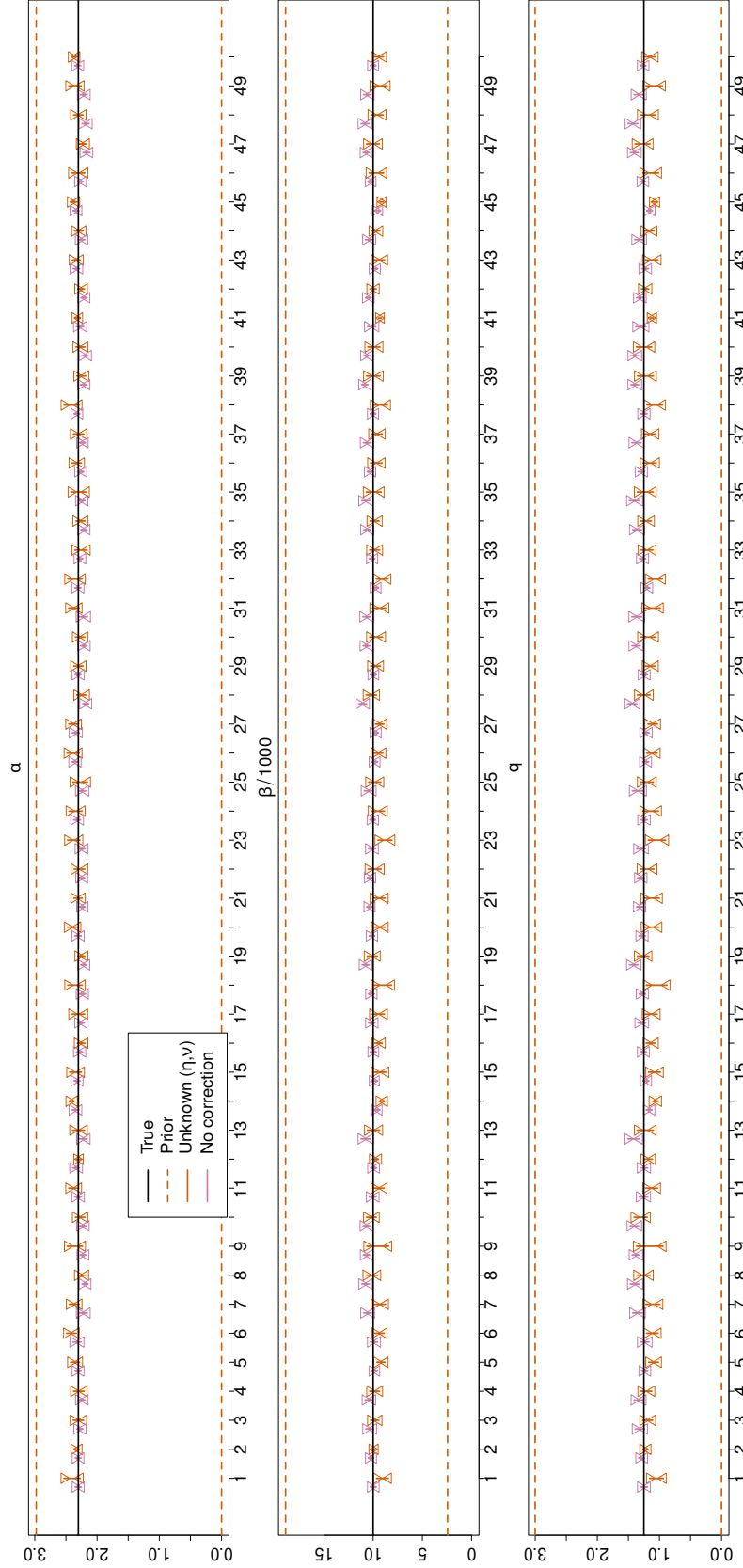


Figure 12: ABC parameter estimates under simulated data setup *II*: Population income distribution parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with  $p$  fixed at  $p = 1$ , corresponding to a Singh-Maddala population income distribution, with parameter values  $(\alpha, \frac{\beta}{1000}, q) = (2.3, 10, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (1, 0, 1)$  introducing no MR issues on the observed incomes. Posterior distribution estimates obtained under model (4) with the constraint  $p = 1$  applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a simple Singh-Maddala), and with prior uncertainty on the MR parameters. In all cases the algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(0.1, 1, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total  $\{\theta^{(j)}\}_{j=1}^{250000}$  in the first case), taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.

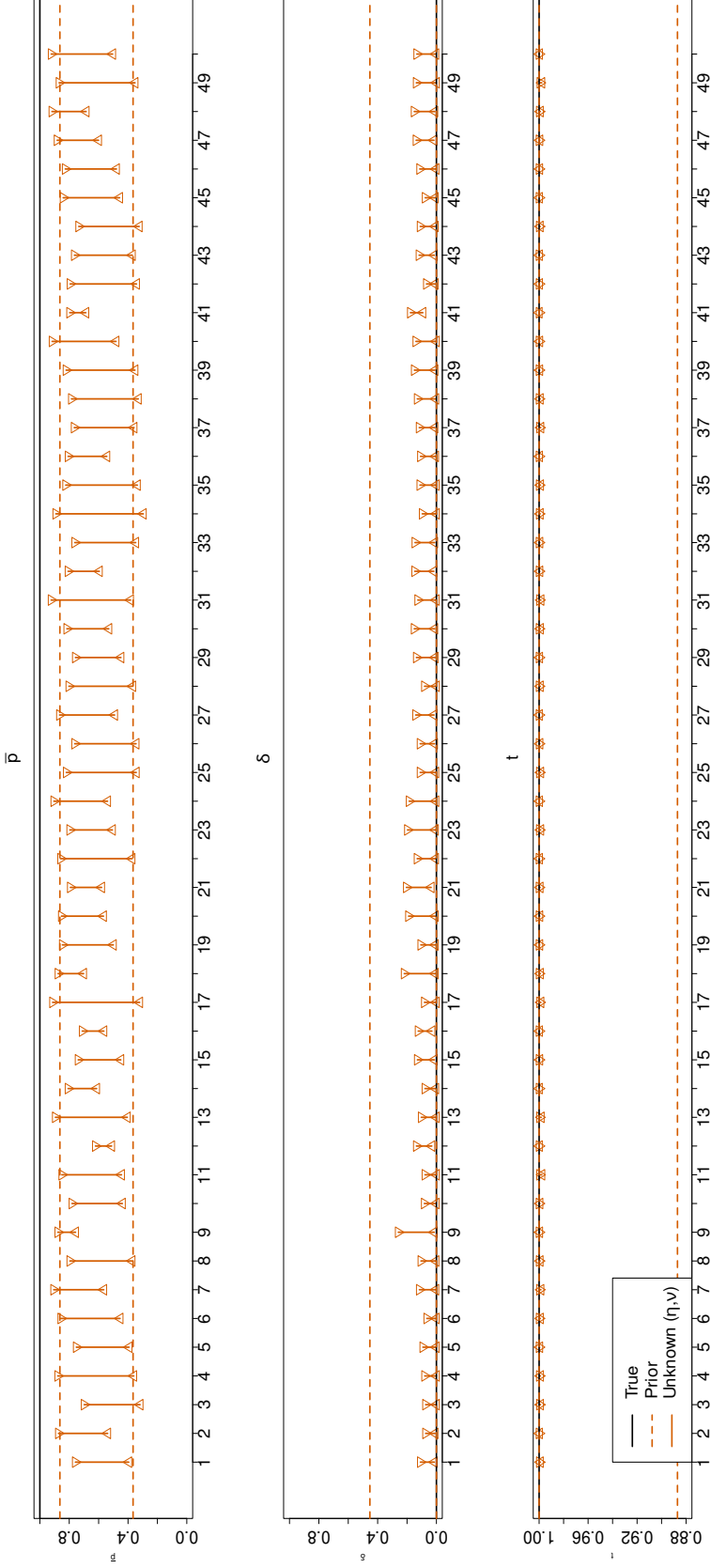


Figure 13: ABC parameter estimates under simulated data setup II: MR parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with  $p$  fixed at  $p = 1$ , corresponding to a Singh-Maddala population income distribution, with parameter values  $(\alpha, \frac{\beta}{1000}, q) = (2.3, 10, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (1, 0, 1)$  introducing no MR issues on the observed incomes. Posterior distribution estimates obtained under model (4) with the constraint  $p = 1$  applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a simple Singh-Maddala), and with prior uncertainty on the MR parameters. In all cases the algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(0.1, 1, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total  $\{\theta^{(j)}\}_{j=1}^{250000}$  in the first case), taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.

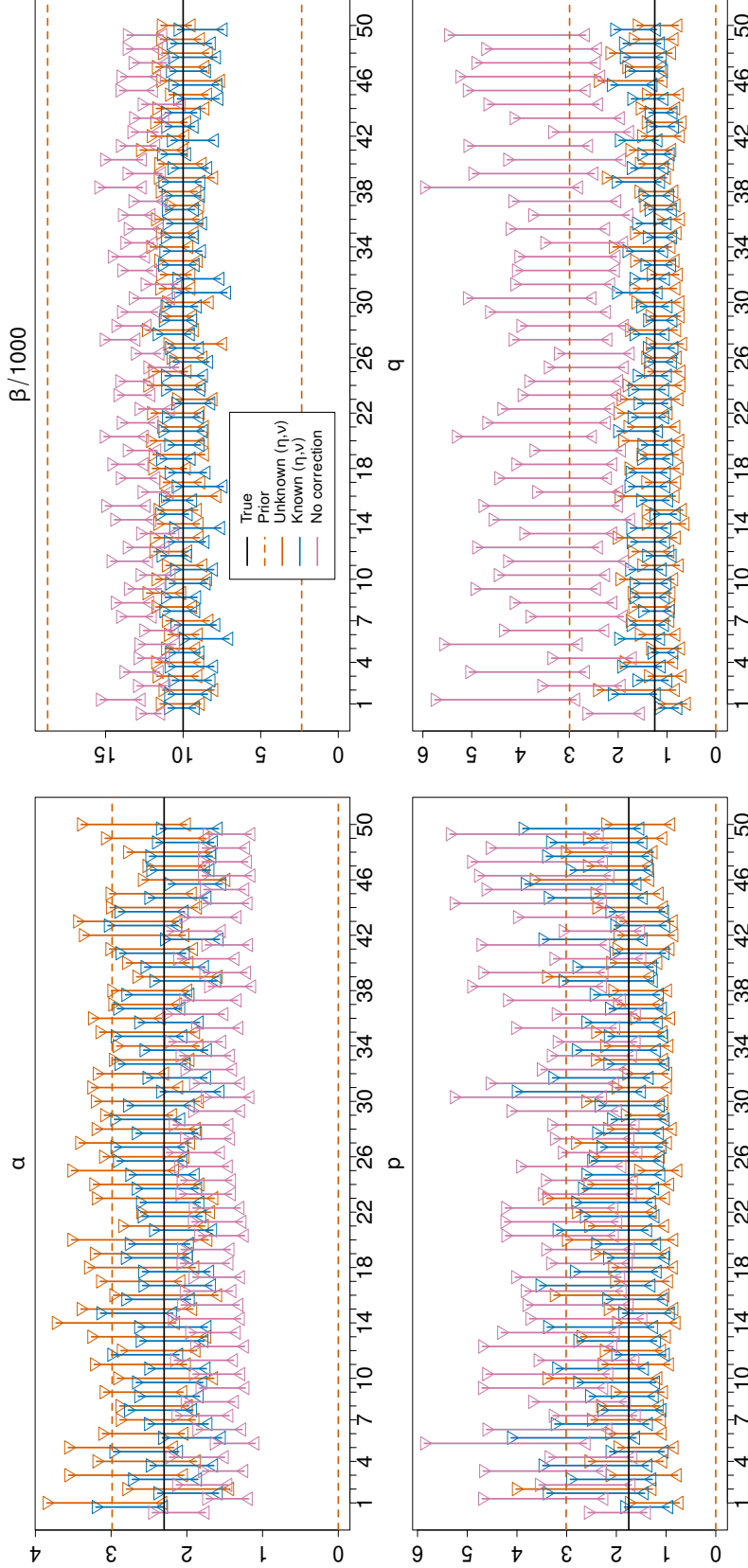


Figure 14: ABC parameter estimates under simulated data setup *III*: Population income distribution parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with parameter values  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Posterior distribution estimates obtained under model (4) applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a GB2), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases the algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total  $(\{\theta^{(j)}\}_{j=1}^{250000})$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.

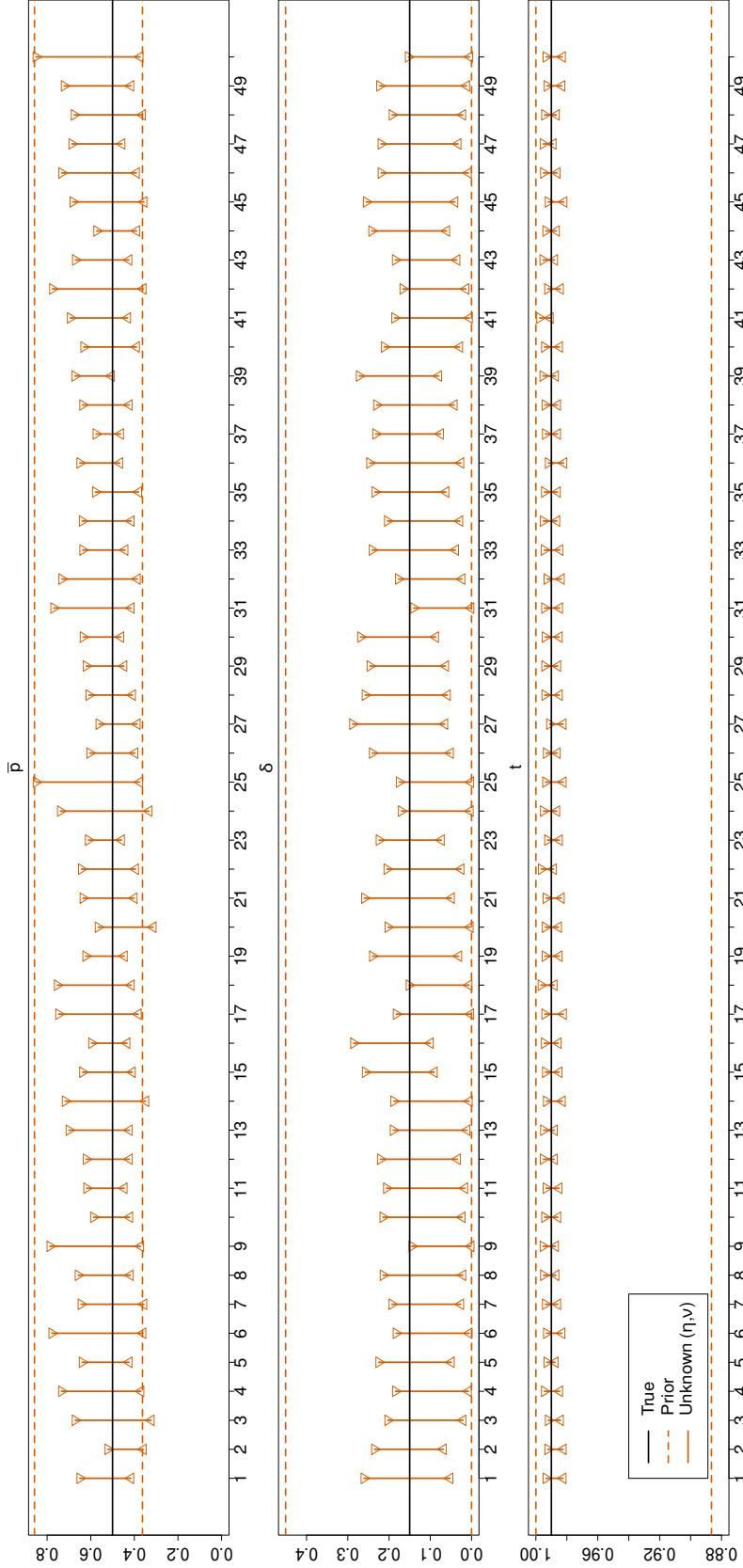


Figure 15: ABC parameter estimates under simulated data setup III: MR parameters.

**Note:** 95% HPDI for ABC parameter estimates over 50 simulated data samples of size  $N/t = 10000$  following (4) with parameter values  $(\alpha, \frac{\beta}{1000}, p, q) = (2.3, 10, 1.75, 1.25)$  and  $(\eta, \nu) = (\bar{p}, \delta, t) = (.5, .15, .99)$ . Posterior distribution estimates obtained under model (4) applying the (ABC-AM) algorithm separately: without MR corrections (i.e.,  $(\bar{p}, \delta, t) = (1, 0, 1)$  corresponding to a GB2), conditional on the true  $(\bar{p}, \delta, t) = (.5, .15, .99)$  correction parameters, and with prior uncertainty on these. In all cases the algorithm was set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\{(\theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)})\}_{j=1}^{250000}$  were obtained in total  $(\{\theta^{(j)}\}_{j=1}^{250000})$  in the first two cases), taking the initial  $J_0 = 50000$  draws as the burn-in period before adaptation and discarding an additional 100000 initial observations. 95% HPDI from the parameter's prior distribution in dashed orange lines. True parameter values in solid black lines.



# Appendix B

## EU-SILC application.

Table 2: ABC posterior distribution estimates for selected EU-SILC samples under (4)

Country	Wave	$\theta$				$\eta$		$\nu$
		$\alpha$	$\frac{\beta}{1000}$	$p$	$q$	$\bar{p}$	$\delta$	$t$ (%)
Austria (AT)	2005	4.4926 [3.4804;5.4754]	17.4491 [16.4642;18.3747]	0.7689 [0.5323;1.0381]	0.6995 [0.4948;0.887]	0.6481 [0.5588;0.7402]	0.1394 [0.0604;0.219]	99.9212 [99.8606;99.9774]
	2007	4.2329 [3.3264;5.6431]	17.497 [16.3194;18.3236]	0.8706 [0.5311;1.1924]	0.7848 [0.5149;0.9549]	0.6917 [0.5881;0.809]	0.1654 [0.0536;0.2412]	99.9987 [99.996;100]
	2011	4.4908 [3.6792;5.4883]	22.5277 [21.8214;23.2466]	0.6266 [0.4651;0.7952]	0.7225 [0.5291;0.8909]	0.6646 [0.5484;0.7882]	0.1145 [0.0464;0.1831]	99.9475 [99.9076;99.9806]
	2016	4.1521 [3.2317;5.3752]	23.2533 [21.3724;24.4954]	0.7947 [0.5194;1.1763]	0.76 [0.5472;0.9347]	0.6726 [0.6233;0.7313]	0.2252 [0.1611;0.2871]	99.9989 [99.9969;100]
Germany (DE)	2005	4.1816 [3.243;5.0834]	15.6159 [14.7938;16.395]	0.8099 [0.5691;1.076]	0.7151 [0.5477;0.8938]	0.5401 [0.4771;0.6122]	0.2318 [0.1642;0.304]	99.9992 [99.9974;100]
	2007	5.1701 [4.464;5.8984]	17.1657 [16.7115;17.6088]	0.5078 [0.4107;0.6081]	0.4937 [0.4123;0.5753]	0.7639 [0.7268;0.7959]	0.2838 [0.2277;0.3407]	99.999 [99.9966;100]
	2011	4.2773 [3.5422;5.1335]	18.1032 [17.45;18.8769]	0.6265 [0.4699;0.8126]	0.5993 [0.4394;0.7026]	0.7678 [0.7408;0.7945]	0.3997 [0.3482;0.4484]	99.9995 [99.9982;100]
	2016	4.3893 [3.8794;4.8718]	20.3698 [19.8252;20.9009]	0.5862 [0.4868;0.6795]	0.5884 [0.5204;0.6677]	0.773 [0.7527;0.7889]	0.3453 [0.3142;0.3776]	99.9995 [99.9985;100]
France (FR)	2005	2.9131 [2.3348;3.4636]	13.4956 [11.8554;14.9544]	1.6066 [1.0129;2.3224]	1.1055 [0.8676;1.3669]	0.6539 [0.573;0.726]	0.1383 [0.0741;0.2019]	99.9923 [99.9772;100]
	2007	3.9296 [3.2279;4.6694]	15.6607 [15.1527;16.1589]	0.9887 [0.7278;1.2634]	0.8818 [0.6438;1.1277]	0.8525 [0.7985;0.9022]	0.1446 [0.0458;0.249]	99.9892 [99.9691;100]
	2011	4.4173 [3.2817;5.1189]	17.8792 [16.9069;18.7775]	0.8398 [0.5991;1.1511]	0.6043 [0.464;0.748]	0.5919 [0.4136;0.716]	0.0779 [0.011;0.1361]	99.9853 [99.9702;99.9978]
	2016	5.4397 [2.9187;8.3704]	18.4294 [16.2994;20.3333]	0.9048 [0.3647;1.6694]	0.5048 [0.2796;0.7936]	0.5747 [0.4268;0.7176]	0.2963 [0.2481;0.3755]	99.9996 [99.9988;100]
Spain (ES)	2005	1.6968 [1.2934;2.1302]	11.0949 [9.2004;12.9569]	2.2101 [1.2547;3.3318]	2.3738 [1.4735;3.414]	0.7627 [0.7141;0.8138]	0.2002 [0.1207;0.2806]	99.9714 [99.9194;100]
	2007	1.6921 [1.2519;2.1322]	12.7219 [10.3523;14.8156]	2.2679 [1.2379;3.5877]	2.5442 [1.6421;3.6412]	0.7057 [0.637;0.7789]	0.1481 [0.0696;0.2237]	99.9733 [99.9272;100]
	2011	2.5284 [2.1636;2.9317]	14.8681 [13.9019;15.8051]	1.0773 [0.8016;1.3675]	1.2021 [0.9435;1.4661]	0.8243 [0.7939;0.856]	0.2938 [0.2314;0.3629]	99.9946 [99.983;100]
	2016	2.5723 [2.2547;2.9011]	17.8008 [16.8279;18.7658]	0.8678 [0.7037;1.0366]	1.4029 [1.1189;1.6883]	0.82 [0.7653;0.8743]	0.1248 [0.0639;0.191]	99.9909 [99.9731;100]
Italy (IT)	2005	3.4764 [2.8887;4.0841]	14.038 [13.2232;14.817]	0.7882 [0.5732;1.0109]	0.7659 [0.613;0.9184]	0.6818 [0.6362;0.7277]	0.2253 [0.165;0.2872]	99.9889 [99.9745;100]
	2007	2.9708 [2.3382;3.6742]	14.818 [13.5246;16.1084]	0.9931 [0.6399;1.4089]	0.972 [0.7172;1.2408]	0.6871 [0.6274;0.7446]	0.2037 [0.132;0.2773]	99.9596 [99.9203;99.9933]
	2011	4.1234 [3.4491;4.8556]	17.1416 [16.4914;17.757]	0.5719 [0.4288;0.7125]	0.6735 [0.5503;0.7979]	0.6788 [0.6229;0.7496]	0.1961 [0.1454;0.2515]	99.9985 [99.9946;100]
	2016	3.3105 [2.319;4.7316]	19.1832 [18.358;20.2492]	0.7476 [0.381;1.1051]	0.9856 [0.5918;1.3342]	0.5769 [0.5323;0.6323]	0.1915 [0.0879;0.2723]	99.9738 [99.9444;99.9994]

**Source:** Own calculations from EU-SILC.

**Note:** ABC posterior mean estimates of all parameters of model (4) for Austria (AT), Germany (DE), France (FR), Spain (ES), and Italy (IT) from 2005, 2007, 2011, and 2016 EU-SILC waves. Only considers households with reported household disposable income (HX090) of at least 1 euro. Respective 95% HPDI in brackets. Estimates obtained applying the (ABC-AM) algorithm in all cases, set with parameters  $\tau = 25$ ,  $M = 10$ ,  $\Sigma^{(0)} = \text{diag}(.01, .1, .01, .01, .01, .01, .01, .01)$ , and  $J = 250000$  MCMC samples  $\left\{ \left( \theta^{(j)}, \bar{p}^{(j)}, \delta^{(j)}, t^{(j)} \right) \right\}_{j=1}^{250000}$  were obtained, taking the initial  $J_0 = 50000$  draws as the burn-in period.