

# STICERD

Suntory and Toyota International Centres  
for Economics and Related Disciplines

Distributional Analysis Research Programme  
Discussion Paper

**Robustness Properties of Inequality  
Measures: The Influence Function and  
the Principle of Transfers**

**Frank A Cowell and Maria-Pia Victoria-Feser**  
October 1993

LSE STICERD Research Paper No. DARP 01

This paper can be downloaded without charge from:

<http://sticerd.lse.ac.uk/dps/darp/darp1.pdf>

# ROBUSTNESS PROPERTIES OF INEQUALITY MEASURES

by

Frank Cowell and Marica—Pia Victoria-Feser<sup>1</sup>  
London School of Economics and Political Science

Discussion Paper  
No.DARP/1  
October 1993

The Toyota Centre  
Suntory and Toyota International Centres for  
Economics and Related Disciplines  
London School of Economics and Political Science  
Houghton Street  
London WC2A 2AE  
Tel.: 020-7955 6678

---

<sup>1</sup> Partially supported by the 'Fond National Suisse pour la Recherche Scientifique'. The authors would like to thank C Dagum, S Howes, S Jenkins, P Lambert and S Yitzhaki for their comments on earlier versions.

## Abstract

Inequality measures are often used to summarise information about empirical income distributions. However, the resulting picture of the distribution and of changes in the distribution can be severely distorted if the data are contaminated. The nature of this distortion will in general depend upon the underlying properties of the inequality measure. We investigate this issue theoretically using a technique based on the influence function, and illustrate the magnitude of the effect using a simulation. We consider both direct nonparametric estimation from the sample, and indirect estimation using a parametric model. In the latter case we demonstrate the application of a robust estimation procedure.

**Keywords:** Inequality measurement, transfer principle, influence function, robust estimation.

**JEL Nos.:** C13 D63.

© by the authors. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Contact address: Frank Cowell, STICERD, London School of Economics and Political Science, Houghton Street, London WC2A 2AE, email: [f.cowell@lse.ac.uk](mailto:f.cowell@lse.ac.uk)

# 1 Introduction

This paper is about the robustness properties of the estimators of inequality measures and other related tools of income distribution analysis.

There are various approaches in the literature to the definition of an inequality measure. For our purposes we see an inequality measure simply as a member of a class of functions that is defined by a set of essential characteristics. The class of functions defines in turn a set of statistics that can be used to characterise distributions of income, wealth, and so on.

The essential characteristics of an inequality measure are the subject of some debate. However, many would accept that the principal property which should be possessed by an inequality measure is respect for the principle of transfers (Dalton 1920). In addition it is often required that the class of admissible inequality measures satisfy properties such as scale independence or decomposability.

There is evidently also room for debate about what constitutes a “good” estimator of a statistic - in particular an inequality statistic. The purpose of this paper is to examine the formal relationship between the economic properties that the inequality measure should fulfil and the statistical properties

that are considered appropriate for the corresponding estimator. This is important because drawing inferences about economic inequality from, for example, income distribution data plays an important part in political debates about economic and social trends, and in a variety of applied studies in the field of welfare economics. However, the statistical basis on which the inferences are drawn is not always spelt out, and so the relationship between the numbers observed in a particular sample and the supposed underlying concept of inequality within the target population may be different from that suggested by superficial appearances.

These empirical inequality measures are actually estimators of the underlying inequality: in a population from a particular set of data, one computes an estimate of the non-measurable “true” value of this inequality. We are particularly interested here in the robustness properties of these estimators. Robustness is a growing field of research in statistics that began with the pioneering paper of Huber (1964). It is a purely statistical concept which in a sense measures a “qualitative” aspect of any estimator, more precisely its stability under non-standard conditions. In particular, we will use the *Influence Function (IF)* (Hampel 1968, Hampel 1974), a measure of robustness which indicates the extent to which an estimator is influenced by an

infinitesimal amount of “errors” either in the specification of the underlying model or in the data. These errors are commonly called data *contaminations*. Moreover, for those who would argue that the “errors” can satisfactorily be determined in a pragmatic or subjective fashion, it is important to point out that the *IF* is also a tool that enables one to control for such errors systematically.

As we will see in section 2, the *IF* of an estimator is related to the bias on this estimator caused by an infinitesimal amount of data contamination: an unbounded *IF* means that this bias can be infinite. In this paper we want to show that in very general cases and for a very large class of inequality measures, we have unbounded *IF*s, meaning not only that the bias on these measures can be very large, but that this bias can actually be caused by a single observation.

The behaviour of the *IF* will depend on the type of statistic under consideration, and so a principal question which we wish to address is whether there is a systematic relationship between the properties that are taken to be defining characteristics of the class of inequality measures and the behaviour of the *IF*.

Why is this important to researchers in the field of inequality analysis? It is well known that economic data in particular are far from being clean; this usually means that some observations

may be present which in a sense have nothing to do with the majority of the data. These rogue data are usually a result of the collection procedures. A simple example is the “decimal point error”: the coder inadvertently puts the decimal point in the wrong place and thus multiplies an observation by a factor of 10. More subtle is the week-month confusion where data are supposedly collected on *weekly* income, but some respondents actually report income per month. We give an example of the effect of these kinds of contamination in subsection 5.1.

If those observations have a negligible impact upon the analysis, then obviously there is nothing to worry about. Unfortunately, in most cases, those observations are “extreme” (for example of a different order of magnitude)<sup>1</sup> and, as we will see below, they can then drive the value of the inequality measure by themselves. Such extreme values may of course be picking up true information; but very often in empirical work a case can be made for dropping an “obviously” inappropriate or suspect observation that may be the result of recording error or other contamination. This type of *ad hoc* procedure is unsatisfactory, but if it is not done then the result of the analysis may be seriously biased so that the inequality measure no

---

<sup>1</sup>In the 1981 and 1986 Family Expenditure Survey for the UK the top-most income in the sample has a value twice that of the next-highest income. It is arguable that an outlier of this sort should be treated as exceptional and dropped from the sample (see Jenkins 1992).

longer represents what it was intended to represent. In other words, the whole effort put into the construction of an inequality measure satisfying important economic properties is simply lost because the estimate deviates substantially from a representation of the “true” value of the inequality measure. For example, it is well known that Atkinson inequality indices, for values of the inequality aversion parameter greater than unity, are extraordinarily sensitive to abnormally small incomes (see for example Cowell 1977, p. 132-134, Pudney and Sutherland 1992); for this reason it has become common practice just to drop unreasonably small incomes (along with zero and negative values<sup>2</sup>) from the sample (see Jenkins 1992). What is perhaps less well recognised is that comparable problems will arise with other inequality measures that bear a family relationship to the Atkinson indices, but are extraordinarily sensitive to abnormally high incomes. However, the aim of this paper is not only to show analytically why different inequality measures are sensitive to extreme observations, but also to propose a statistical procedure for the estimation of inequality measures which avoids subjective pre-screening of the data.

We believe that robust inequality measures can provide the

---

<sup>2</sup>The Atkinson index is not defined for negative incomes. For inequality aversion greater than unity, the Atkinson index tends to its maximum value as any sample observation approaches zero.



researcher with a useful supplement of information concerning the true underlying structure of inequality for a given data sample. Indeed, the effect of data contamination in the tails of the distribution can result in serious confusion between quite different underlying income distributions, as shown in figure 1. Suppose the shape of the income distribution is represented by the continuous frequency distribution in part A of the figure, but that in a sample from the population there are some rogue observations represented by the point mass labelled “contamination”. Then, according to inequality statistics that are fairly sensitive to the top end of the distribution, the implied picture of the income distribution will be indistinguishable from that represented in part B of the figure (see Victoria-Feser 1993 for an empirical example). “Common sense” might suggest that the rogue observations be dropped from the sample, in such a situation, but “common sense” may not be an adequate guide to sound statistical practice when the data set is complex.

This paper is organized as follows: in the next section we introduce the  $IF$ . In section 3, we study the robustness properties of a general class of inequality measures and show that the principle of transfers alone is insufficient to determine whether or not the  $IF$  of these measures is bounded. We then consider the properties of an important subclass, decomposable

inequality measures: in section 4 we analyse the behaviour of this class under mean-preserving contaminations and show that under some important cases, the corresponding  $IF$  is unbounded (see proposition 1). Moreover, in section 5 we address the problem of arbitrary contaminations (where the mean is also affected by model deviations), and show that under these circumstances the  $IF$  is always unbounded. In section 6, the case of parametric models for the distribution of income is analysed through the generalized entropy class of income inequality measures. In section 7 we propose a robust method of computing inequality measures and present some simulation results. Section 8 concludes.

## 2 The influence function

The use of the  $IF$  to assess the robustness properties of any estimator was originated by Hampel (1968), Hampel (1974) and further developed in Hampel, Ronchetti, Rousseeuw, and Stahel (1986). It is defined as the influence of an infinitesimal proportion of “bad” observations on the value of the estimate.

We first introduce some notation. Let  $I$  be the (true) inequality measure that we estimate by means of a sample  $x_1, \dots, x_n$ , where the  $x_i$  are realizations of a random variable  $X$ . For example,  $X$  is the income variable when measuring income inequal-

ity. We denote by  $I(F)$  the functional version of  $I$ , depending on the true distribution  $F$ . An estimator of  $I(F)$  is obtained when replacing  $F$  by either  $F_n$ , the empirical distribution of  $X$  given by

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \Delta_x(x_i) \quad (1)$$

where  $\Delta_x$  is a point mass in  $x$ , or  $F_{\hat{\theta}}$ , an estimated parametric model, such that  $X \sim F_{\hat{\theta}}$ . Let us also define the following *mixture distribution*

$$G_\varepsilon = (1 - \varepsilon)F + \varepsilon H \quad (2)$$

where  $0 < \varepsilon < 1$  and  $H$  is some *perturbation distribution*. For example,  $H$  could be the distribution  $\Delta_z$  which puts a point mass 1 at any point  $z$ . Then  $G_\varepsilon$  is the mixture model from which an observation has probability  $(1 - \varepsilon)$  of being generated by  $F$  and a probability  $\varepsilon$  of being an arbitrary value  $z$ . In our case, as we will see later, it is more convenient to consider a distribution  $H$  with corresponding probability distribution

$$dH(x) = \begin{cases} \alpha_1 & \text{if } x = z_1 \\ \dots & \\ \alpha_m & \text{if } x = z_m \end{cases} \quad (3)$$

$\forall i, \alpha_i \geq 0$ , and  $\sum \alpha_i = 1$ .

The influence of an infinitesimal model deviation (or model contamination) on the estimate is then given by

$$\lim_{\varepsilon \rightarrow 0} \left[ \frac{I(G_\varepsilon) - I(F)}{\varepsilon} \right] \quad (4)$$

or, when the derivative exists, by

$$\left. \frac{\partial}{\partial \varepsilon} I(G_\varepsilon) \right|_{\varepsilon=0} \quad (5)$$

It should be stressed that (5) is a slightly different definition from the usual  $IF$ . Indeed, when  $H = \Delta_z$ , then it is equal to the  $IF$ , denoted by  $IF(z; I, F)$ . It gives the influence on the estimator  $I$  of an infinitesimal amount<sup>3</sup> of contamination at the point  $z$ . The  $IF$  then depends on the position of  $z$  with respect to the position of the majority of the data. Usually, the  $IF$  is maximal when  $z$  approaches extreme values such as  $\infty$ ,  $-\infty$  or  $0$ .

When  $H$  is any distribution, then (5) can be called an *integrated IF* ( $IIF$ ) because it is equal to  $\int IF(z; I, F)dH(z)$ . However, in our examples,  $H$  will be a distribution which puts different point masses at different points depending on a common point  $z$ . Therefore, by extension of the definition of the  $IF$ , we will consider that with such a distribution  $H$ , (5) is

---

<sup>3</sup>Infinitesimal means here that the probability  $\varepsilon$  that this contamination occurs tends to zero.

still an  $IF$ .

The  $IF$  is a very useful tool not only as a measure of the influence of errors in the model on the value of the estimator, but also because it can be seen as a first order approximation of the bias of the estimator when the model is misspecified as  $G_\varepsilon$  rather than  $F$ . Indeed, if we were able to compute the maximum bias  $\sup_{G_\varepsilon} |I(G_\varepsilon) - I(F)|$  as a function of  $\varepsilon$ , then  $\sup_z |IF(z; I, F)|$  would appear to be nothing else but the slope of this function at  $\varepsilon = 0$  (see figure 2 and Hampel, Ronchetti, Rousseeuw, and Stahel 1986). Therefore, we can see why it is important for an estimator to have at least a bounded  $IF$  which ensures a bounded bias, near  $\varepsilon = 0$ .

### 3 Robustness properties of a general class of inequality measures

Although the results presented in this paper apply to any inequality measure for the distribution of any random variable, we shall find it convenient throughout to refer to income distribution. We first consider whether inequality measures satisfying the principle of transfers will automatically yield unbounded  $IF$ s under very general conditions.

We define a general class of income inequality measures  $I(F)$

by the set of all  $I(F)$  which satisfies the principle of transfers (Dalton 1920). In other words if, say,  $I^*(F)$  belongs to this class, then the transfer of an arbitrary positive amount of income from a poorer income receiver to a richer income receiver (such that the mean of the distribution is preserved), increases the value of  $I^*(F)$ .

In order to study the effect on the estimator of an infinitesimal amount of contamination, we suppose that the underlying distribution lies in a neighbourhood of the model as defined in (2).  $H$  is in principle any perturbation distribution. In our case, an appropriate perturbation distribution is given by  $H^{(z)}$  with corresponding probability distribution

$$dH^{(z)}(x) = \begin{cases} 0.5 & \text{if } x = \mu - z = x_1(z) \\ 0.5 & \text{if } x = \mu + z = x_2(z) \end{cases} \quad (6)$$

where  $0 \leq z < \mu$ . Let us denote by  $G_\varepsilon^{(z)}$  the mixture distribution

$$G_\varepsilon^{(z)} := (1 - \varepsilon)F + \varepsilon H^{(z)}$$

then

$$\lim_{\varepsilon \rightarrow 0} \left[ \frac{I(G_\varepsilon^{(z)}) - I(F)}{\varepsilon} \right] = IF(z; I, F)$$

This kind of perturbation distribution is convenient because it

preserves the mean of the distribution<sup>4</sup>:

$$\begin{aligned}
 \mu(G_\varepsilon^{(z)}) &= \int x dG_\varepsilon^{(z)}(x) \\
 &= (1 - \varepsilon)\mu(F) + \frac{1}{2}\varepsilon(\mu(F) - z) + \frac{1}{2}\varepsilon(\mu(F) + z) \\
 &= \mu(F)
 \end{aligned} \tag{7}$$

It might be thought that the principle of transfers alone would give rise to an unboundedness problem: imagine a data contamination simultaneously at  $\mu - z$  and  $\mu + z$  and then consider increasing  $z$ . The principle of transfers implies that the empirical value of the inequality measure must increase and the question arises whether, for an indefinitely large value of  $z$ , the contamination will totally dominate the inequality measure. In other words, is  $\sup_z |IF(z; I, F)| = \lim_{z \rightarrow \infty} IF(z; I, F) = \infty$ , only by virtue of the principle of transfers?

Suppose that  $z^* > z$ , then by the principle of transfers we have

$$I(G_\varepsilon^{(z^*)}) > I(G_\varepsilon^{(z)}) \tag{8}$$

The  $IF$  of these inequality measures are then obtained by subtracting  $I(F)$ , dividing by  $\varepsilon$  and taking the limit when  $\varepsilon \rightarrow$

---

<sup>4</sup>We drop this restriction on the mean-preserving nature of the contamination in section 5 below.

0. We have

$$\lim_{\varepsilon \rightarrow 0} \frac{I(G_\varepsilon^{(z^*)}) - I(F)}{\varepsilon} \geq \lim_{\varepsilon \rightarrow 0} \frac{I(G_\varepsilon^{(z)}) - I(F)}{\varepsilon} \quad (9)$$

i.e.

$$IF(z^*; I, F) \geq IF(z; I, F) \quad (10)$$

However, we cannot conclude from (10) that the  $IF$  of an inequality measure which satisfies the principle of transfers necessarily is unbounded. Even if  $z$  may take any real value, the  $IF$  may actually be everywhere increasing and bounded; furthermore the admissible values of  $z$  may be restricted because the relevant support of  $F$  is - for some inequality measures - not the entire real line. In order to obtain sharper results, it is useful to consider a restricted class of measures. Under this restriction we will find that the principle of transfers can indeed give rise to problems of unboundedness for certain well-defined types of contamination.

#### 4 Decomposable inequality measures under mean-preserving contaminations

In this section we restrict our considerations to mean-preserving contaminations. We show that any inequality measure belonging to the class of decomposable inequality measures satisfy-



ing the principle of transfers and with corresponding income variable defined under an unrestricted domain must have an unbounded  $IF$ . Although we have restricted the class of inequality measures, the results found here are applicable to a great number of inequality measures, such as the coefficient of variation, the relative mean and median deviation, the variance of the logarithm of income (Gibrat 1931), members of the generalized entropy family, Hirschman's index (Hirschman 1945), Atkinson's index (Atkinson 1970), Kolm's index (Kolm 1976a, Kolm 1976b) etc. A notable exception to this class is the Gini coefficient. However, we will see below that the same conclusions can be drawn for this measure too.

#### 4.1 General properties

An inequality measure fulfilling the property of decomposability can be written in the following form

$$I(F) = \psi[J(F, \mu(F)) , \mu(F)] \tag{11}$$

where

$$J(F, \mu) = \int \phi(x; \mu) dF(x) , \tag{12}$$

$\phi$  and  $\psi$  are functions  $\mathcal{R}^2 \rightarrow \mathcal{R}$  and  $\psi$  is monotonic increasing in its first argument.

In the case of decomposable measures given by (11) and (12) where  $\phi$  is a differentiable function, the principle of transfers implies

$$\phi_1(x_2; \mu) - \phi_1(x_1; \mu) > 0 \text{ for any } x_1 < x_2 \quad (13)$$

where  $\phi_1$  is the derivative of  $\phi$  with respect to its first argument.

This result can be shown by considering an infinitesimal transfers of income  $dx > 0$  from a poorer income receiver to a richer income receiver. If the income distribution before the transfer is represented by  $F$ , then after the transfer, the new distribution  $F^*$  is the same as  $F$  except that

$$\exists x_1^*, x_2^* \text{ such that } x_1^* = x_1 - dx, x_2^* = x_2 + dx \text{ and } x_1 < x_2$$

By construction,  $\mu(F^*) = \mu(F) = \mu$ . The effect of this transfer on the inequality measure is given by

$$dI = I(F^*) - I(F) = \psi \left[ \int \phi(x; \mu) dF^*(x), \mu \right] - \psi \left[ \int \phi(x; \mu) dF(x), \mu \right] \quad (14)$$

The limiting case gives

$$\lim_{dx \rightarrow 0} \frac{1}{dx} \left\{ \psi \left[ \int \phi(x; \mu) dF^*(x), \mu \right] - \psi \left[ \int \phi(x; \mu) dF(x), \mu \right] \right\}$$

$$= \psi_1 [J(F, \mu), \mu] \left\{ -\phi'(x_1; \mu) + \phi'(x_2; \mu) \right\} \quad (15)$$

where  $\psi_1 > 0$  is the derivative of  $\psi$  with respect to its first argument. If  $I$  satisfies the principle of transfers, then we have that

$$\phi_1(x_2; \mu) - \phi_1(x_1; \mu) > 0 \quad \text{for any } x_1 < x_2 \quad (16)$$

that is,  $\phi$  must be strictly convex in its first argument.

## 4.2 Robustness properties

Let us now derive the  $IF$  of any  $I(F)$  satisfying the principle of transfers for the perturbation distribution defined in (6).

Given that

$$\begin{aligned} I(G_\varepsilon^{(z)}) &= \psi[J(G_\varepsilon^{(z)}, \mu(G_\varepsilon^{(z)})) , \mu(G_\varepsilon^{(z)})] \\ &= \psi[J(G_\varepsilon^{(z)}, \mu(F)) , \mu(F)] \end{aligned} \quad (17)$$

where

$$\begin{aligned} J(G_\varepsilon^{(z)}, \mu) &= (1 - \varepsilon) \int \phi(x; \mu) dF(x) + \frac{1}{2} \varepsilon \phi(x_1(z); \mu) + \\ &\quad \frac{1}{2} \varepsilon \phi(x_2(z); \mu) \end{aligned} \quad (18)$$

The  $IF$  of the inequality measure is given by

$$\begin{aligned}
 IF(z; I, F) &= \lim_{\varepsilon \rightarrow 0} \left[ \frac{I(G_\varepsilon^{(z)}) - I(F)}{\varepsilon} \right] \\
 &= \left. \frac{\partial}{\partial J} \psi[J, \mu] \cdot \frac{\partial}{\partial \varepsilon} J(G_\varepsilon^{(z)}, \mu) \right|_{\varepsilon=0} \\
 &= \psi_1[J, \mu] \cdot \left[ -J(F, \mu) + \frac{1}{2} \phi(x_1(z); \mu) + \right. \\
 &\quad \left. \frac{1}{2} \phi(x_2(z); \mu) \right] \tag{19}
 \end{aligned}$$

The behaviour of the  $IF$  is directly related to the properties of the function  $\phi$ . Since  $\phi$  is strictly convex in its first argument,  $\phi(x; \cdot)$  is unbounded either when  $x \rightarrow -\infty$  or when  $x \rightarrow \infty$ . Therefore, as  $z \rightarrow \infty$ , unless  $\phi$  is symmetric in  $\mu$ , the  $IF$  of  $I$  is unbounded.

However, it should be stressed that the above statement is valid only if we consider a random variable which takes its values on the whole set of the real numbers. It could be argued that in the case of *income* distributions an additional a priori restriction on the values of the random variable is appropriate, namely  $x \geq 0$ . If so, then the above analysis needs to be revised. Consider instead of (6) the following definition

$$dH^{(z)}(x) = \begin{cases} \frac{1}{1+z} & \text{if } x = z \cdot \mu \\ \frac{z}{1+z} & \text{if } x = \frac{1}{z} \cdot \mu \end{cases} \tag{20}$$

then the  $IF$  of  $I$  is given by

$$IF(z; I, F) = \psi_1[J, \mu] \left\{ -J(F, \mu) + \frac{1}{1+z} \phi(\mu \cdot z; \mu) + \frac{z}{1+z} \phi\left(\mu \frac{1}{z}; \mu\right) \right\} \quad (21)$$

The behaviour of the  $IF$  as  $z \rightarrow \infty$  depends on the behaviour of  $\phi(x; \mu)$  as  $x \rightarrow 0$  or as  $x \rightarrow \infty$ . We have the following situations. If

$$\lim_{x \rightarrow \infty} \phi(x; \mu) < \infty \text{ and } \phi(0; \mu) < \infty$$

then the  $IF$  is bounded. But if

$$\lim_{x \rightarrow \infty} \phi(x; \mu) = \infty \text{ or } \lim_{x \rightarrow 0} \phi(x; \mu) = \infty$$

then the  $IF$  is unbounded.

In the following three subsections, we will analyse some special cases of inequality measures more closely: the first two of these form an important subclass of the class of fully decomposable measures; the third special case is the Gini coefficient.

### 4.3 Kolm's index

Kolm's (Kolm 1976a, Kolm 1976b) index belongs to the class of decomposable inequality measures and is given by

$$I_K^\kappa(F) := \frac{1}{\kappa} \log \left[ \int e^{\kappa(\mu(F)-x)} dF(x) \right] \quad (22)$$

where  $\kappa$  is a parameter. Its function  $\phi$  is given by

$$\phi(x; \mu) = e^{\kappa(\mu-x)} \quad (23)$$

Since  $\lim_{x \rightarrow \infty} \phi(x; \mu) = 0$  and  $\phi(0; \mu) < \infty$ , the  $IF$  of Kolm's index with mean-preserving contaminations is bounded. This is also easily shown analytically.

Consider again the mixture model  $G_\varepsilon^{(z)}$  with the perturbation function given in (20). We have

$$I_K^\kappa(G_\varepsilon) = \frac{1}{\kappa} \log \left[ (1 - \varepsilon) \int e^{\kappa(\mu-x)} dF(x) + \frac{1}{1+z} \varepsilon e^{\kappa\mu(1-z)} + \frac{z}{1+z} \varepsilon e^{\kappa\mu(1-\frac{1}{z})} \right] \quad (24)$$

The  $IF$  of the Kolm index is given by

$$IF(z; I_K^\kappa, F) = \frac{1}{\kappa \int e^{\kappa(\mu-x)} dF(x)} \cdot \left\{ - \int e^{\kappa(\mu-x)} dF(x) + \frac{1}{1+z} e^{\kappa\mu(1+z)} + \frac{z}{1+z} e^{\kappa\mu(1+\frac{1}{z})} \right\} \quad (25)$$

Therefore, as  $\lim_{z \rightarrow \infty} IF(z; I_K^z, F) < \infty$ .

#### 4.4 Inequality measures from the generalized entropy family

Members of the ‘generalized entropy’ family are defined by

$$I_E^\beta := \frac{1}{\beta(\beta + 1)} \int \left[ \left( \frac{x}{\mu} \right)^{\beta+1} - 1 \right] dF(x). \quad (26)$$

where  $\beta \in (-\infty; +\infty)$ . This family has very convenient properties for the study of inequality: each inequality measure derived from it can be interpreted as a measure of the distance between the distribution of the income and the distribution in which every economic unit receives the mean income  $\mu$  (see Cowell 1980).

Theil’s (Theil 1967) inequality measures correspond to the limiting cases when  $\beta \rightarrow -1$  and  $\beta \rightarrow 0$ . They are given by

$$I_E^{-1} = - \int \log \left( \frac{x}{\mu} \right) dF(x) \quad (27)$$

$$I_E^0 = \int \left( \frac{x}{\mu} \right) \log \left( \frac{x}{\mu} \right) dF(x) \quad (28)$$

The corresponding  $\phi$  functions are given by

$$\phi(x; \mu) = \begin{cases} \frac{1}{\beta(\beta+1)} \left[ \left(\frac{x}{\mu}\right)^{\beta+1} - 1 \right] & \text{for } \beta \neq -1, 0 \\ -\log\left(\frac{x}{\mu}\right) & \text{for } \beta = -1 \\ \left(\frac{x}{\mu}\right) \log\left(\frac{x}{\mu}\right) & \text{for } \beta = 0 \end{cases} \quad (29)$$

For the  $IF$ , we have the following situation:

1.  $\beta \geq 0$ :  $\lim_{x \rightarrow \infty} \phi(x; \mu) = \infty$ , then the  $IF$  is unbounded.
2.  $-1 < \beta < 0$ :  $\lim_{x \rightarrow \infty} \phi(x; \mu) = -\infty$  and  $\phi(0; \mu) < \infty$ , then the  $IF$  is bounded.
3.  $\beta \leq -1$ :  $\lim_{x \rightarrow 0} \phi(x; \mu) = \infty$ , then the  $IF$  is unbounded.

These results can also be found by analysing the  $IF$  with mean-preserving contaminations. If we consider the perturbation function given by (20), the  $IF$  becomes

$$IF(z; I_E^{\beta \neq -1, 0}, F) = -I_E^\beta(F) + \frac{1}{\beta(\beta+1)} \left\{ \frac{z^{\beta+1}}{1+z} + \frac{z}{(1+z)z^{\beta+1}} - \right\} \quad (30)$$

$$IF(z; I_E^{-1}, F) = -I_E^{-1}(F) - \frac{1-z}{1+z} \log(z) \quad (31)$$

$$IF(z; I_E^0, F) = -I_E^0(F) - \frac{1-z}{1+z} \log(z) \quad (32)$$

We can see then that the  $IF$  with mean-preserving contaminations is bounded for the members of the generalized entropy family having the parameter  $-1 < \beta < 0$ .

However, we should emphasize that this result is only valid



for contaminations that leave the mean of the distribution unaffected. As we shall see in section 5, the possibility that the contamination affects the mean can have a dramatic impact on the  $IF$ .

#### 4.5 The Gini index

The Gini index does not belong to the class of decomposable inequality measures, although it is "non-overlapping" decomposable, see Cowell (1988) and Ebert (1988). However, we will show that with mean-preserving contaminations, the  $IF$  of this statistic is unbounded.

The Gini concentration ratio is given by

$$I_G(F) := 1 - 2 \int_0^1 q_F(\alpha) d\alpha \quad (33)$$

where  $q_F(\alpha)$  is the Lorenz ordinate of  $F$  given by

$$q_F(\alpha) = \mu^{-1}(F) \int_{-\infty}^{F^{-1}(\alpha)} u dF(u) \quad (34)$$

and can be interpreted as the proportion of income belonged by the proportion  $\alpha$  of the poorest income receivers.

We consider here the mixture distribution (2) which ensures that  $\mu(G_\varepsilon) = \mu(F) := \mu$ . The influence on  $I_G$  of this misspeci-

fication is then given by

$$\lim_{\varepsilon \rightarrow 0} \left[ \frac{I_G(G_\varepsilon) - I_G(F)}{\varepsilon} \right] = 2 \cdot \lim_{\varepsilon \rightarrow 0} \left[ \frac{R(F) - R(G_\varepsilon)}{\varepsilon} \right] \quad (35)$$

where

$$R(F) := \int_0^1 q_F(\alpha) d\alpha \quad (36)$$

We have

$$\begin{aligned} R(G_\varepsilon) &= \int_0^1 q_{G_\varepsilon}(\alpha) d\alpha \\ &= \int_0^1 \frac{1}{\mu} \int_0^{G_\varepsilon^{-1}(\alpha)} u dG_\varepsilon(u) d\alpha \\ &= \int_{-\infty}^{\infty} \frac{1}{\mu} \int_{-\infty}^x u dG_\varepsilon(u) dG_\varepsilon(x) \\ &= \int_{-\infty}^{\infty} \frac{1}{\mu} \left\{ (1 - \varepsilon) \int_{-\infty}^x u dF(u) + \varepsilon \int_{-\infty}^x u dH(u) \right\} dG_\varepsilon(x) \\ &= \frac{(1 - \varepsilon)^2}{\mu} \int_{-\infty}^{\infty} \int_{-\infty}^x u dF(u) dF(x) + \\ &\quad \frac{\varepsilon(1 - \varepsilon)}{\mu} \int_{-\infty}^{\infty} \int_{-\infty}^x u dF(u) dH(x) + \\ &\quad \frac{\varepsilon(1 - \varepsilon)}{\mu} \int_{-\infty}^{\infty} \int_{-\infty}^x u dH(u) dF(x) + \\ &\quad \frac{\varepsilon^2}{\mu} \int_{-\infty}^{\infty} \int_{-\infty}^x u dH(u) dH(x) \end{aligned} \quad (37)$$

Let us now consider the perturbation distribution  $H^{(z)}$  given

by (6)<sup>5</sup>. We then have

$$\begin{aligned}
 R(G_\varepsilon^{(z)}) &= (1 - \varepsilon)^2 R(F) + \frac{\varepsilon(1 - \varepsilon)}{\mu} \left\{ \frac{1}{2} \int_{-\infty}^{\mu-z} u dF(u) + \right. \\
 &\quad \left. \frac{1}{2} \int_{-\infty}^{\mu+z} u dF(u) \right\} + \\
 &\quad \frac{\varepsilon(1 - \varepsilon)}{\mu} \left\{ \frac{\mu - z}{2} \int_{\mu-z}^{\mu+z} dF(x) + \frac{\mu + z}{2} \int_{\mu+z}^{\infty} dF(x) \right\} + \\
 &\quad \frac{\varepsilon^2}{\mu} \left\{ \frac{1}{2} \int_{-\infty}^{\mu-z} u dH^{(z)}(u) + \frac{1}{2} \int_{-\infty}^{\mu+z} u dH^{(z)}(u) \right\} \quad (38)
 \end{aligned}$$

The  $IF$  is finally given by

$$\begin{aligned}
 IF(z; I_G, F) &= 4R(F) - \frac{1}{\mu} \left[ \int_{-\infty}^{\mu-z} u dF(u) + \int_{-\infty}^{\mu+z} u dF(u) \right] - \\
 &\quad \left[ \frac{\mu - z}{\mu} (F(\mu + z) - F(\mu - z)) + \right. \\
 &\quad \left. \frac{\mu + z}{\mu} (1 - F(\mu + z)) \right] \quad (39)
 \end{aligned}$$

Now consider the behaviour of the  $IF$  as  $z$  varies.

$$\frac{\partial}{\partial z} IF(z; I_G, F) = \frac{1}{\mu} [2F(\mu + z) - F(\mu - z) - 1] - \frac{\mu - z}{\mu} f(\mu + z) \quad (40)$$

As  $z \rightarrow \infty$ ,  $\frac{\partial}{\partial z} IF(z; I_G, F) \rightarrow \frac{1}{\mu}$  and therefore  $IF(z; I_G, F) \rightarrow \infty$ .

So even, if we restrict attention to the special case where

---

<sup>5</sup>Since for the Gini index the underlying income variable is defined for the whole set of real numbers, we choose this type of mean-preserving contamination. However, choosing for  $H^{(z)}$  the perturbation distribution given by (20) would lead to the same conclusion.

contamination does not affect the mean of the distribution we find that the  $IF$  of the Gini coefficient is unbounded<sup>6</sup>.

## 5 Robustness properties with arbitrary contaminations

So far we have considered the issue of robustness of inequality measures given a very special kind of contamination: one which preserves the mean of the distribution. This approach would be relevant to cases in which inequality was defined in terms of income shares rather than incomes, and observations were available on shares. Under such circumstances the population mean would be known by definition. But this sort of situation is exceptional. In other cases we have to assume either that the impact of the contamination on the mean is negligible, or that the impact on inequality of variability in the mean is negligible, neither of which is satisfactory.

However when we allow for arbitrary perturbations to a distribution the resulting impact on inequality measures is going to yield an expression involving both transfer effects and a change in the mean, which are difficult to interpret analytically without more restriction on the class of inequality measures.

---

<sup>6</sup>Unboundedness of the  $IF$  with arbitrary contaminations has been shown by Monti (1992).

In this section, we derive the  $IF$  of an inequality measure belonging to the class of decomposable measures. We consider here the mixture model  $G_\varepsilon$  where the perturbation distribution is the point mass 1 at an arbitrary income level  $z$ . The  $IF$  is given by

$$IF(z; I, F) = \psi_1[J(F, \mu), \mu] \cdot \left. \frac{\partial J(G_\varepsilon, \mu(G_\varepsilon))}{\partial \varepsilon} \right|_{\varepsilon=0} + \psi_2[J(F, \mu), \mu] \cdot (z - \mu) \quad (41)$$

where  $\psi_j$ ,  $j = 1, 2$ , is the derivative of  $\psi$  with respect to its  $j^{\text{th}}$  argument. For most inequality measures, the relevant part is given by the first term in (41). We have

$$\begin{aligned} \left. \frac{\partial J(G_\varepsilon, \mu(G_\varepsilon))}{\partial \varepsilon} \right|_{\varepsilon=0} &= \left. \frac{\partial}{\partial \varepsilon} \left\{ (1 - \varepsilon) \int \phi(x; \mu(G_\varepsilon)) dF(x) + \varepsilon \phi(z; \mu(G_\varepsilon)) \right\} \right|_{\varepsilon=0} \\ &= - \int \phi(x; \mu) dF(x) + \int \frac{\partial}{\partial \mu} \phi(x; \mu) dF(x) \cdot [z - \mu] + \phi(z; \mu) \end{aligned} \quad (42)$$

We may make the following remarks: (a) The first term is independent of  $z$ , (b) the second term is due to the effect of the contamination on the mean and (c) the third term would also be present if were to suppose that the mean be given.

If a decomposable measure has an unbounded  $IF$  for the mean-preserving contamination case, then it must also be unbounded in the arbitrary contamination case by virtue of (c) above.

However, some inequality measures which have bounded  $IF$  in the mean-preserving contamination case will be unbounded in the arbitrary contamination case because of (b). This can easily be shown in the case of the class of decomposable measures. Indeed, consider the same perturbation function as above and suppose that the mean  $\mu$  is given. The  $IF$  of the members of the generalized entropy family with parameter  $-1 < \beta < 0$  (for which the  $IF$  with mean-preserving contamination is bounded), is given by

$$IF(z; I_\beta, F) = -I_\beta(F) + \frac{1}{\beta(\beta + 1)} \left[ \left( \frac{z}{\mu} \right)^{\beta+1} - 1 \right] \quad (43)$$

We can see that the  $IF$  is unbounded.

Drawing together our results here and in section 4 we may state:

**Proposition 1:**

*If the mean has to be estimated from the sample then all scale independent, translation independent, and decomposable inequality measures have an unbounded  $IF$ .*

The proposition follows from our results on the Kolm indices which form the entire class of decomposable measures that are translation independent and on the Generalized Entropy indices which form the entire class of decomposable measures that are scale independent. It also includes all “intermediate” cases as defined by Bossert and Pfingsten (1990).

If we try to take account of the case where contamination affects the mean (which is probably more realistic in most applications) the problem of an unbounded  $IF$  arises for a large class of commonly used inequality measures. This presents a serious problem for the empirical analysis of income inequality from micro data, since it means that a few “rogue” low observations (in the case of bottom-sensitive inequality measures) or high observations (for the other types of inequality measure) would drive the estimated value of the inequality index on their own. The situation may not be so bad for the rare cases where the mean does not itself have to be estimated, but this is small comfort.

What can be done about this situation? One approach is to screen pragmatically by eye. This is frequently done and usually for good reasons. However, although the researcher’s judgment may be very good in a particular instance, the procedure is very arbitrary. Moreover, since the only option is

whether or not to drop the questionable observation from the sample, there can be situations in which this procedure is regarded as too drastic. We will suggest an alternative procedure in section 6 and 7 below. First let us take a look at the importance of contamination on empirical estimates of inequality.

## 5.1 Simulation study

For this first simulation study, we computed the Theil index (see equation (28)) which as many other inequality measures, may be determined empirically by only a few observations. This is shown in table 1. We computed 100 samples of 200 observations generated by a Lognormal distribution given by

$$F_{\mu,\sigma}(t) = \int_0^t \frac{1}{x\sqrt{\sigma^2 2\pi}} e^{-\frac{1}{2\sigma^2}(\log(x)-\mu)^2} dx \quad (44)$$

with parameters  $\mu = 1.0$  and  $\sigma = 0.8$ . We contaminated a percentage of those observations by simply multiplying them by 10. The differences between the Theil indexes are then due to those small proportions of observations.

We can see that with only 5% of “rogue” data, Theil’s index becomes twice its initial value! However, one could argue that this type of contamination, i.e. a random proportion of contamination multiplied by 10, is too extreme. But as said before, this is the type of error that can occur during the recording



Degree of contamination	Theil Index	(SD)
0%	0.312	(0.002)
1%	0.430	(0.024)
2%	0.515	(0.030)
3%	0.586	(0.027)
4%	0.649	(0.029)
5%	0.700	(0.031)

Table 1: Empirical Theil index when a random proportion of data are multiplied by 10

process, that is a “decimal point error”. The other type of error we consider here is the week-month confusion defined in the introduction. Although this type of error is less large, we can see by means of a simulation study that its effect on inequality measures is quite important. In table 2 we show the computed Theil’s index when a random proportion of data are multiplied by 4.

Degree of contamination	Theil Index	(SD)
0%	0.312	(0.002)
1%	0.334	(0.003)
2%	0.357	(0.004)
3%	0.379	(0.010)
4%	0.400	(0.010)
5%	0.429	(0.012)

Table 2: Empirical Theil index when a random proportion of data are multiplied by 4

## 6 The estimation of inequality through parametric distributions

As an alternative to direct estimation of inequality, we will consider here the “parametric approach”, i.e. the analysis of inequality measures estimated through a parametric model  $F = F_\theta$ , where  $\theta$  is a vector of parameters. Although the following results apply to any type of inequality measure, we will examine the particular case of the generalized entropy family.

The parametric approach has several advantages. First it allows one to build robust estimators which have some optimality properties (for example minimal variance among estimators belonging to the same class). Second, it allows one to distinguish systematically between extreme values that “belong” to the estimate of inequality and those that do not. In effect it permits one to construct an automatic procedure for detecting observations that are very “far” (in an appropriate sense) from the bulk of the data and that should be treated separately. This is particularly important in parts of the sample where the observations are sparse. A good example of the use of the parametric approach arises when estimating the inequality of wealth distributions: in such cases the common practice of fitting a Pareto distribution to the upper tail permits the re-

searcher to build in information about the structure of wealth distributions in general in order to make the best use of one particular sample.

It may be argued that the use of a parametric approach has the inbuilt problem that inappropriate model choice can of itself bias the estimators, and that non-parametric methods would avoid this type problem. However it is possible to mitigate model-selection bias by prior use of a selection procedure, and robust methods for such a procedure are available - see Victoria-Feser (1993).

An estimator of one of the members of the generalized entropy family (henceforth called GEPFE) is obtained by replacing  $F$  in (26) by  $F_\theta$  where  $\theta$  is estimated from the sample.

Let  $T$  be an estimator of  $\theta$  for the parametric model  $F_\theta$ .  $T$  can be written as a functional  $T(F_\theta)$  of the distribution. Hence, for the contaminated model  $G_\epsilon$  given by

$$G_\epsilon = (1 - \epsilon)F_\theta + \epsilon\Delta_z \quad (45)$$

$T(G_\epsilon)$  can also be influenced by model deviations. The  $IF$  of a GEPFE is given by

$$IF(z; I_\beta, F_\theta) = A(\beta, F_\theta) \cdot IF(z; T, F_\theta) \quad (46)$$

where

$$A(\beta, F_\theta) = \frac{1}{\beta(\beta+1)} \int \left( \frac{x}{\mu(F_\theta)} \right)^{\beta+1} s(x, \theta)^T dF_\theta(x) - \frac{1}{\mu(F_\theta)} \frac{1}{\beta} \int \left( \frac{x}{\mu(F_\theta)} \right)^{\beta+1} dF_\theta(x) \int x s(x, \theta)^T dF_\theta(x)$$

and  $s$  is the scores function given by

$$s(x; \theta) = \frac{\partial}{\partial \theta} \log f(x; \theta) \quad (48)$$

$f(x; \theta)$  being the density function corresponding to  $F_\theta$ . For the particular cases when  $\beta = 0$  and  $\beta = -1$ , we have

$$\begin{aligned} A(0, F_\theta) &= \int \frac{x}{\mu(F_\theta)} \log \left( \frac{x}{\mu(F_\theta)} \right) s(x; \theta)^T dF_\theta(x) - \\ &\int \frac{x}{\mu(F_\theta)} s(x; \theta)^T dF_\theta(x) \cdot \\ &\int \frac{x}{\mu(F_\theta)} \log \left( \frac{x}{\mu(F_\theta)} \right) dF_\theta(x) + \\ &\int \frac{x}{\mu(F_\theta)} s(x; \theta)^T dF_\theta(x) \end{aligned} \quad (49)$$

and

$$\begin{aligned} A(-1, F_\theta) &= \int \log \left( \frac{x}{\mu(F_\theta)} \right) s(x; \theta)^T dF_\theta(x) + \\ &\int \frac{x}{\mu(F_\theta)} s(x; \theta)^T dF_\theta(x) \end{aligned} \quad (50)$$

Since  $A(\beta, F_\theta)$  does not depend on  $z$ , the  $IF$  of the GEPFE is

proportional to the  $IF$  of the estimator of  $\theta$  given by

$$IF(z; T, F_\theta) = \left. \frac{\partial}{\partial \varepsilon} T(G_\varepsilon) \right|_{\varepsilon=0} \quad (51)$$

Therefore, if  $IF(z; T, F_\theta)$  is unbounded then so is  $IF(z; I_\beta, F_\theta)$ .

Typically,  $\theta$  is estimated by the maximum likelihood estimator (MLE). The  $IF$  of this estimator is proportional to the scores function (see Hampel, Ronchetti, Rousseeuw, and Stahel 1986. For almost all the models used in income distributions, the scores function is unbounded (see Victoria-Feser 1993). Hence, GEPFE computed through MLE have an unbounded  $IF$ .

For example, consider the Pareto distribution with density function

$$f(x; \theta) = \theta x^{-(\theta+1)} x_0^\theta \quad (52)$$

where  $0 < x_0 \leq x < \infty$ . The corresponding scores function is given by

$$s(x; \theta) = \frac{1}{\theta} - \log(x) + \log(x_0) \quad (53)$$

Then  $IF(z; I_\beta, F_\theta)$  is proportional to  $-\log(z)$  which is clearly unbounded.

## 7 Robust income inequality measures

### 7.1 Optimal B-robust estimators

In this section, we show that a reasonable way of obtaining robust estimators of income inequality measure is through the specification of a parametric model. Parametric models allow one to compute robust estimators in an optimal way, which yield robust estimated income inequality measures. In other words, given an appropriate specification of a parametric model for the income distribution, we are able to compute robust inequality measures through the robust estimators for the parameters of the model. In this section, we first present robust estimators for the parameters of income distribution models which are optimal in a sense defined below and present some simulation results.

For the estimator of  $\theta$ , we propose taking the optimal B-robust estimators (OBRE) defined in Hampel, Ronchetti, Rousseeuw, and Stahel (1986). These estimators are M-estimators, i.e. they belong to the class of estimators defined by

$$\sum_{i=1}^n \psi(x_i; \theta) = 0 \quad (54)$$

where  $\psi$  is any function  $\mathcal{R} \rightarrow \mathcal{R}^p$ ,  $p = \dim(\theta)$ . For example, the MLE belong to this class, and in this case  $\psi$  is equal to

the scores function. However, although MLE are the most efficient estimators, they are in general non-robust. By defining the wider class of M-estimators, one can obtain robust estimators optimally. The trade-off is between efficiency and robustness (bounded  $IF$ ). OBRE are the optimal solution of this trade-off. They are called  $B$ -robust for *Bias*-robust because their bias measured by the  $IF$  is bounded. One should stress that there are several versions of OBRE that depend on the way one chooses to minimize the asymptotic covariance matrix. We propose here standardized OBRE (see Hampel, Ronchetti, Rousseeuw, and Stahel 1986).

Their expression is similar to that of MLE. They are defined implicitly by

$$\sum_{i=1}^n \psi(x_i; \theta) = \sum_{i=1}^n [s(x_i; \theta) - a(\theta)] \cdot W_c(x_i; A(\theta), a(\theta)) = 0 \quad (55)$$

In (55),  $W_c(x_i; \dots)$  are actually weights attributed to each observation according to its influence on the estimator. They depend on a constant  $c$  which can be seen as a regulator between efficiency and robustness: the lower  $c$  is the more robust is the estimator. At  $c = \infty$  we have the MLE.

The  $p \times p$  matrix  $A(\theta)$  and the  $p \times 1$  vector  $a(\theta)$  used in (55)

are defined implicitly by

$$E [\psi(x; \theta)\psi(x; \theta)^T] = [A(\theta)^T A(\theta)]^{-1} \quad (56)$$

$$E [\psi(x; \theta)] = 0 \quad (57)$$

so that the constraints of efficiency and consistency are satisfied. The solution of (55) is not straightforward and in general OBRE have to be computed iteratively. For a more precise definition of OBRE and a simple algorithm, see Victoria-Feser and Ronchetti (1993).

## 7.2 Simulation results

In order to show how large the bias on the income inequality measure can be when the data are contaminated, we computed the Theil index for different samples. The results come from 100 simulated samples of 200 observations from a Gamma distribution given by

$$F_{\alpha, \lambda}(t) = \frac{\lambda^\alpha}{\Gamma(\alpha)} \int_0^t x^{\alpha-1} e^{-\lambda x} dx \quad (58)$$

where  $\Gamma(\alpha) = \int_0^1 t^{\alpha-1} e^{-t} dt$ . The mixture distribution is

$$G_\varepsilon = (1 - \varepsilon)F_{\alpha, \lambda} + \varepsilon F_{\alpha, 0.1 \cdot \lambda} \quad (59)$$



where  $\alpha = 3.0$  and  $\lambda = 1.0$ . That is, the contaminations are observations from the same distribution but multiplied by ten. The results are presented in table 3.

We can see that only 3% of contaminations push the value of the Theil index from 0.155 to 0.27, and hence determine alone the value of the estimator. This is certainly not a desirable property. However, it is possible to avoid those situations with robustified inequality measures.

		Bias of the MLE	MSE of the MLE	Theil's Index
No contamination	$\alpha$	0.05	0.07	0.155
	$\lambda$	0.01	0.01	
3% contamination	$\alpha$	-1.33	1.89	0.270
	$\lambda$	-0.54	0.32	
5% contamination	$\alpha$	-1.72	3.09	0.342
	$\lambda$	-0.70	0.5	

Table 3: MLE and Theil index with and without data contamination

In table 4 we presents the same computations as in table 3 except that instead of using the MLE, we calculated the OBRE ( $c = 1.5$ ). More details on computational aspects are given in Victoria-Feser (1993).

		Bias of the OBRE	MSE	Theil's Index
No contamination	$\alpha$	0.06	0.07	0.155
	$\lambda$	0.02	0.01	
3% contamination	$\alpha$	0.15	0.12	0.165
	$\lambda$	0.07	0.02	
5% contamination	$\alpha$	-0.22	0.16	0.169
	$\lambda$	-0.11	0.03	

Table 4: OBRE and Theil index with and without data contamination

The results show that with OBRE, the value of the Theil index is almost unaffected by the contaminations, even with a contamination proportion as high as 5%. Therefore, there are potential gains by using OBRE instead of MLE when computing estimates of the underlying parametric model which then serves to compute robust estimates of inequality.

## 8 Conclusions

Inequality measures should convey practical information about income distributions. Whether or not they do this effectively depends of course upon the reliability of income distribution data; it also depends upon the method of estimation. It is possible that under certain estimation procedures the apparent picture of inequality is strongly influenced by data contamination.

The conventional properties of inequality measures also play a rôle in determining the extent to which the picture of income distribution is distorted by data contamination: the  $IF$  is a useful device for quantifying this effect. One might have supposed that the fundamental property of mainstream inequality analysis - the transfer principle - is a sufficient condition for unboundness of the  $IF$ . We have shown that this is not the case. However, this negative result is not actually very encour-

aging because we have shown that as soon as one introduces other standard features of inequality measures (scale independence, decomposability) or a more realistic specification of the estimation problem (where the mean itself has to be estimated rather than being specified *a priori*) the problem of unboundness of the *IF* re-emerges. As we have shown, the impact upon measured inequality of quite small amounts of contamination in the tails of the distribution can be disastrous. Controlling for this contamination in practice can be tricky, and *ad hoc* methods are likely to be unreliable.

One way of dealing with this problem is to adopt a parametric approach to income distribution analysis and inequality measurement, and to estimate inequality through robust estimates of the parameters of the income distribution model. Even if the robust estimates of inequality are not used on their own, they should provide a useful supplementary check against estimates of inequality computed by classical methods. Where discrepancies between the results for the two approaches emerge and are attributable to a small number of observations, this information should be taken into account in drawing conclusions about the “true” picture of the underlying income distribution.

## References

- Atkinson, A. B. (1970). On the measurement of inequality. *Journal of Economic Theory* 2, 244–263.
- Bossert, W. and A. Pfingsten (1990). Intermediate inequality: Concepts, indices and welfare implications. *Mathematical Social Sciences* 19, 117–134.
- Cowell, F. A. (1977). *Measuring Inequality* (First ed.). Oxford: Phillip Allan.
- Cowell, F. A. (1980). Generalized entropy and the measurement of distributional change. *European Economic Review* 13, 147–159.
- Cowell, F. A. (1988). Poverty measures, inequality and decomposability. In D. Bös, M. Rose, and C. Seidl (Eds.), *Welfare and Efficiency in Public Economics*, pp. 149–166. Springer-Verlag.
- Dalton, H. (1920). The measurement of the inequality of incomes. *Economic Journal* 30, 348–361.
- Ebert, U. (1988). Measurement of inequality: An attempt at unification and generalization. *Social Choice and Welfare* 5, 147–169.
- Gibrat, R. (1931). *Les Inégalités Economiques*. Paris: Sirey.

- Hampel, F. R. (1968). *Contribution to the Theory of Robust Estimation*. Ph. D. thesis, University of California, Berkeley.
- Hampel, F. R. (1974). The influence curve and its role in robust estimation. *Journal of the American Statistical Association* 69, 383–393.
- Hampel, F. R., E. Ronchetti, P. J. Rousseeuw, and W. A. Stahel (1986). *Robust Statistics: The Approach Based on Influence Functions*. New York: John Wiley.
- Hirschman, A. O. (1945). *National Power and the Structure of Foreign Trade*. Berkeley: University of California Press.
- Huber, P. J. (1964). Robust estimation of a location parameter. *Annals of Mathematical Statistics* 35, 73–101.
- Jenkins, S. P. (1992). Accounting for inequality trends: Decomposition analyses for the UK, 1971–1986. Mimeo, University College of Swansea, UK.
- Kolm, S. C. (1976a). Unequal inequalities. *Journal of Economic Theory* 12, 416–442. Part I.
- Kolm, S. C. (1976b). Unequal inequalities. *Journal of Economic Theory* 13, 82–111. Part II.
- Monti, A. C. (1992). The study of the Gini concentration ratio by means of the influence function. *Statistica* 51,

561-577.

Pudney, S. and H. Sutherland (1992). The statistical reliability of micro-simulation estimates: Results for a UK tax-benefit model. *Journal of Public Economics*. To appear.

Theil, H. (1967). *Economics and Information Theory*. Amsterdam: North- Holland.

Victoria-Feser, M.-P. (1993). *Robust Methods for Personal Income Distribution Models*. Ph. D. thesis, University of Geneva.

Victoria-Feser, M.-P. and E. Ronchetti (1993). Robust methods for personal income distribution models. *The Canadian Journal of Statistics*. To appear.

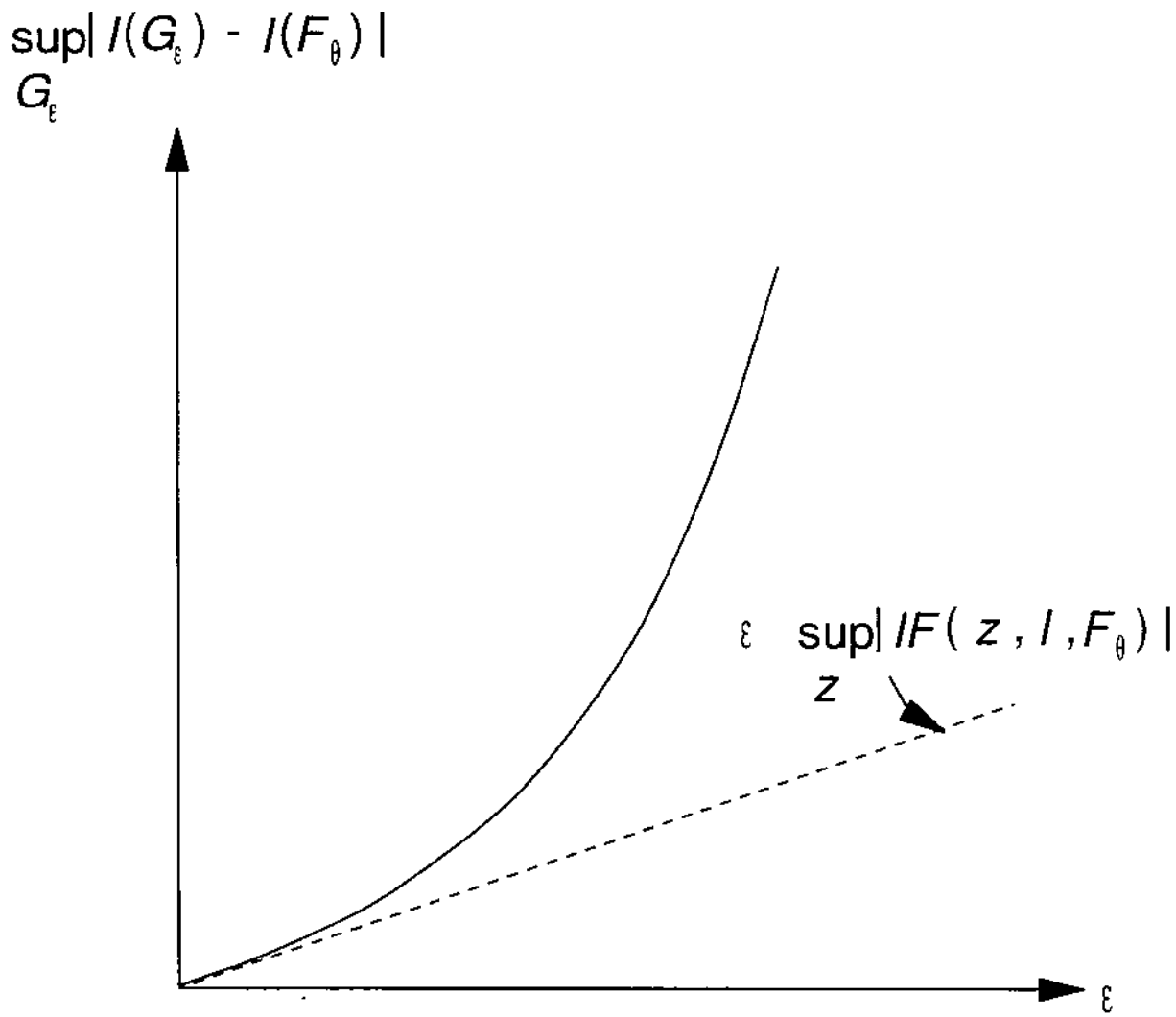


Figure 2. Bias on the estimator.

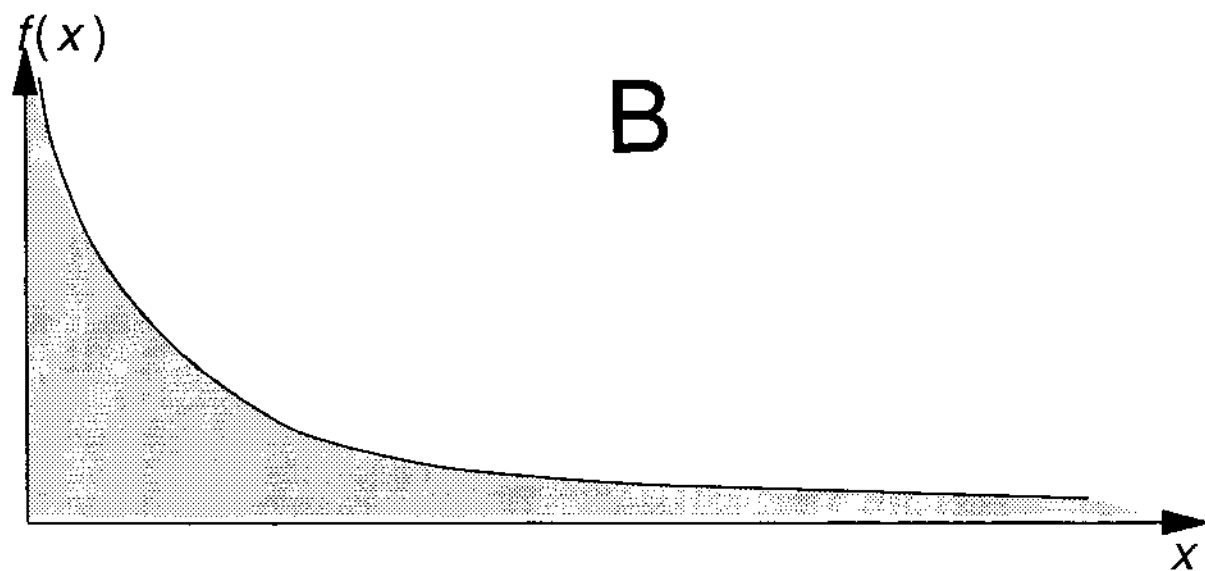
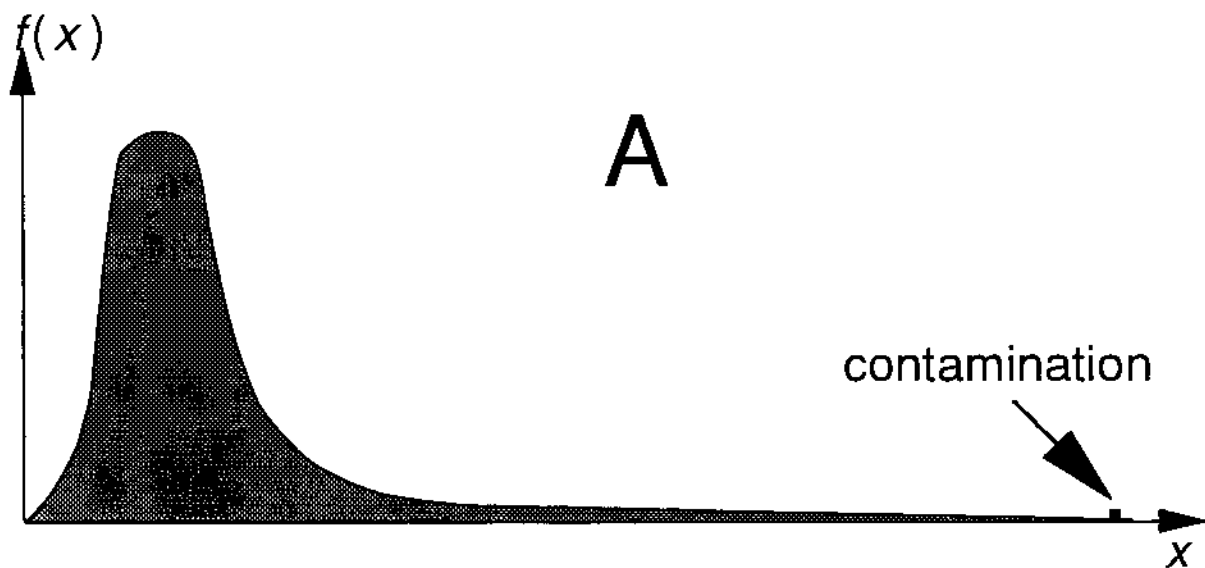


Figure 1. Confusion caused by contamination.