

# Permutation tests for comparing inequality measures\*

Jean-Marie Dufour<sup>†</sup>  
McGill University

Emmanuel Flachaire<sup>‡</sup>  
Aix-Marseille University

Lynda Khalaf<sup>§</sup>  
Carleton University

December 2015

---

\* This work was supported by the William Dow Chair in Political Economy (McGill University), the Bank of Canada (Research Fellowship), a Guggenheim Fellowship, a Konrad-Adenauer Fellowship (Alexander-von-Humboldt Foundation, Germany), the Institut de finance mathématique de Montréal (IFM2), the Social Sciences and Humanities Research Council of Canada, the Fonds de recherche sur la société et la culture (Government of Québec), the Institut Universitaire de France and the A\*MIDEX project (ANR-11-IDEX-0001-02) funded by the "Investissements d'Avenir" French Government program, managed by the French National Research Agency (ANR).

<sup>†</sup> William Dow Professor of Economics, McGill University, Centre interuniversitaire de recherche en analyse des organisations (CIRANO), and Centre interuniversitaire de recherche en économie quantitative (CIREQ). Mailing address: Department of Economics, McGill University, Leacock Building, Room 519, 855 Sherbrooke Street West, Montréal, Québec H3A 2T7, Canada. TEL: (1) 514 398 8879; FAX: (1) 514 398 4938; e-mail: jean-marie.dufour@mcgill.ca . Web page: <http://www.jeanmariedufour.com>

<sup>‡</sup> Aix-Marseille University (Aix-Marseille School of Economics), EHESS, CNRS and Institut Universitaire de France, 2 rue de la charité, 13002, Marseille, France. Email: [emmanuel.flachaire@univ-amu.fr](mailto:emmanuel.flachaire@univ-amu.fr).

<sup>§</sup> Economics Department, Carleton University, Centre interuniversitaire de recherche en économie quantitative (CIREQ), and Groupe de recherche en économie de l'énergie, de l'environnement et des ressources naturelles (GREEN), Université Laval. Mailing address: Economics Department, Carleton University, Loeb Building 1125 Colonel By Drive, Ottawa, Ontario, K1S 5B6 Canada. Tel (613) 520-2600-8697; FAX: (613)-520-3906. email: [Lynda.Khalaf@carleton.ca](mailto:Lynda.Khalaf@carleton.ca).

## ABSTRACT

Asymptotic and bootstrap tests based on inequality measures are known to perform poorly in finite samples when the underlying distribution is heavy-tailed. We propose Monte-Carlo permutation and bootstrap methods for the problem of testing the equality of inequality measures between two samples. Results cover the Generalized Entropy class which includes Theil's index, the Atkinson class of indices and the Gini index. We analyze finite-sample and asymptotic conditions for the validity of proposed methods, and introduce a convenient rescaling to improve finite-sample performance. Simulation results show that size correct inference can be obtained with our proposed methods despite heavy tails if the underlying distributions are sufficiently close in the upper tails. Substantial reduction in size distortion is achieved more generally. Studentized rescaled Monte Carlo permutation tests outperform the competing methods we consider in terms of power.

**Key words:** permutation test; bootstrap; inequality measures; income distribution

**Journal of Economic Literature classification:** C15, D31, D63.

# 1 Introduction

Income and wealth distributions are typically non-normal and can take various shapes. In view of this, distribution-free approaches are especially well suited to the task of comparing inequality measures. However, despite a sizable literature, non-parametric methods for inference on such measures perform poorly in finite samples. As the sample size grows, concern shifts from small sample distortions to asymptotic problems caused by the failure of the assumptions needed to ensure size control. These problems are often associated with heavy tails, a common situation in applied work. In economics for example, income inequalities are of primary interest and income distributions are characterized by a prominently heavy right tail. In addition, inequality measures can be equal even if the underlying distributions differs, which also confounds inference.

Consider two variables  $x$  and  $y$  assumed drawn from two distributions,  $F_x$  and  $F_y$ . We study distribution-free tests of

$$H_0 : \theta(F_x) = \theta(F_y) \tag{1}$$

where  $\theta(\cdot)$  is some functional on some subset  $\mathcal{F}$  of distributions (further structure is provided below). Inequality indices that are special cases of (1) provide the motivation for our work. Formally, we analyze centered and uncentered moments, the Generalized Entropy (**GE**) class of inequality measures which includes Theil's index, the Atkinson class of inequality indices (Atkinson, 1970) and the Gini index.

While bootstraps offer a natural alternative to standard asymptotics for this problem, Davidson and Flachaire (2007), Cowell and Flachaire (2007), Schluter and van Garderen (2009) or Davidson (2009, 2012) show that heavy tails also cause bootstrap failures. A few improvements have been proposed. Davidson and Flachaire (2007) consider a bootstrap data generating process (DGP) which combines a parametric estimate of the upper tail with a non-parametric estimate of the rest of the distribution. Schluter and van Garderen (2009) propose normalizing transformations of inequality measures using Edgeworth expansions, to adjust asymptotic Gaussian approximations. Such corrections which seem effective in

specific instances, for example when the null hypothesis takes the form  $H_1 : \theta(F_x) = \delta_0$  with  $\delta_0$  known, still fail with heavy-tailed distributions.

This paper analyzes permutation methods for inference on (1) and offers compelling arguments that permutation-based Monte Carlo (**MC**) test methods (Pitman, 1937; Dwass, 1957; Dufour, 2006) are one large step in the right direction. We ask what finite and large-sample assumptions are needed to support reliable permutations, focusing on the specificities of commonly used inequality measures. Our analysis applies and extends the theoretical setups of Romano (1990), Dufour (2006) and Chung and Romano (2013).

We first consider a baseline problem which restricts (1) to the case where available samples are drawn from the same distribution, that is when  $F_x = F_y$ . We show that MC permutation tests provide exact inference in finite samples, even if the common distribution is heavy-tailed. Our result allows for continuous and discrete distributions, and does not require any regularity condition on the form of the functional  $\theta(\cdot)$ . We also allow for exchangeable (as opposed to i.i.d.) observations, hence covering the case of random draws without replacement from a finite population. To the best of our knowledge, although restrictive, this special case provides the only available exact solution for the problem at hand.

The fact remains that (1) does not guarantee that  $F_x = F_y$ . A simple counter-example is provided in Figure 1 which depicts Singh-Maddala distributions (Burr XII) with density  $f(u) = aqu^{a-1}/(b^a[1 + (u/b)^a]^{1+q})$ , for two choices of  $a, b$  and  $q$ : 2.8, 0.1930698, 1.7 [depicted as  $\mathbf{F}_x$ ] and 4.8, 0.1930698, 0.6366578 [depicted as  $\mathbf{F}_y$ ]. The latter choice yields a distribution much more heavy-tailed than the former, yet both distributions share the same value of the Theil index which equals 0.1401. In this case, the use of permutation tests is not justified from an exact perspective. Romano (1990) shows that, when  $F_x \neq F_y$ , permutation tests of the hypothesis in (1) are asymptotically valid - in the sense that the probability of a Type I error tends asymptotically to the nominal level - for some very specific cases, but are not valid in general. For instance, permutations work using differences of sample means if the samples are of the same size, but are invalid with differences

of medians. We suggest a convenient rescaling that validates permutations for several inequality measures. A bootstrap method for this null hypothesis is also proposed.

More recently, Chung and Romano (2013) showed that permutation tests are asymptotically valid in a more general setting if the underlying statistic is studentized. The importance of studentization is well-known for bootstrapping to achieve asymptotic refinements (Hall, 1992). In contrast, with permutation tests, studentization may be required for validity. In particular, when comparing medians, studentized statistics will work while non-studentized counterparts are invalid. Although Chung and Romano (2013) do not analyze inequality measures, their general statistical setup validates comparing these measures using studentized criteria. The rescaling we introduce may not be necessary for size control with studentized criteria, at least asymptotically. Yet we show that it matters from the power perspective.

Simulation experiments are provided to study the finite-sample properties of proposed tests when the samples are drawn from similar and different distributions. We consider excessively heavy-tailed cases to analyze size within a worst-case scenario design. In terms of Type 1 error or size distortion, our results show that when the samples are drawn from two (strongly) heavy-tailed distributions which are not too different, permutation tests perform very well in finite samples. When distributions differ dramatically particularly in their tails, while size distortions are not completely eradicated, permutation tests outperform the standard asymptotic and bootstrap tests. In terms of power, our results show that permutation tests based on rescaled samples perform better in small samples than permutation tests based on original samples.

The paper is organized as follows. Section 2 sets a general framework and presents our proposed inference methods. In section 3, we show how exact simulation-based permutation tests for the hypothesis of equal distributions can be obtained using statistics comparing general functionals of empirical distribution functions. In section 4, we consider the problem of testing the equality of general functionals when the distributions of the two populations can differ. In section 5, we study specific cases based on moments and

commonly used inequality measures. Simulation experiments are reported in section 6. We conclude in section 7.

## 2 Framework

In this section, we set notation, define the test statistics, and present alternative permutational and simulation-based  $p$ -values for comparing a general functional  $\theta(\cdot)$  on two different populations. The specific treatment of inequality measures is deferred to section 5.

We consider two samples  $X = \{X_1, X_2, \dots, X_n\}$  and  $Y = \{Y_1, Y_2, \dots, Y_m\}$  each of which comprises independent and identically distributed observations with cumulative distribution functions  $F_x$  and  $F_y$  respectively. We wish to test general hypotheses of type  $H_0$  as in (1). A natural statistic for such a problem is given by

$$T = \theta(\hat{F}_x) - \theta(\hat{F}_y) \quad (2)$$

where  $\hat{F}_x$  and  $\hat{F}_y$  are the empirical distribution functions (**EDFs**) of the samples  $X$  and  $Y$ , and  $N$  is the total number of observations,  $N = n + m$ . On studentizing  $T$ , we get the *studentized* test statistic

$$S = \frac{\theta(\hat{F}_x) - \theta(\hat{F}_y)}{\sqrt{\hat{V}[\theta(\hat{F}_x)] + \hat{V}[\theta(\hat{F}_y)]}} \quad (3)$$

where  $\hat{V}[\cdot]$  denotes an estimate of the variance of the indices in question.

Suppose the asymptotic distribution of  $S$  under  $H_0$  is  $N[0, 1]$  (as  $m, n \rightarrow \infty$ ), and consider a critical region of the form  $|S| > c$  where  $c$  is a critical value. Then an asymptotic  $p$ -value for this test can be obtained as follows:

$$p_\alpha = G_\Phi(|S_0|) = 2 \min\left(\Phi(S_0); 1 - \Phi(S_0)\right) \quad (4)$$

where  $S_0$  is the observed value of  $S$ ,

$$G_{\Phi}(x) := \mathbb{P}[|Z| > |x|] = \Phi(-|x|) + 1 - \Phi(|x|) = 2 \min\left(\Phi(x), 1 - \Phi(x)\right), \quad (5)$$

$Z \sim N[0, 1]$ , and  $\Phi(\cdot) = \mathbb{P}[Z \leq x]$  is the standard Normal distribution function. Note the identities in (5) depend on the continuity and symmetry of the normal distribution.

We introduce three  $p$ -values, described in the following subsections. For clarity, in what follows,  $p_*$  refers to the MC-permutation  $p$ -value,  $p_b$  to the bootstrap  $p$ -value, and  $p_{\bullet}$  denotes its counterpart that imposes the null hypothesis.

## 2.1 Monte-Carlo permutational $p$ -values

The permutational distribution based on the test statistic  $T$  in (2), also known as the randomization distribution, is the distribution obtained by permuting in all possible ways the  $N = n + m$  observations of the combined sample

$$Z = (X_1, X_2, \dots, X_n, Y_1, Y_2, \dots, Y_m)'. \quad (6)$$

We denote  $\mathcal{P}(Z)$  the set of all vectors obtained by permuting the components of  $Z$ . Two permutations are viewed as distinct as soon as they correspond to different orderings of the components of  $Z$  (even if some of the observations are numerically equal), so the total number of different permutations in  $\mathcal{P}(Z)$  is  $(m + n)!$ . Under the assumption that the  $m + n$  observations in  $Z$  are i.i.d., the  $(m + n)!$  permutations in  $\mathcal{P}(Z)$  are equally probable, which in turn determines the permutational distribution of  $T$  (or  $S$ ). However, the total number of permutations  $(m + n)!$  to consider rapidly becomes prohibitive large as the sample sizes  $m$  and  $n$  increase.

Following the suggestion of Dwass (1957), we draw at random  $B$  permutations of  $Z$  from the set  $\mathcal{P}(Z)$ . These may be drawn with or without replacement [in  $\mathcal{P}(Z)$ ]. When draws are taken with replacement, the random permutations are i.i.d.; when taken without

replacement, they are exchangeable. In this paper, we focus on the case where the permutations are generated with replacement. Along with the actual data, this yields  $B + 1$  random permutations of  $Z$ :  $Z_{*1}, \dots, Z_{*B}$ . From each permuted sample, the corresponding EDFs  $\hat{F}_{x*j}$  and  $\hat{F}_{y*j}$  are computed, and the value of the test statistic as defined in (2):

$$T_{*j} = \theta(\hat{F}_{x*j}) - \theta(\hat{F}_{y*j}), \quad j = 1, \dots, B. \quad (7)$$

Using the above simulated permutational test statistics, we can then compute the following  $p$ -value functions:

$$\hat{p}_{*B}^-(x) = \frac{\sum_{j=1}^B \mathbf{1}[T_{*j} \leq x] + 1}{B + 1}, \quad \hat{p}_{*B}^+(x) = \frac{\sum_{j=1}^B \mathbf{1}[T_{*j} \geq x] + 1}{B + 1}, \quad (8)$$

where the indicator function  $\mathbf{1}(A)$  is equal to one if  $A$  is true, and zero otherwise. We can then obtain one-sided tests of  $H_0$  against  $H_1^- : \theta(F_x) < \theta(F_y)$  and  $H_1^+ : \theta(F_x) > \theta(F_y)$ , by taking the following critical regions respectively:

$$\hat{p}_{*B}^-(T) \leq \alpha, \quad (9)$$

$$\hat{p}_{*B}^+(T) \leq \alpha, \quad (10)$$

where  $\alpha$  is the level of the test and  $T$  is the observed value of the test statistic. To get a two-sided test, we can reject  $H_0$  against  $H_1 : \theta(F_x) \neq \theta(F_y)$  when either one of the one-sided tests is significant at level  $\alpha/2$ :

$$\hat{p}_{*B}^-(T) \leq \alpha/2 \quad \text{or} \quad \hat{p}_{*B}^+(T) \leq \alpha/2 \quad (11)$$

or equivalently

$$\hat{p}_{*B}^c := 2 \min\{\hat{p}_{*B}^-(T), \hat{p}_{*B}^+(T)\} \leq \alpha. \quad (12)$$

Another way of building a two-sided test consists in working with the absolute value of



the test statistic: setting

$$\hat{p}_{*B}^a(x) = \frac{\sum_{j=1}^B \mathbf{1}(|T_{*j}| \geq x) + 1}{B + 1}, \quad (13)$$

$H_0$  is rejected against  $H_1$  when

$$\hat{p}_{*B}^a(|T|) \leq \alpha. \quad (14)$$

The two critical regions in (12) and (14) are not generally equivalent. In this paper, we focus on two-sided tests of type (12). In (8)-(14), the statistic  $T$  can be replaced by its studentized version  $S$ , in which case

$$S_{*j} = \frac{\theta(\hat{F}_{x*j}) - \theta(\hat{F}_{y*j})}{\sqrt{\hat{V}[\theta(\hat{F}_{x*j})] + \hat{V}[\theta(\hat{F}_{y*j})]}}, \quad j = 1, \dots, B. \quad (15)$$

Of course, tests based on  $T$  or  $S$  are not generally equivalent.

## 2.2 Conventional bootstrap $p$ -values

A bootstrap test is computed by resampling the original data with replacement. A bootstrap sample, of the same size as the observed sample, is obtained by making  $n$  draws with replacement from the  $n$  observed realizations  $\{X_1, \dots, X_n\}$ , where each  $X_i$  has probability  $1/n$  of being selected on each draw, and then making, independently,  $m$  draws with replacement from the  $m$  observed realizations  $\{Y_1, \dots, Y_m\}$ , where each  $Y_i$  has probability  $1/m$  of being selected on each draw. Let  $(X_b, Y_b)$  refer to the bootstrap sample so obtained and denote by  $\hat{F}_{x_b}$  and  $\hat{F}_{y_b}$  the associated EDFs respectively. The bootstrap statistic is computed as was  $S$  in (3), except that the null hypothesis tested is that the difference between the two indices is equal to  $\theta(\hat{F}_x) - \theta(\hat{F}_y)$  rather than to 0. Formally, the adjusted bootstrap statistic takes the form

$$S_b = \frac{\left(\theta(\hat{F}_{x_b}) - \theta(\hat{F}_{y_b})\right) - \left(\theta(\hat{F}_x) - \theta(\hat{F}_y)\right)}{\sqrt{\hat{V}[\theta(\hat{F}_{x_b})] + \hat{V}[\theta(\hat{F}_{y_b})]}}. \quad (16)$$

This modification ensures that the hypothesis pertaining to the bootstrap statistics holds true for the population the bootstrap samples are drawn from, that is, the original sample. Let  $S_{bj}$ ,  $j = 1, \dots, B$ , refer to the series of bootstrap statistics. The bootstrap  $p$ -value is the proportion of the bootstrap samples for which the absolute value of the bootstrap statistic is more extreme than the statistic computed from the original data. Thus, for a two-tailed test, the bootstrap  $p$ -value is

$$p_b = 2 \min \left( \frac{1}{B} \sum_{j=1}^B \mathbf{1}(S_{bj} \leq S_0); \frac{1}{B} \sum_{j=1}^B \mathbf{1}(S_{bj} > S_0) \right). \quad (17)$$

### 2.3 Bootstrap $p$ -values under the null hypothesis

The permutation approach does not differ radically from the bootstrap approach. For example, a sample obtained by permuting elements of the combined sample  $Z$  defined in (6) is equivalent to resampling *without* replacement  $N$  observations from  $Z$ . It thus makes sense to resample *with* replacement from  $Z$  to form an alternative bootstrap sample that respects the null hypothesis. One can proceed as follows. Draw with replacement  $n$  observations in  $Z$  to form a sample denoted  $X_\bullet$  and then draw with replacement  $m$  other observations in  $Z$  to form a sample denoted  $Y_\bullet$ . Again, let  $\hat{F}_{x_\bullet}$  and  $\hat{F}_{y_\bullet}$  refer to the associated EDFs respectively. The bootstrap statistic can be computed as

$$S_\bullet = \frac{\theta(\hat{F}_{x_\bullet}) - \theta(\hat{F}_{y_\bullet})}{\sqrt{\hat{V}[\theta(\hat{F}_{x_\bullet})] + \hat{V}[\theta(\hat{F}_{y_\bullet})]}} \quad (18)$$

with no further adjustments since the sampling scheme imposes the null hypothesis. Let  $S_{\bullet j}$ ,  $j = 1, \dots, B$ , refer to the series of bootstrap statistics so obtained, leading to the two-tailed bootstrap  $p$ -value:

$$p_\bullet = 2 \min \left( \frac{1}{B} \sum_{j=1}^B \mathbf{1}(S_{\bullet j} \leq S_0); \frac{1}{B} \sum_{j=1}^B \mathbf{1}(S_{\bullet j} > S_0) \right). \quad (19)$$

### 3 Exact Monte-Carlo permutation tests

Before we consider the general problem of testing  $H_0$  when the populations considered may have otherwise different distributions, it will be of interest to study the problem of testing  $H_0$  against  $H_1^-$ ,  $H_1^+$  or  $H_1$ . This is equivalent to testing  $F_x = F_y$  against alternatives where  $\theta(F_x) \neq \theta(F_y)$ . This relatively restrictive null hypothesis appears naturally when subsets of a wider population are considered. We will show here that both the level and the size of permutation tests based on general statistics of the form  $T$  or  $S$  can be controlled, irrespective whether the distribution  $F_x$  (or  $F_y$ ) is discrete or continuous, without any restriction on the form of the functional  $\theta(\cdot)$ . We also allow for exchangeable (as opposed to i.i.d.) observations, hence covering the case of random draws without replacement from a finite population (in addition to i.i.d. observations). Thus, the result given here can be viewed as an extension of the basic finding of Dwass (1957) who considered tests that compare arithmetic means of i.i.d. random variables with continuous distribution.

Since most estimated inequality measures rely on statistics based on EDFs, which are not continuous, a tie-breaking procedure may be needed to control test size. For this purpose, we propose to use the randomized ordering described in Dufour (2006), which leads to the following procedure. We first draw by simulation  $B + 1$  i.i.d. random variables  $U_0, U_1, \dots, U_B$ , according to a uniform distribution on  $(0, 1)$ , independently of  $(T, T_{*1}, \dots, T_{*B})$ . Then we compute  $p$ -value functions similar to those described in Section 2.1, with the difference that the pairs  $(T_{*j}, U_j)$ ,  $j = 0, \dots, B$ , are ordered according to the lexicographic order:

$$(T_{*i}, U_i) \leq (T_{*j}, U_j) \Leftrightarrow [T_{*i} < T_{*j} \quad \text{or} \quad (T_{*i} = T_{*j} \quad \text{and} \quad U_i \leq U_j)] \quad (20)$$

where  $T_{*0} = T$  is the statistic obtained from the actual data. More precisely, this yields the following modified (*tie-adjusted*)  $p$ -value functions:

$$\tilde{p}_{*B}^-(x) = \frac{\sum_{i=1}^B \mathbf{1}[(T_{*i}, U_i) \leq (x, U_0)] + 1}{B + 1}, \quad (21)$$

$$\tilde{p}_{*B}^+(x) = \frac{\sum_{i=1}^B \mathbf{1}[(x, U_0) \leq (T_{*i}, U_i)] + 1}{B + 1}, \quad (22)$$

$$\tilde{p}_{*B}^a(x) = \frac{\sum_{i=1}^B \mathbf{1}[(x, U_0) \leq (|T_{*i}|, U_i)] + 1}{B + 1}. \quad (23)$$

The tests are performed as before on replacing  $\hat{p}$  by  $\tilde{p}$ . We can then establish the following theorem.

**Theorem 1** *Suppose the  $n + m$  random variables  $X_1, X_2, \dots, X_n, Y_1, Y_2, \dots, Y_m$  are exchangeable. Then, for  $0 < \alpha < 1$ ,*

$$\mathbb{P}[\hat{p}_{*B}^-(T) \leq \alpha] \leq \mathbb{P}[\tilde{p}_{*B}^-(T) \leq \alpha] = \alpha, \quad (24)$$

$$\mathbb{P}[\hat{p}_{*B}^+(T) \leq \alpha] \leq \mathbb{P}[\tilde{p}_{*B}^+(T) \leq \alpha] = \alpha, \quad (25)$$

$$\mathbb{P}[\hat{p}_{*B}^a(|T|) \leq \alpha] \leq \mathbb{P}[\tilde{p}_{*B}^a(|T|) \leq \alpha] = \alpha, \quad (26)$$

where the  $p$ -value functions are defined in (8) and (21) - (23).

**Proof:** By the exchangeability assumption, different permutations of the components of  $Z$  are equally probable, so randomly selected permutations (either with or without replacement) are themselves exchangeable. Consequently, the random variables  $T, T_{*1}, \dots, T_{*B}$  are exchangeable. On applying Proposition 2.4 of Dufour (2006), we then get

$$\mathbb{P}[\tilde{p}_{*B}^-(T) \leq \alpha] = \mathbb{P}[\tilde{p}_{*B}^+(T) \leq \alpha] = \mathbb{P}[\tilde{p}_{*B}^a(|T|) \leq \alpha] = \alpha. \quad (27)$$

Finally, the inequalities in (24) - (26) follow on observing that

$$\tilde{p}_{*B}^-(x) \leq \hat{p}_{*B}^-(x), \quad \tilde{p}_{*B}^+(x) \leq \hat{p}_{*B}^+(x), \quad \tilde{p}_{*B}^a(x) \leq \hat{p}_{*B}^a(x). \quad (28)$$

■

Theorem 1 means that the critical regions  $\tilde{p}_{*B}^-(T) \leq \alpha$ ,  $\tilde{p}_{*B}^+(T) \leq \alpha$  and  $\tilde{p}_{*B}^a(|T|) \leq \alpha$  have size  $\alpha$  for testing  $H_0$ , while the critical regions  $\hat{p}_{*B}^-(T) \leq \alpha$ ,  $\hat{p}_{*B}^+(T) \leq \alpha$  and

$\hat{p}_{*B}^a(|T|) \leq \alpha$  are typically conservative, so they still have level  $\alpha$  for testing  $H_0$ . Clearly, the same result holds if the test statistic  $T$  is replaced by its studentized version  $S$ .

Note also the events  $\tilde{p}_{*B}^-(T) \leq \alpha/2$  and  $\tilde{p}_{*B}^+(T) \leq \alpha/2$  are mutually exclusive [and similarly for  $\hat{p}_{*B}^-(T) \leq \alpha/2$  and  $\hat{p}_{*B}^+(T) \leq \alpha/2$ ] for  $0 < \alpha < 1$ , so that

$$\mathbb{P}[\hat{p}_{*B}^-(T) \leq \alpha/2 \text{ or } \hat{p}_{*B}^+(T) \leq \alpha/2] \leq \mathbb{P}[\tilde{p}_{*B}^-(T) \leq \alpha/2 \text{ or } \tilde{p}_{*B}^+(T) \leq \alpha/2] = \alpha \quad (29)$$

under  $H_0$ , for  $0 < \alpha < 1$ . Thus, on setting

$$\hat{p}_{*B}^c(x) := 2 \min\{\hat{p}_{*B}^-(x), \hat{p}_{*B}^+(x)\}, \quad \tilde{p}_{*B}^c(x) := 2 \min\{\tilde{p}_{*B}^-(x), \tilde{p}_{*B}^+(x)\}, \quad (30)$$

we have

$$\mathbb{P}[\hat{p}_{*B}^c(T) \leq \alpha] \leq \mathbb{P}[\tilde{p}_{*B}^c(T) \leq \alpha] = \alpha \quad (31)$$

under  $H_0$ , so  $\hat{p}_{*B}^c(T) \leq \alpha$  and  $\tilde{p}_{*B}^c(T) \leq \alpha$  constitute two-sided critical regions with level  $\alpha$  for  $H_0$  ( $0 < \alpha < 1$ ). The latter are not in general equivalent to  $\hat{p}_{*B}^a(|T|) \leq \alpha$  or  $\tilde{p}_{*B}^a(|T|) \leq \alpha$ .

## 4 Permutation tests for comparing linear functionals: different distributions

In this section, we identify specific finite-sample permutation test problems and extend the results in Romano (1990) on the asymptotic validity of permutation tests (2) to the case where  $\theta(\cdot)$  is a linear functional. The latter result will be used in the next section to show that permutation tests are asymptotically valid with inequality measures if the samples are previously rescaled adequately. Permutations based on studentized statistics [such as (3) here] have been shown to be valid under rather general regularity conditions by Chung and Romano (2013). The following analysis which focuses on (2) is nevertheless informative, as the properties we derive are new to this literature. More to the point, the

transformation we introduce to validate (3) ends up enhancing power when combined with studentization.

In general, if the assumption  $F_x = F_y$  does not hold, the permutation test of  $H_0$  in (1) is no longer of  $\alpha$ -level. However, as pointed out by Romano (1990), permutation tests can be asymptotically valid in some specific cases, in the sense that under the null hypothesis, the rejection frequency of the permutation test tends to the nominal level  $\alpha$  as the sample size increases.

Romano (1990) shows that, in a general two-sample problem of testing  $\theta(F_x) = \theta(F_y)$ , the asymptotic validity of the permutation test follows when the asymptotic variance of the original statistic  $n^{1/2}[\theta(\hat{F}_x) - \theta(\hat{F}_y)]$  and the asymptotic variance of the permutation statistic are the same (both statistics converge weakly to a Gaussian distribution with mean 0). This result requires that  $\theta(\cdot)$  is appropriately differentiable in the sense of Gill (1988). It leads him to conclude that the validity of the permutation test does not hold in general for two sample problems, while it is generally asymptotically valid for certain one-sample problems. In his Theorem 3.3, Romano (1990) shows that, with two independent samples,  $\{X_1, \dots, X_n\}$  and  $\{Y_1, \dots, Y_m\}$ , drawn from the probability distributions  $F_x$  and  $F_y$ , the test statistic  $n^{1/2}[\theta(\hat{F}_x) - \theta(\hat{F}_y)]$  is asymptotically Gaussian with mean 0 and variance

$$V_{as} \left[ \theta(\hat{F}_x) \right] + \frac{1 - \lambda}{\lambda} V_{as} \left[ \theta(\hat{F}_y) \right] \quad (32)$$

where  $V_{as}[\theta(\hat{F}_x)]$  and  $V_{as}[\theta(\hat{F}_y)]$  are the asymptotic variances of, respectively,  $n^{1/2}[\theta(\hat{F}_x) - \theta(F_x)]$  and  $m^{1/2}[\theta(\hat{F}_y) - \theta(F_y)]$ , and  $m/(m + n) \rightarrow \lambda$  as  $n \rightarrow \infty$ ,  $\lambda > 0$ . Moreover, he shows that the distribution of the permutation test behaves asymptotically as Gaussian with mean 0 and variance

$$V_{as} \left[ \theta \left( (1 - \lambda)\hat{F}_x + \lambda\hat{F}_y \right) \right] + \frac{1 - \lambda}{\lambda} V_{as} \left[ \theta \left( (1 - \lambda)\hat{F}_x + \lambda\hat{F}_y \right) \right]. \quad (33)$$

This result is useful to consider asymptotic validity of a permutation test when the vari-

ances (32) and (33) are the same, that is, when the following condition holds:

$$V_{as} \left[ \theta \left( (1 - \lambda) \hat{F}_x + \lambda \hat{F}_y \right) \right] = \lambda V_{as} \left[ \theta(\hat{F}_x) \right] + (1 - \lambda) V_{as} \left[ \theta(\hat{F}_y) \right]. \quad (34)$$

If this condition holds, the permutation test is asymptotically valid, in the sense that the distribution of the permutation test behaves asymptotically as the unconditional sampling distribution of  $n^{1/2}[\theta(\hat{F}_x) - \theta(\hat{F}_y)]$ . The distribution of the permutation test can then be used to compute a critical value or a  $p$ -value.

Looking at situations where (34) holds, Romano shows that for comparing means, the permutation test is asymptotically valid only if the sample sizes are (approximately) the same or the distributions have (approximately) the same variances. Moreover, he argues that this condition does not hold in general. Looking more closely at situations where (34) holds, we derive the following result, which will be useful to show that permutation tests may be valid in several useful situations, such as comparing moments or inequality measures between two samples.

This result assumes that the considering functionals are mixture-linear and the available estimators are asymptotically linear. These concepts are defined in what follows:

**Definition 1** *A mixture-linear functional is a mapping from a subset of distributions on the set of real numbers,  $\theta : D \rightarrow \mathbb{R}$ , satisfying the following condition*

$$\theta \left( \sum_{k=1}^K \lambda_k F_k \right) = \sum_{k=1}^K \lambda_k \theta(F_k) \quad (35)$$

where  $F_1, \dots, F_K$  are distributions and  $\lambda_1, \dots, \lambda_K$  are fixed positive scalars such that  $\sum_{k=1}^K \lambda_k = 1$ .

**Definition 2** *Given a sample  $Z = \{Z_1, \dots, Z_n\}$  of independent observations drawn from distribution  $F_z$ , an estimator of a functional  $\theta(\hat{F}_z)$  is asymptotically linear if*

$$n^{1/2}[\theta(\hat{F}_z) - \theta(F_z)] = \frac{1}{\sqrt{n}} \sum_{i=1}^n g_z(Z_i) + o_{F_z}(1) \quad (36)$$

where the function  $g_z$  can depend on the underlying distribution.

**Theorem 2** Consider the problem of testing

$$H_0 : \theta(F_x) = \theta(F_y) \quad (37)$$

from two samples  $X = \{X_1, \dots, X_n\}$  and  $Y = \{Y_1, \dots, Y_m\}$  of independent observations, drawn, respectively, from distributions  $F_x$  and  $F_y$ , using estimators  $\theta(\hat{F}_x)$  and  $\theta(\hat{F}_y)$  that are asymptotically linear, that is, satisfy (36), not just for i.i.d. samples under  $F_x$  and  $F_y$ , but also when sampling i.i.d. observations from the mixture distribution  $\lambda F_x + (1 - \lambda)F_y$ , and assume finite second moments. A permutation test (7) is asymptotically valid if  $\theta(\cdot)$  is a mixture-linear functional, as in (35), and, in addition, at least one of the following conditions hold:

$$m/(m+n) \xrightarrow[n \rightarrow \infty]{} \lambda = 1/2 \quad (38)$$

$$V_{as}[\theta(\hat{F}_x)] = V_{as}[\theta(\hat{F}_y)] \quad (39)$$

where  $V_{as}[\theta(\hat{F}_x)]$  and  $V_{as}[\theta(\hat{F}_y)]$  are the asymptotic variances of, respectively,  $n^{1/2}[\theta(\hat{F}_x) - \theta(F_x)]$  and  $m^{1/2}[\theta(\hat{F}_y) - \theta(F_y)]$ .

**Proof:** Romano (1990) shows that the distribution of the permutation test based on the statistic  $n^{1/2}[\theta(\hat{F}_x) - \theta(\hat{F}_y)]$  is asymptotically Gaussian.<sup>1</sup> With a linear functional  $\theta(\cdot)$  and independent samples, we have:<sup>2</sup>

$$V \left[ \theta \left( \sum_{k=1}^K \lambda_k \hat{F}_k \right) \right] = V \left[ \sum_{k=1}^K \lambda_k \theta(\hat{F}_k) \right] = \sum_{k=1}^K \lambda_k V \left[ \theta(\hat{F}_k) \right].$$

---

<sup>1</sup>Romano (1990) relies on differentiability conditions, whereas Chung and Romano (2013) use the weaker condition (36) instead.

<sup>2</sup>With independent samples,  $\text{Cov}[\theta(\hat{F}_k), \theta(\hat{F}_l)] = 0$ , for  $k \neq l$ ,  $k, l = 1, \dots, K$ .



For  $K = 2$ ,  $F_1 = F_x$  and  $F_2 = F_y$ , we have

$$V \left[ \theta \left( \lambda_1 \hat{F}_x + \lambda_2 \hat{F}_y \right) \right] = \lambda_1 V[\theta(\hat{F}_x)] + \lambda_2 V[\theta(\hat{F}_y)], \quad (40)$$

Let us define  $\lambda_2 = m/(m+n) \rightarrow \lambda$  as  $n \rightarrow \infty$  and  $\lambda_1 = n/(m+n) = 1 - \lambda_2$ , we obtain

$$V \left[ \theta \left( \lambda_1 \hat{F}_x + \lambda_2 \hat{F}_y \right) \right] \xrightarrow{n \rightarrow \infty} (1 - \lambda) V_{as}[\theta(\hat{F}_x)] + \lambda V_{as}[\theta(\hat{F}_y)], \quad (41)$$

The condition defined in (34) is then satisfied if  $\lambda = 1/2$  or  $V_{as}[\theta(\hat{F}_x)] = V_{as}[\theta(\hat{F}_y)]$ .

■

Note that the (arithmetic) mean is a linear functional, but the quantile is not a linear functional. Comparing means from samples of similar size with a permutation test is then asymptotically valid, even if the underlying distributions are not identical, while comparing quantiles with a permutation test is no longer valid, in general, if the underlying distributions are not identical.

## 5 Comparing inequality measures

This section focuses on identifying functionals  $\theta(\cdot)$  of interest for which permutation tests (7) are valid in the sense of Theorem 2. We study moments and inequality measures for which condition (34) would hold.

### 5.1 Centered and uncentered moments

Consider the functional

$$\theta(F_z) = \int \phi(z) dF_z(z), \quad (42)$$

where  $\phi(\cdot)$  is any function in  $\mathbb{R}$  for which  $E[\phi(\cdot)]$  exists. For any random variable  $w$  that follows a mixture of  $K$  distributions

$$w \sim \sum_{k=1}^K \lambda_k F_k(w),$$

if we consider  $K$  random variables  $w_1, \dots, w_K$  from the  $K$  component distributions then

$$\theta(F_w) = E[\phi(w)] = \sum_{k=1}^K \lambda_k E[\phi(w_k)] = \sum_{k=1}^K \lambda_k \theta(F_{w_k}). \quad (43)$$

The functional (42) is linear and Theorem 2 applies: permutation tests are asymptotically valid if either (38) or (39) holds.

Then, comparing *uncentered* moments between two samples with permutation tests is asymptotically valid, it corresponds to the special case of (42) with  $\phi(z) = z^r$ , where  $r$  is a positive integer. The mean corresponds to the case  $\phi(z) = z$ . Turning to *centered* moments, consider

$$\theta(F_z) = \int [z - E(z)]^r dF_z(z), \quad (44)$$

where  $r$  is an integer greater than 1. With  $w$  and  $w_1, \dots, w_K$  as in (43), we have

$$\theta(F_w) = E\left([w - E(w)]^r\right) \quad (45)$$

$$= \sum_{k=1}^K \lambda_k E\left(\left[w_k - E(w_k) + E(w_k) - E(w)\right]^r\right). \quad (46)$$

The two last terms in parenthesis vanish if they are equal leading to:

$$\theta(F_w) = \sum_{i=1}^K \lambda_k E\left([w_k - E(w_k)]^r\right) \quad (47)$$

$$= \sum_{i=1}^K \lambda_k \theta(F_{w_k}) \quad \text{if } E(w_k) = E(w), \forall k. \quad (48)$$

This result suggest that (44) is not a linear functional, unless the component distributions

share a common mean. From Theorem 2, comparing centered moments from two samples with a permutation test (7) is then invalid, unless the samples come from distributions with the same mean,  $\mu(F_x) = \mu(F_y)$ , and either (38) or (39) holds.

However, centered moments are translation invariant: calculating centered moments from the original samples or from the *centered* samples,

$$\left\{X_1 - \mu(F_x), \dots, X_n - \mu(F_x)\right\} \quad \text{and} \quad \left\{Y_1 - \mu(F_y), \dots, Y_m - \mu(F_y)\right\}, \quad (49)$$

gives the same result. The main issue here is that the centered samples have a common mean, which is equal to zero, and the statistic (44) is a linear functional in this particular case. Comparing centered moments from the two *centered* samples rather than from the original samples makes no difference, while it validates (asymptotically) the use of permutation test.

In practice,  $\mu(F_x)$  and  $\mu(F_y)$  are often unknown. Permutation tests can nevertheless be applied on the combined sample

$$Z_c = \{X_1 - \bar{X}, \dots, X_n - \bar{X}, Y_1 - \bar{Y}, \dots, Y_m - \bar{Y}\} \quad (50)$$

where  $\bar{X} = n^{-1} \sum_{i=1}^n X_i$  and  $\bar{Y} = m^{-1} \sum_{i=1}^m Y_i$  are the sample means of each sample. This procedure is known to perform well in finite sample when testing the equality of variances from two samples, see Lim and Loh (1996) and Boos and Brownie (2004).

## 5.2 The generalized entropy class

We consider the important **GE** class of inequality measures, defined by equations (51)-(53):

$$\theta_{\text{GE}}^\zeta(F) = \frac{1}{\zeta^2 - \zeta} \left[ \int \left[ \frac{y}{\mu(F)} \right]^\zeta dF(y) - 1 \right], \quad \zeta \in \mathbb{R}, \zeta \neq 0, 1, \quad (51)$$

$$\theta_{\text{GE}}^0(F) = - \int \log \left[ \frac{y}{\mu(F)} \right] dF(y), \quad (52)$$

$$\theta_{\text{GE}}^1(F) = \int \frac{y}{\mu(F)} \log \left[ \frac{y}{\mu(F)} \right] dF(y). \quad (53)$$

The parameter  $\zeta$  of the GE class characterizes the sensitivity to differences over different segments of the distribution. The more positive (negative)  $\zeta$  is, the more sensitive is the inequality measure to differences at the top (bottom) of the distribution. The Mean Logarithmic Deviation (MLD) index,  $\theta_{\text{GE}}^0(F)$ , is the limiting case when  $\zeta = 0$ . The Theil index,  $\theta_{\text{GE}}^1(F)$ , is the limiting case of the GE when  $\zeta = 1$ .

The GE class of inequality measures is decomposable, that is, it can be expressed as a simple additive function of within-group and between-group inequality. Let there be  $K$  groups and let the proportion of the population falling in group  $k$  be  $\lambda_k$ ; then the class of GE indices is equal to

$$\theta_{\text{GE}}^\zeta(\hat{F}_w) = \sum_{k=1}^K \lambda_k \left[ \frac{\bar{w}_k}{\bar{w}} \right]^\zeta \theta_{\text{GE}}^\zeta(\hat{F}_{w_k}) - \frac{1}{\zeta^2 - \zeta} \left( \sum_{k=1}^K \lambda_k \left[ \frac{\bar{w}_k}{\bar{w}} \right]^\zeta - 1 \right) \quad (54)$$

where  $\bar{w}_k$  is the mean income in group  $k$ ,  $\bar{w}$  is the mean income of the population

$$\bar{w} = K^{-1} \sum_{k=1}^K \lambda_k \bar{w}_k$$

and  $\theta_{\text{GE}}^\zeta(\hat{F}_{w_k})$  is the GE index in group  $k$ , see Cowell (2011). It is clear that

$$\theta_{\text{GE}}^\zeta(\hat{F}_w) = \sum_{k=1}^K \lambda_k \theta_{\text{GE},k}^\zeta(\hat{F}_{w_k}) \quad \text{if} \quad \bar{w}_k = \bar{w}, \forall k. \quad (55)$$

It follows that  $\theta_{\text{GE}}^\zeta(\hat{F}_w)$  is not a linear functional, unless the mean in each group is the same. From Theorem 2, comparing GE inequality measures from two samples with permutation tests (7) is then valid only if the samples come from distributions with the same mean,  $\mu(F_x) = \mu(F_y)$ , and either (38) or (39) holds.

As is clear from equations (51)-(53), the GE class of inequality measures is scale invariant, which suggests to base a permutation test on the *rescaled* samples, where the

observations are divided by their distributional mean,

$$\left\{ \frac{X_1}{\mu(F_x)}, \dots, \frac{X_n}{\mu(F_x)} \right\} \quad \text{and} \quad \left\{ \frac{Y_1}{\mu(F_y)}, \dots, \frac{Y_n}{\mu(F_y)} \right\}. \quad (56)$$

Comparing Generalized Inequality indices from these rescaled samples rather than from the original samples makes no differences, while it validates (asymptotically) the use of permutation test. In practice, distributional means are often unknown; we thus use sample means  $\bar{X}$  and  $\bar{Y}$  instead, so the permutation test is based on the following combined sample

$$Z_s = \left\{ \frac{X_1}{\bar{X}}, \dots, \frac{X_n}{\bar{X}}, \frac{Y_1}{\bar{Y}}, \dots, \frac{Y_m}{\bar{Y}} \right\}, \quad \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i. \quad (57)$$

It is worth noting that when we consider the *rescaled* samples (56), the GE inequality measures can be rewritten as a moment  $\theta_{\text{GE}}^\zeta(F_z) = \int \phi(z) dF_z(z)$ , as defined in (42), where

$$\phi(z) = \begin{cases} (z^\zeta - 1)/(\zeta^2 - \zeta) & \text{for } \zeta \neq 0, 1 \\ -\log z & \text{for } \zeta = 0 \\ z \log z & \text{for } \zeta = 1 \end{cases} \quad (58)$$

which leads us back to the results of section 5.1.

Clearly the same approach can be applied to the Atkinson class of inequality indices (Atkinson, 1970),

$$\theta_{\text{Atk}}^\zeta(F) = 1 - \left[ \int \left[ \frac{y}{\mu(F)} \right]^\zeta dF(y) \right]^{1/\zeta}, \quad \zeta < 1 \quad (59)$$

which can be rewritten as a function of the GE class of inequality measures

$$\theta_{\text{Atk}}^\zeta(F) = \begin{cases} 1 - [(\zeta^2 - \zeta)\theta_{\text{GE}}^\zeta(F) + 1]^{1/\zeta}, & \zeta \neq 0, \\ 1 - \exp(-\theta_{\text{GE}}^0(F)), & \zeta = 0. \end{cases} \quad (60)$$

### 5.3 The Gini coefficient

The Gini index can be expressed in a number of different forms. Let us consider the following expressions,

$$\theta_{\text{Gini}}(F) = \frac{1}{2\mu(F)} \iint |y_1 - y_2| dF(y_1)dF(y_2) = \frac{E(|y_1 - y_2|)}{2\mu(F)}, \quad (61)$$

$$\theta_{\text{Gini}}(F) = 1 - 2 \int_0^1 L(F; q) dq, \quad (62)$$

where  $y_1$  and  $y_2$  are two random variables independently drawn from  $F$ , and  $L(F; q)$  is the  $q$ th ordinate of the Lorenz curve. Equation (61) presents the Gini as the normalized average absolute difference between all the possible pairs of incomes in the population, while equation (62) shows that the Gini index is twice the area between the Lorenz curve and the  $45^\circ$  line.

The Gini index is also closely related to a measure of dispersion of a distribution. The most popular measure of dispersion is the standard deviation, which is the square root of the variance that can be rewritten as follows:

$$V(y) = E[(y - \mu(F))^2] = E \left[ \frac{1}{2}(y_1 - y_2)^2 \right]. \quad (63)$$

Another well-known measure of dispersion is the Gini's mean difference,

$$\Delta(F) = E(|y_1 - y_2|). \quad (64)$$

Both measures of dispersion are translation invariant. In section 5.1, we prove that testing the equality of variances from two samples with different means can be done with permutation tests based on the combined sample of the recentered individual samples. Boos et al. (1989) proved that this procedure is asymptotically correct for a large class of  $U$ -statistics, from which the Gini's mean difference is a special case. We can then use the relationship between the Gini (inequality) index and the Gini's mean difference to justify

asymptotically the use of permutation test with the Gini index. Indeed, we have

$$\Delta(F) = 2\mu(F)\theta_{\text{Gini}}(F). \quad (65)$$

With  $\mu(F) = 1$ , the Gini's mean difference is twice the Gini index. Comparing the Gini's mean difference or the Gini index from two samples is then equivalent if the underlying distributions share a common mean equal to one.

The last condition does not hold in general. However, the Gini index is scale invariant. Then, calculating Gini index from the original samples or from the *rescaled* samples, where the observations are divided by their distributional mean,

$$\left\{ \frac{X_1}{\mu(F_x)}, \dots, \frac{X_n}{\mu(F_x)} \right\} \quad \text{and} \quad \left\{ \frac{Y_1}{\mu(F_y)}, \dots, \frac{Y_n}{\mu(F_y)} \right\}, \quad (66)$$

gives the same results. The main issue here is that these rescaled samples share a common mean, equals to one. Comparing Gini inequality measures from the two *rescaled* samples in (66) rather than from the original samples makes no difference for scale invariant statistic, while it validates asymptotically the use of permutation test. In practice, distributional means are replaced by sample means and permutation tests are based on the combined sample of empirically rescaled individual data  $Z_s$  as defined in (57).

## 6 Simulation study

Overall, we focus our simulation study to extreme cases of (very) heavy-tailed distributions in (very) small sample to stress-test the methods employed in testing. The heavy-tailed distribution used as a benchmark in previous studies is a more favorable case here, and we use much more heavy-tailed distributions with a very small number observations in each samples. In our experiments, we test the equality of Gini and Theil inequality measures between two samples.

## 6.1 Model design

We make use of simulated data sets drawn from the Singh-Maddala distribution, which can quite successfully mimic observed income distributions in various countries (McDonald, 1984; Kleiber and Kotz, 2003). The CDF of the Singh-Maddala distribution can be written as

$$\text{SM}(x; a, b, q) : \quad F(x) = 1 - \left[ 1 + \left( \frac{x}{b} \right)^a \right]^{-q}, \quad (67)$$

where  $a, b, q$  are positive,  $b$  is a scale parameter and  $a, q$  are shape parameters;  $q$  only affects the right tail, whereas  $a$  affects both tails. The  $k$ th moment exists for  $-a < k < aq$ . The upper-tail of the Singh-Maddala distribution behaves like a Pareto distribution with a tail index equal to  $\xi = aq$  (Schluter and Trede, 2002). Smaller is  $\xi$ , heavier is the upper tail of the distribution.

As a benchmark, we use the parameter values  $a = 2.8$ ,  $b = 100^{-\frac{1}{2.8}}$ ,  $q = 1.7$ , a choice that closely mimics the net income distribution of German households (Brachman et al., 1996). This distribution is used in Davidson and Flachaire (2007) and Cowell and Flachaire (2007) to show poor finite-sample performance of asymptotic and bootstrap inference. Its tail index is equal to  $\xi = aq = 4.76$ . We will depart from this distribution using heaviest-tailed distributions (Singh-Maddala distributions with smaller tail parameters  $\xi$ ) for which we know that bootstrap inference is poorest.

We compute the Theil index as

$$\theta_{\text{GE}}^1(\hat{F}_y) = \frac{1}{n} \sum_{i=1}^n \frac{y_i}{\hat{\mu}} \log \left( \frac{y_i}{\hat{\mu}} \right) \quad (68)$$

where  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$ . The variance of the Theil index is computed as  $\widehat{\text{var}}(\theta_{\text{GE}}^1(\hat{F}_y)) = \frac{1}{n^2} \sum_{i=1}^n (Z_i - \bar{Z})^2$ , where

$$Z_i = \frac{y_i}{\hat{\mu}} \left[ \log \left( \frac{y_i}{\hat{\mu}} \right) - \theta_{\text{GE}}^1(\hat{F}_y) - 1 \right], \quad (69)$$



and  $\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i$ , see Cowell and Flachaire (2014).

We compute the Gini index as

$$\theta_{\text{Gini}}(\hat{F}_y) = \frac{2 \sum_{i=1}^n i y_{(i)}}{\hat{\mu} n (n-1)} - \frac{n+1}{n-1} \quad (70)$$

where the  $y_{(i)}$ ,  $i = 1, \dots, n$  are the order statistics ( $y_{(1)} \leq \dots \leq y_{(n)}$ ) and  $\hat{\mu}$  is the sample mean,  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$ . The variance of the Gini index is computed as  $\widehat{\text{var}}(\theta_{\text{Gini}}(\hat{F}_y)) = \frac{1}{(n\hat{\mu})^2} \sum_{i=1}^n (Z_i - \bar{Z})^2$ , where

$$Z_i = - \left( \theta_{\text{Gini}}(\hat{F}_y) + 1 \right) y_{(i)} + \frac{2i-1}{n} y_{(i)} - \frac{2}{n} \sum_{j=1}^i y_{(j)} \quad (71)$$

and  $\bar{Z} = \frac{1}{n} \sum_{i=1}^n Z_i$ , see Davidson (2009) and Cowell and Flachaire (2014).

Our results are presented with figures, with the following legend:

- **asymptotic**: asymptotic test.
- **bootstrap**: standard bootstrap test  $S_b$  (defined in (16))
- **Perm T, rescaled**: permutation test  $T_*$  based on  $Z_s$  (defined in (7) and (57))
- **Perm S, rescaled**: permutation test  $S_*$  based on  $Z_s$  (defined in (15) and (57))
- **Perm S, standard**: permutation test  $S_*$  based on  $Z$  (defined in (15) and (6))
- **Boot S, rescaled**: bootstrap test  $S_\bullet$  based on  $Z_s$  (defined in (18) and (57))
- **Boot S, standard**: bootstrap test  $S_\bullet$  based on  $Z$  (defined (18) and (6))

The number of replications is equal to 10,000. The number of bootstrap and permutation samples are similar,  $B = 999$ . The sample size of both samples  $X$  and  $Y$  is the same,  $n = m$ . The permutation and bootstrap  $p$ -values are obtained as described above. We compute the rejection probability, or rejection frequency, as the proportion of  $p$ -value less than a nominal level equals to 0.05.

## 6.2 Size

In the experiments, we consider several Singh-Maddala distributions for which the Theil inequality measure index is the same and the tail index varies,  $\xi \in [2.9, 6.26]$ .<sup>3</sup> The Singh-Maddala distribution with  $\xi = 2.9$  is then the heaviest-tailed distribution considered here. Similar experiments are conducted for the Gini inequality measure, with slightly different tail parameters,  $\xi \in [2.59, 6.6]$ .<sup>4</sup> Inference is exact if the rejection probability is equal to 0.05.

### 6.2.1 Identical distributions

Figure 2 shows empirical rejection frequencies for asymptotic, bootstrap and permutation tests for the Theil index, when  $F_x = F_y$ , as the upper tail becomes heavier. The sample size is very small  $n = 50$ . Figure 3 shows similar results for the Gini index. When the upper-tail of the distribution becomes heavier (as  $\xi_y$  decreases), asymptotic and standard bootstrap tests perform very poorly, while permutation and bootstrap under the null tests based on the studentized statistic (**perm S rescaled**, **perm S standard**, **boot S rescaled** and **boot S standard**) provide empirical frequencies almost equal to 0.05. Note that studentized permutation tests based on the combined original samples (**perm S standard**) provides exact inference - as shown by (Chung and Romano, 2013), not permutation tests based on the combined rescaled samples (**perm S rescaled**): it is because samples are previously divided by sample means rather than by distributional means.

---

<sup>3</sup>Singh-Maddala distributions with parameters  $(a, q)$  equal to (2.5, 2.502199), (2.6, 2.149747), (2.7, 1.894309), (2.8, 1.7), (3.0, 1.4223847), (3.2, 1.2320215), (3.4, 1.0922125), (3.8, 0.8984488), (4.8, 0.6366578) and (5.8, 0.4996163), share the same (scale-invariant) Theil index, equal to 0.1401151. The tail parameters are, respectively, equal to  $\xi = 6.26, 5.59, 5.11, 4.76, 4.27, 3.94, 3.71, 3.41, 3.06, 2.9$ .

<sup>4</sup>Singh-Maddala distributions with parameters  $(a, q)$  equal to (2.5, 2.640350), (2.6, 2.218091), (2.7, 1.920967), (2.8, 1.7), (3.0, 1.3921126), (3.2, 1.1866026), (3.4, 1.0388049), (3.8, 0.8387663), (4.8, 0.5784599) and (5.8, 0.4473111), share the same (scale-invariant) Gini index, equals to 0.2887138. The tail parameters are, respectively, equal to  $\xi = 6.6, 5.77, 5.19, 4.76, 4.18, 3.80, 3.53, 3.19, 2.78, 2.59$ .

### 6.2.2 Different distributions

We then generate samples from different distributions,  $F_x \neq F_y$ , with the same value of the inequality index. Figure 4 shows rejection frequencies for asymptotic, bootstrap and permutation tests for the Theil index, as the distribution  $F_y$  moves away from  $F_x$ . The distribution  $F_x$  is fixed, with a tail index  $\xi_x = 4.76$ , while  $F_y$  has varying tail indices. When the tail index of  $F_y$  is smaller (higher) than that of  $F_x$ , that is, when  $\xi_y < \xi_x$ ,  $F_y$  is more (less) heavy-tailed than  $F_x$ . Figure 5 shows similar results for the Gini index. From these Figures, we can see that the results deteriorate when  $F_y$  tends to be much more heavy-tailed than  $F_x$ , that is, when  $\xi_y < 3.5$ . Overall, permutation and bootstrap under the null tests based on the studentized statistic (perm S rescaled, perm S standard, boot S rescaled and boot S standard) perform similarly and they outperform other methods. They perform very well when  $\xi_y > 3.5$ , that is, when  $F_y$  is not much more heavy-tailed than  $F_x$ .

### 6.2.3 Sample size

Figure 6 shows rejection frequencies for the Theil measure, as the sample size increases ( $n = 50, \dots, 10.000$ ), with identical distributions ( $F_x = F_y$ ), in the worst case previously studied ( $\xi_x = \xi_y = 2.9$ ). Figure 7 shows similar results for the Gini index (with  $\xi_x = \xi_y = 2.59$ ). We can see that the rejection frequencies decrease slowly as the sample size increases with asymptotic tests, and, even more slowly with standard bootstrap tests. In contrast, permutation tests and bootstrap under the null perform very well in all cases when they are based on studentized statistic: rejection frequencies are always almost equal to 0.05 for perm S rescaled, perm S standard, boot S rescaled and boot S standard.

Figure 8 shows rejection frequencies for the Theil measure, as the sample size increases ( $n = 50, \dots, 10.000$ ) with different distributions ( $F_x \neq F_y$ ), in the worst cases previously studied ( $\xi_x = 4.76, \xi_y = 2.9$ ). Figure 9 shows similar results for the Gini index (with  $\xi_x = 4.76, \xi_y = 2.59$ ). We can see that, for each method, the rejection frequencies decrease very slowly as the sample size increases. Moreover, permutation and bootstrap under the

null tests based on a studentized statistic outperform other methods.

### 6.3 Power

To study the power, we test the equality of an inequality measure between two samples, when the samples come from two distributions with different values of the inequality measure. From the study on the size, studentized permutation and bootstrap under the null tests outperform other methods. They also perform similarly when the null hypothesis is true, we can thus compare power between these methods.

In our experiments, the distribution  $F_x$  is fixed and the distribution  $F_y$  varies:

$$F_x = \text{SM}(x; 2.8, 0.1930698, 1.7)$$

$$F_y = \text{SM}(y; 2.8, 0.1930698, q) \quad \text{where } q \in [0.8; 31.7]$$

As  $q$  increases, the tail index and the inequality measure increase.<sup>5</sup>

Figure 10 shows rejection frequencies for testing the equality of the Theil measure between two samples, when the true null hypothesis,  $\theta(F_y) - \theta(F_x)$ , goes away from 0. We consider a very small sample ( $n = 50$ , on the left) and a moderate sample ( $n = 200$ , on the right). In the Figure, we shall consider different cases:

- $\theta(F_y) - \theta(F_x) = 0$ : the two distributions are identical and  $H_0$  is true (Size).
- $\theta(F_y) - \theta(F_x) \neq 0$ : the two distributions are different and  $H_0$  is not true (Power).
- $\theta(F_y) - \theta(F_x) < 0$ :  $F_y$  is less heavy-tailed than  $F_x$  ( $\xi_y > \xi_x$ ).
- $\theta(F_y) - \theta(F_x) > 0$ :  $F_y$  is more heavy-tailed than  $F_x$  ( $\xi_y < \xi_x$ ).

---

<sup>5</sup>We take  $q = 0.8, 0.9, 1, 1.1, 1.2, 1.3, 1.4, 1.5, 1.7, 1.9, 2.2, 2.7, 3.7, 5.7, 31.7$  from which we have the tail indices  $\xi = 2.8q \in [2.24; 88.76]$  and the true null hypothesis, respectively, equals to  $H_0 : \theta(F_y) - \theta(F_x) = 0.21, 0.143, 0.1, 0.071, 0.051, 0.035, 0.023, 0.014, 0, -0.01, -0.02, -0.03, -0.041, -0.049, -0.06$  for the Theil and to  $0.12, 0.09, 0.068, 0.052, 0.038, 0.028, 0.019, 0.012, 0, -0.009, -0.018, -0.029, -0.041, -0.052, -0.066$  for the Gini.

Power comparison of the considered permutation and bootstrap methods are valid since rejection probabilities under the null hypothesis  $\theta(F_y) - \theta(F_x) = 0$  are close to the nominal level (here 0.05) which in Figure 10 is represented via the dashed horizontal line.

From Figure 10, we can see that the curves are asymmetric around 0. When  $F_y$  is less heavy-tailed than  $F_x$  ( $\theta(F_y) - \theta(F_x) < 0$ ), the null is quickly rejected as the true null hypothesis moves away from 0. On the other side, when  $F_y$  is more heavy-tailed than  $F_x$  ( $\theta(F_y) - \theta(F_x) > 0$ ), the null is slowly rejected as the true null hypothesis moves away from 0. Overall, we can see that the permutation approach (**perm S rescaled and standard**) is more powerful than the bootstrap under the null approach (**boot S rescaled and standard**), the difference between the two approaches being resampling without replacement rather than with replacement. In addition, the studentized permutation tests based on the combined rescaled samples (**perm S rescaled**) outperforms other methods. It rejects the null much more faster than other methods, especially when  $F_y$  is heavier-tailed than  $F_x$  ( $\theta(F_y) - \theta(F_x) > 0$ ). Figure 11 shows similar results for the Gini index.

## 7 Conclusion

We study Monte-Carlo permutation and bootstrap methods for the problem of testing the equality of inequality measures between two samples. For scale-independent measures, as the Gini, Theil, Generalized Entropy and Atkinson indices, we introduce a convenient rescaling to validate and enhance performance. Our simulation results show that permutation tests control size regardless of tail thickness, when underlying distributions are not too distant (with respect to scale). When underlying distributions differ substantially in their upper tails, proposed permutation methods still provide significant improvement over standard asymptotic and bootstrap tests. In addition, results suggest that rescaling observations by sample means before permutation improves power in finite samples.

## References

- Atkinson, A. B. (1970). On the measurement of inequality. *Journal of Economic Theory* 2, 244–263.
- Boos, D. D. and C. Brownie (2004). Comparing variances and other measures of dispersion. *Statistical Science* 19, 571–578.
- Boos, D. D., P. Janssen, and N. Veraverbeke (1989). Resampling from centered data in the two-sample problem. *Journal of Statistical Planning and Inference* 21, 327–345.
- Brachman, K., A. Stich, and M. Trede (1996). Evaluating parametric income distribution models. *Allgemeines Statistisches Archiv* 80, 285–298.
- Chung, E. and J. P. Romano (2013). Exact and asymptotically robust permutation tests. *Annals of Statistics* 41, 484–507.
- Cowell, F. A. (2011). *Measuring Inequality*. Oxford University Press.
- Cowell, F. A. and E. Flachaire (2007). Income distribution and inequality measurement: The problem of extreme values. *Journal of Econometrics* 141, 1044–1072.
- Cowell, F. A. and E. Flachaire (2014). Statistical methods for distributional analysis. In A. B. Atkinson and F. Bourguignon (Eds.), *Handbook of Income Distribution*, Volume 2, Chapter 7. Elsevier.
- Davidson, R. (2009). Reliable inference for the Gini index. *Journal of Econometrics* 150(1), 30–40.
- Davidson, R. (2012). Statistical inference in the presence of heavy-tails. *Econometrics Journal* 15, C31–C53.
- Davidson, R. and E. Flachaire (2007). Asymptotic and bootstrap inference for inequality and poverty measures. *Journal of Econometrics* 141, 141–166.

- Dufour, J.-M. (2006). Monte Carlo tests with nuisance parameters: A general approach to finite sample inference and nonstandard asymptotics. *Journal of Econometrics* 133, 443–477.
- Dwass, M. (1957). Modified randomization tests for nonparametric hypotheses. *Annals of Mathematical Statistics* 28, 181–187.
- Gill, R. (1988). Non- and semiparametric maximum likelihood estimators and the Von-Mises method (part i). *Scandinavian Journal of Statistics* 16, 97–128.
- Hall, P. (1992). *The Bootstrap and Edgeworth Expansion*. Springer Series in Statistics. New York: Springer Verlag.
- Kleiber, C. and S. Kotz (2003). *Statistical Size Distributions in Economics and Actuarial Sciences*. Wiley Series in Probability and Statistics.
- Lim, T.-S. and W.-Y. Loh (1996). A comparison of tests of equality of variances. *Computational Statistics and Data Analysis* 22, 287–301.
- McDonald, J. B. (1984). Some generalized functions for the size distribution income. *Econometrica* 52, 647–663.
- Pitman, E. J. G. (1937). Significance tests which may be applied to samples from any populations. *Journal of the Royal Statistical Society, Series A* 4, 119–130.
- Romano, J. P. (1990). On the behavior of randomized tests without a group invariance assumption. *Journal of the American Statistical Association* 85, 686–692.
- Schluter, C. and M. Trede (2002). Tails of Lorenz curves. *Journal of Econometrics* 109, 151–166.
- Schluter, C. and K. J. van Garderen (2009). Edgeworth expansions and normalizing transforms for inequality measures. *Journal of Econometrics* 150, 16–29.

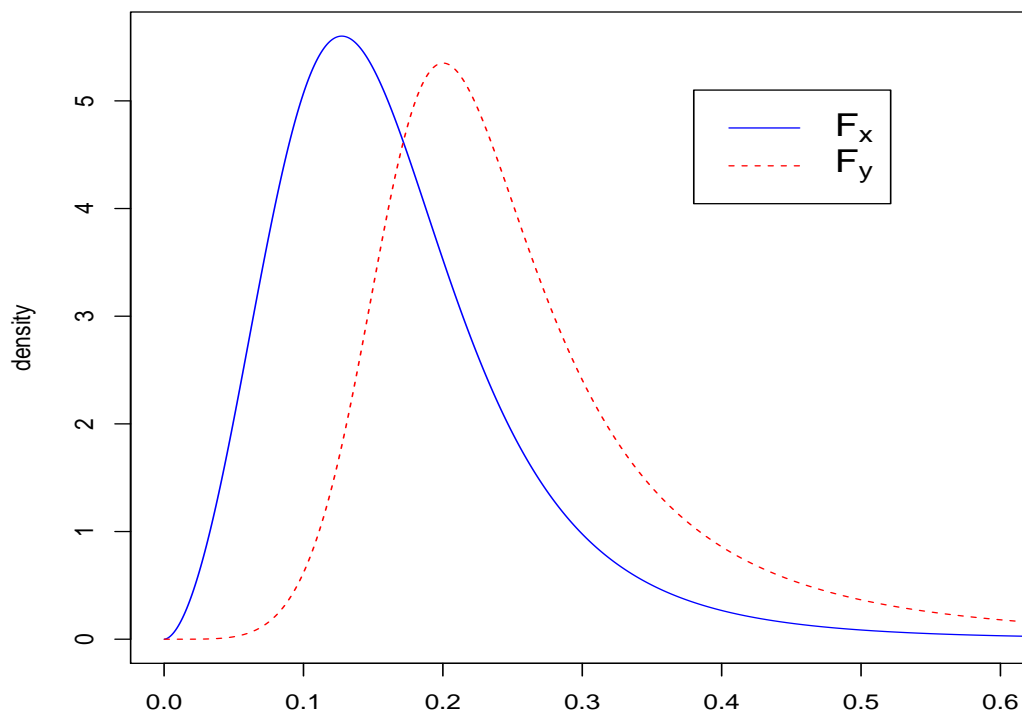


Figure 1: Two very different Singh-Maddala distributions,  $F_x = \text{SM}(x; 2.8, b, 1.7)$  and  $F_y = \text{SM}(y; 5.8, b, 0.5)$  where  $b = 0.1930698$ . The tail parameters are, respectively, equal to  $\xi_x = 2.8 \times 1.7 = 4.76$  and  $\xi_y = 5.8 \times 0.5 = 2.9$ . Thus,  $F_y$  is much more heavy-tailed than  $F_x$ . These two distributions share the same Theil index, equals to 0.1401151.



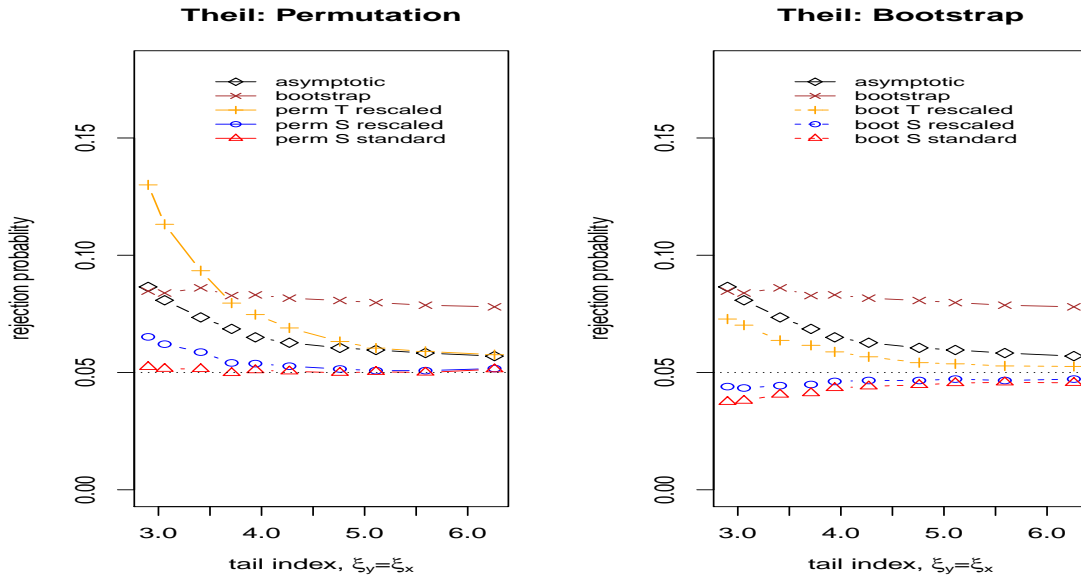


Figure 2: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Theil** inequality measures between two samples. The two distributions are identical,  $\mathbf{F}_x = \mathbf{F}_y$ . The upper tail is heavier as  $\xi_y$  decreases, with  $\xi_y = \xi_x \in [2.9; 6.26]$  and  $n = 50$ .

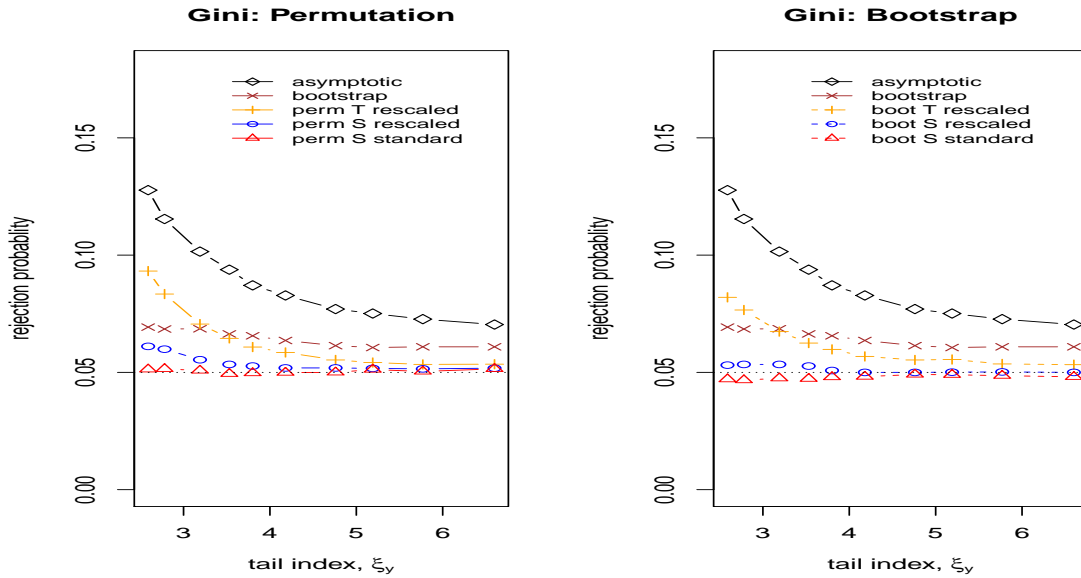


Figure 3: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Gini** inequality measures between two samples. The two distributions are identical,  $\mathbf{F}_x = \mathbf{F}_y$ . The upper tail is heavier as  $\xi_y$  decreases, with  $\xi_y = \xi_x \in [2.59; 6.6]$  and  $n = 50$ .

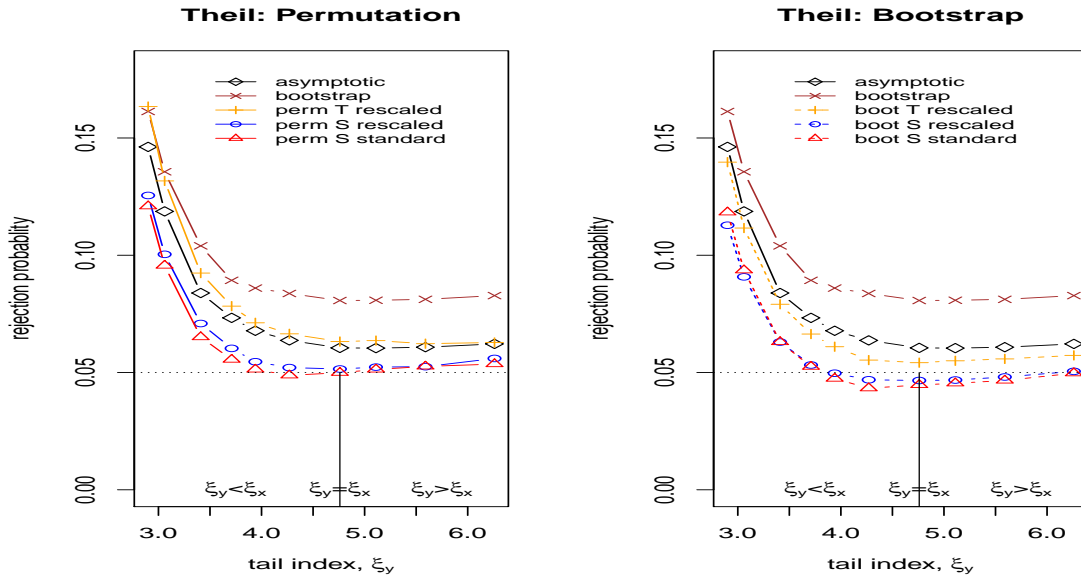


Figure 4: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Theil** inequality measures between two samples. The distribution  $F_x$  is fixed ( $\xi_x = 4.76$ ) and  $\mathbf{F}_x \neq \mathbf{F}_y$ . The distribution  $F_y$  goes away from  $F_x$ , being heavier tailed as  $\xi_y$  decreases, with  $\xi_y \in [2.9; 6.26]$  and  $n = 50$ .

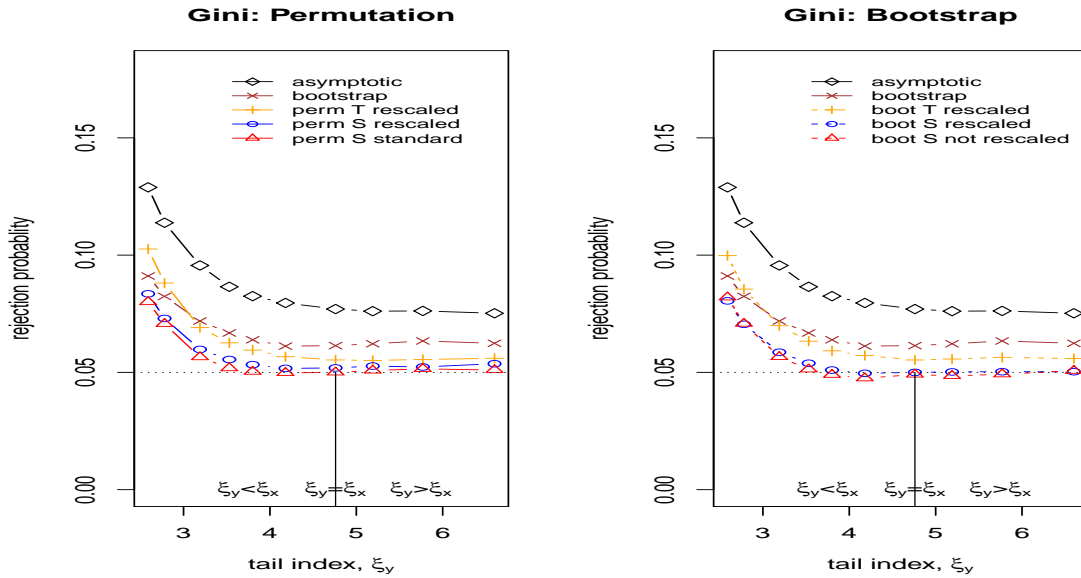


Figure 5: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Gini** inequality measures between two samples. The distribution  $F_x$  is fixed ( $\xi_x = 4.76$ ) and  $\mathbf{F}_x \neq \mathbf{F}_y$ . The distribution  $F_y$  goes away from  $F_x$ , being heavier tailed as  $\xi_y$  decreases, with  $\xi_y \in [2.59; 6.6]$  and  $n = 50$ .

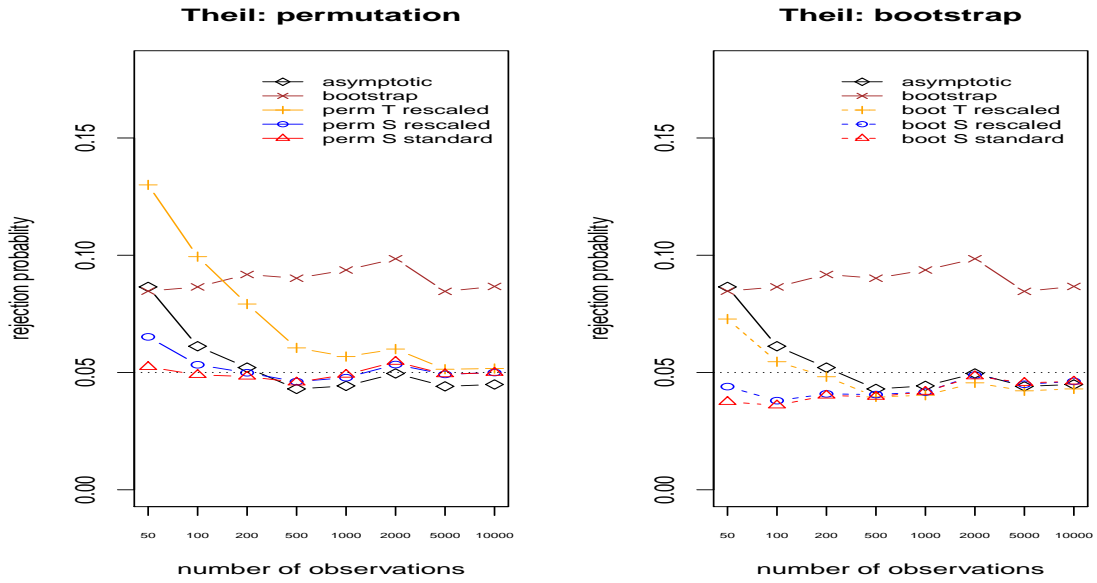


Figure 6: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Theil** inequality measures between two samples, as the **sample size increases**. The two distributions are identical,  $\mathbf{F}_x = \mathbf{F}_y$ , and very heavy-tailed,  $\xi_y = \xi_x = 2.9$ .

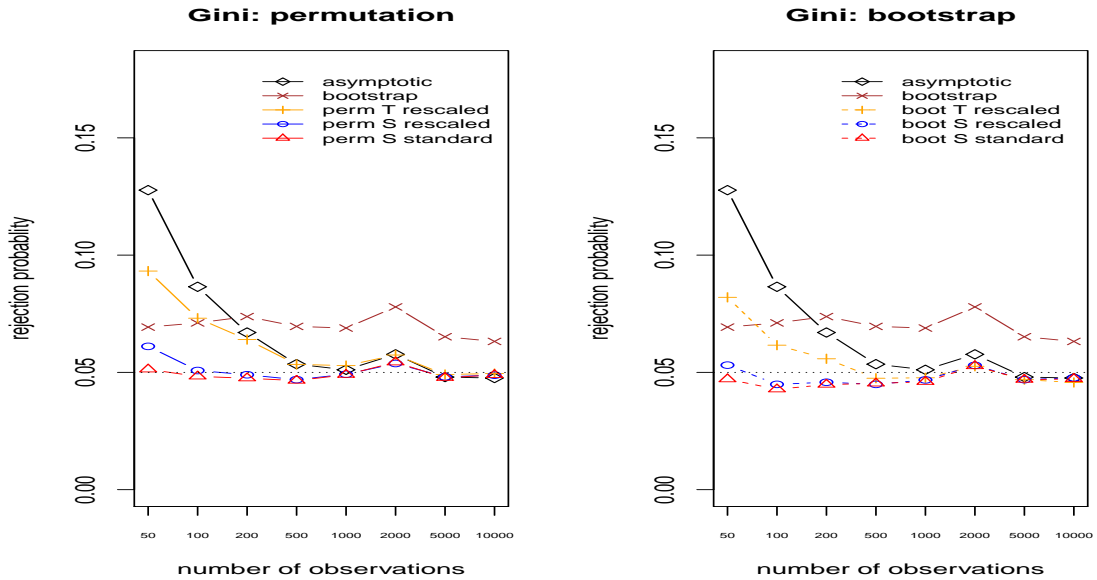


Figure 7: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Gini** inequality measures between two samples, as the **sample size increases**. The two distributions are identical,  $\mathbf{F}_x = \mathbf{F}_y$ , and very heavy-tailed,  $\xi_y = \xi_x = 2.59$ .

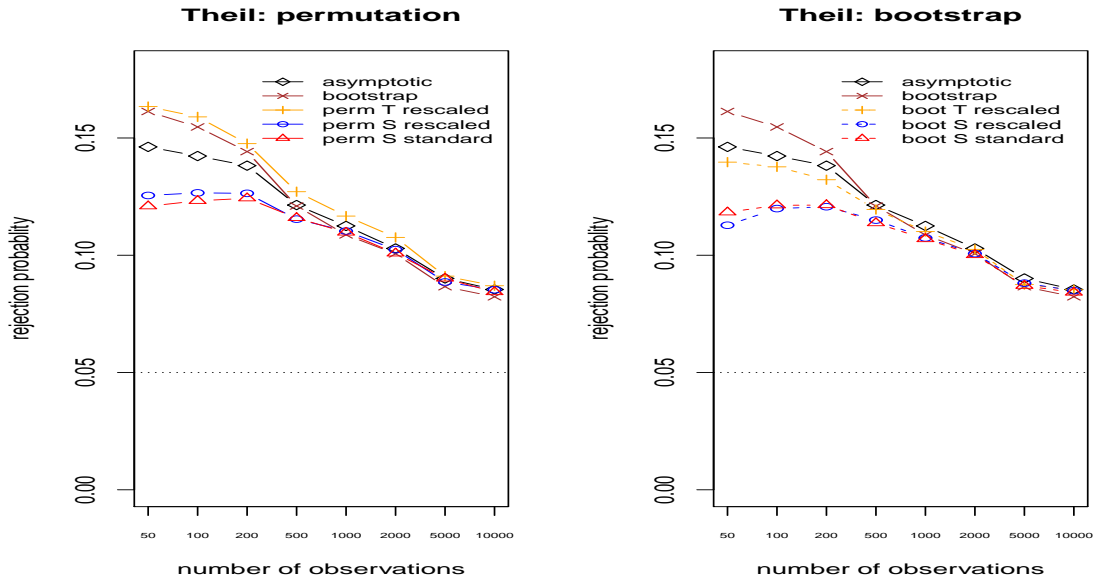


Figure 8: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Theil** inequality measures between two samples, as the **sample size increases**. The two distributions are very different in their upper tails,  $\mathbf{F}_x \neq \mathbf{F}_y$ , with tail parameters equal to  $\xi_x = 4.76$  and  $\xi_y = 2.9$ .

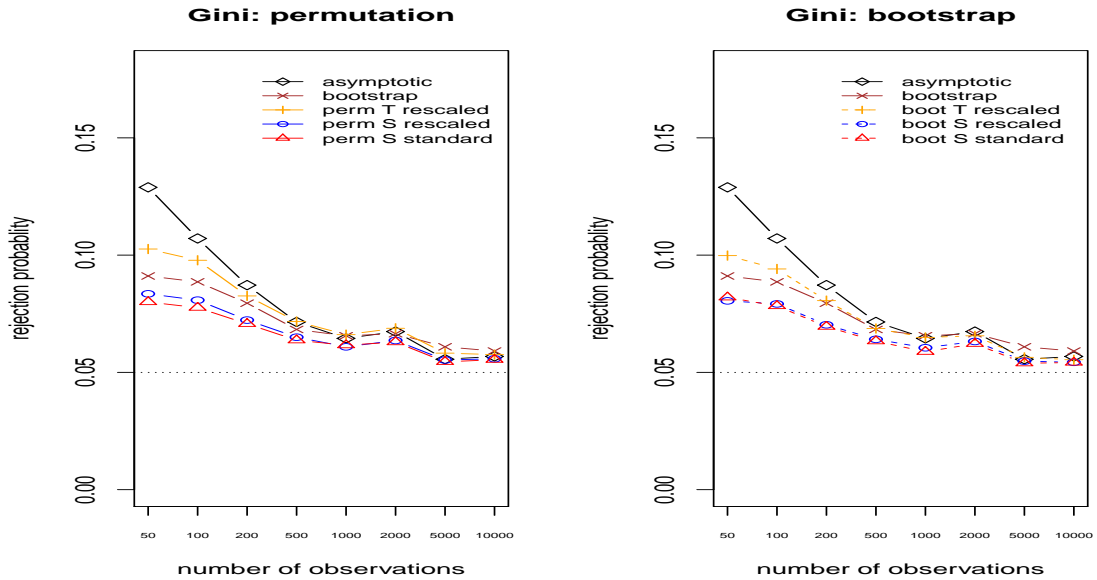


Figure 9: SIZE: Rejection frequencies of asymptotic, permutation and bootstrap tests for the problem of testing the equality of **Gini** inequality measures between two samples, as the **sample size increases**. The two distributions are very different in their upper tails,  $\mathbf{F}_x \neq \mathbf{F}_y$ , with tail parameters equal to  $\xi_x = 4.76$  and  $\xi_y = 2.59$ .

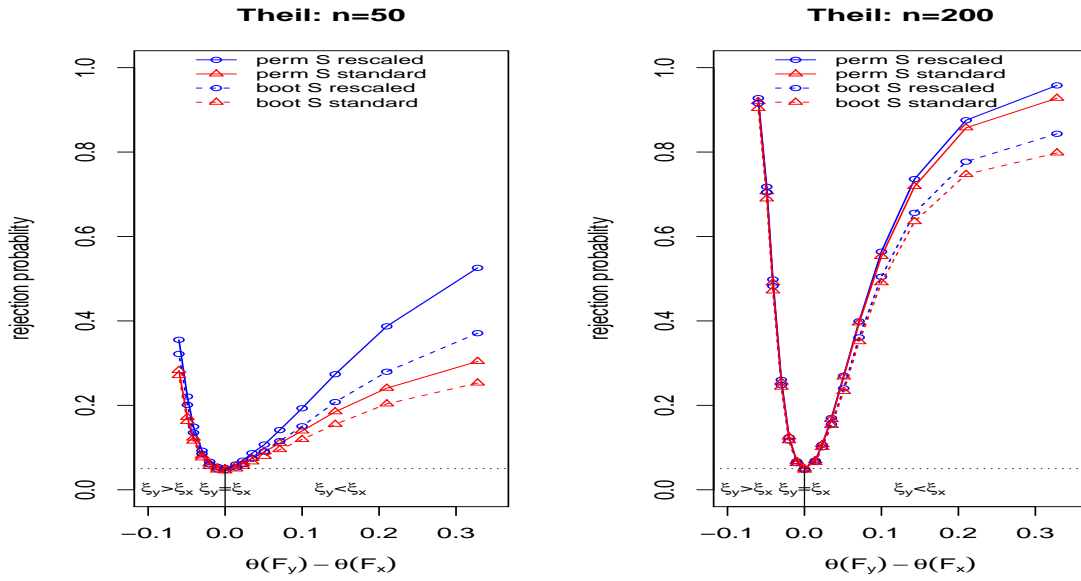


Figure 10: POWER: Rejection frequencies of permutation and bootstrap tests for the problem of testing the equality of **Theil** inequality measures between two samples, when the true null hypothesis is equal to  $\theta(F_y) - \theta(F_x)$ . The distribution  $F_x$  is fixed and the distribution  $F_y$  is heavier tailed as  $\theta(F_y) - \theta(F_x)$  increases.

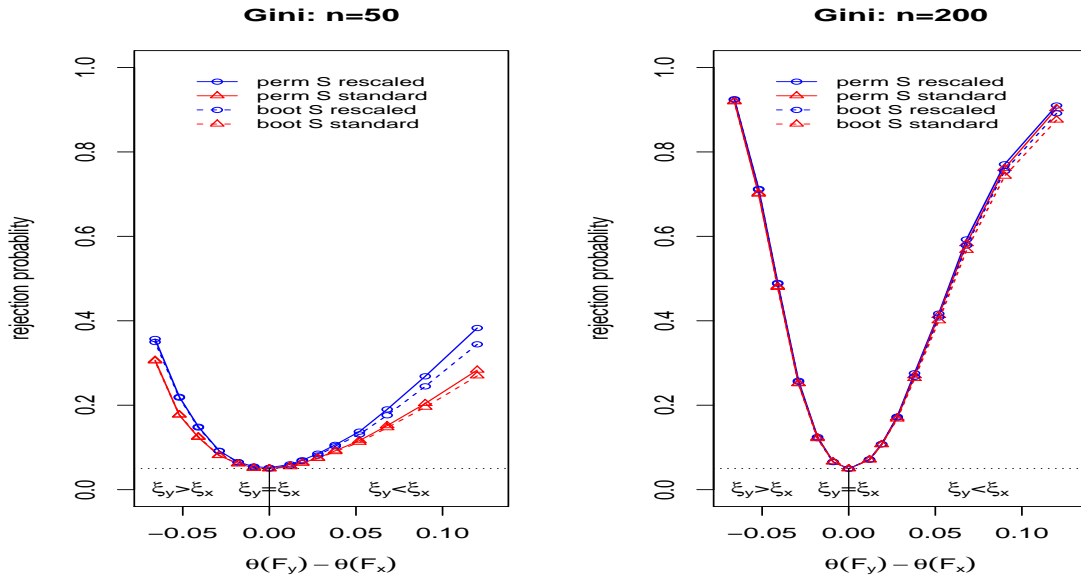


Figure 11: POWER: Rejection frequencies of permutation and bootstrap tests for the problem of testing the equality of **Gini** inequality measures between two samples, when the true null hypothesis is equal to  $\theta(F_x) - \theta(F_y)$ . The distribution  $F_x$  is fixed and the distribution  $F_y$  is heavier tailed as  $\theta(F_y) - \theta(F_x)$  increases.