

Identification-robust methods for comparing inequality with an application to economic convergence *

Jean-Marie Dufour [†]
McGill University

Emmanuel Flachaire [‡]
Aix-Marseille Université

Lynda Khalaf [§]
Carleton University

Abdallah Zalgout [¶]
Carleton University

First version: July 2018

Revised: October 2019, March 2020, June 2020, August 2020, January 2021,
November 2021

This version: November 2021

Compiled: November 16, 2021, 13:44

* This work was supported by the William Dow Chair of Political Economy (McGill University), the Bank of Canada (Research Fellowship), the Toulouse School of Economics (Pierre-de-Fermat Chair of excellence), the Universidad Carlos III de Madrid (Banco Santander de Madrid Chair of excellence), the Natural Sciences and Engineering Research Council of Canada, the Social Sciences and Humanities Research Council of Canada, the Fonds de recherche sur la société et la culture (Québec), and by project ANR-16-CE41-0005 managed by the French National Research Agency (ANR).

[†] 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 414, 855 Sherbrooke Street West, Montréal, Québec H3A 2T7, Canada. TEL: (1) 514 398 6071; FAX: (1) 514 398 4800; email: jeanmarie.dufour@mcgill.ca. Web page: <http://www.jeanmariedulfour.com>

[‡] Aix-Marseille Université, CNRS, EHESS, Centrale Marseille, AMSE. Mailing address: GREQAM-EHESS, Aix-Marseille Université, 2 rue de la charité, 13002, Marseille, France. Email: emmanuel.flachaire@univ-amu.fr

[§] Corresponding author, Economics Department, Carleton University, K1S 5B6, ON, Canada. 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. Email: Lynda.Khalaf@carleton.ca

[¶] Economics Department, Carleton University. abdallahzalgout@cmail.carleton.ca

ABSTRACT

Fieller-type methods are proposed for set inference on the generalized entropy (GE) family of inequality indices. This family satisfies a set of key axiomatic principles and is widely used in practice. We study the general comparison problem of testing differences between indices, with independent or dependent samples. Solutions are analytically tractable and cover tests for any given value of the difference - i.e. not just zero - allowing the construction of confidence sets through test inversion. These sets are robust to the fact that GE indices involve possibly weakly identified parameter ratios. Simulation results illustrate the superiority of our solutions relative to available counterparts, including simulation-based permutation methods. Improvements are especially notable for indices that put more weight on the right tail of the distribution. Proposed methods are applied to study economic convergence across U.S. states and non-OECD countries. We document the fragility of decisions that rely on traditional interpretations of - significant or insignificant - comparisons when the tested differences can be weakly identified. With reference to the growth literature which typically uses the variance of log per-capita income to measure dispersion, results confirm the importance of accounting for micro-founded axioms and shed new light on enduring controversies surrounding convergence.

Keywords: inequality; generalized entropy; two samples; Fieller; identification-robust; economic convergence; per capita income.

Journal of Economic Literature classification: C100, C120, D63, I3, E0, O4.

MSC2020 classification: 62, 62P20, 62P25, 62F03, 62F05, 62F25.

Contents

1	Introduction	1
2	Fieller-type confidence sets for Generalized Entropy inequality measures	4
3	Simulation evidence	9
4	Application: Regional economic convergence	15
5	Conclusion	19
A	Proof of Theorem 2.1	A-1
B	Figures	A-3
B.1	Experiment I; Design (I-a) – Independent samples: $n = m, F_X = F_Y, \Delta_0 = 0$	A-3
B.2	Experiment I; Design (I-b) – Independent samples: $n = m, F_X \neq F_Y, \Delta_0 = 0$	A-4
B.3	Experiment I; Design (I-c) – Independent samples: $n = m, F_X \neq F_Y, \Delta_0 \neq 0$	A-5
B.4	Comparing Fieller’s method and the permutation method	A-6
B.5	Behavior with respect to the sensitivity parameter γ	A-7
B.6	Robustness to the shape of the null distributions	A-7
C	Tables	A-8
C.1	Effect of right tail thickness	A-8
C.2	Effect of left tail thickness	A-9
C.3	Boundedness and width of the confidence intervals	A-10

List of Figures

B.1	Design (I-a) – Size and power of Delta and Fieller-type tests for GE_1 comparisons $H_0: GE_1(X) = GE_1(Y)$, <i>Nominal size</i> = 0.05 [Opt]	A-3
B.2	Design (I-a) – Size and power of Delta and Fieller-type tests for GE_2 comparisons $H_0: GE_2(X) = GE_2(Y)$, <i>Nominal size</i> = 0.05	A-3
B.3	Design (I-b) – Size and power of Delta and Fieller-type tests for GE_1 comparisons.	A-4
B.4	Design (I-b) – Size and power of Delta and Fieller-type tests for GE_2 comparisons	A-4
B.5	Design (I-c) – Size and power of Delta and Fieller-type tests for GE_1 comparisons	A-5
B.6	Design (I-c) – Size and power of Delta and Fieller-type tests for GE_2 comparisons	A-5
B.7	Size and Power of two-sample tests	A-6
B.8	Size and Power of two-sample tests	A-6
B.9	Rejection frequencies of the tests inverted to derive the Delta method and Fieller's confidence sets over the sensitivity parameter γ . <i>Nominal size</i> = 0.05	A-7
B.10	Rejection frequencies of the tests inverted to derive the Delta method and Fieller's confidence sets over the tail index ξ_y . <i>Nominal size</i> = 0.05	A-7

List of Tables

4.1	Estimates and confidence intervals of the change in inequality across U.S. states between 1946 and 2016.	18
4.2	Estimates and confidence intervals of the change in inequality across non-OECD countries	18
C.3	Rejection frequencies of Delta and Fieller methods: effect of right-tail thickness; n=50.	A-8
C.4	Rejection frequencies of Delta and Fieller methods: effect of left-tail thickness; n=50.	A-9
C.5	Rejection probabilities and widths of confidence sets based on the Delta and Fieller-type methods: One-sample problem	A-10
C.6	Rejection probabilities and widths of confidence sets based on the Delta and Fieller-type methods	A-10

1 Introduction

Economic inequality can be broadly defined in terms of the distribution of economic variables, which include income (predominantly), and other variables such as consumption or health. Inequality can be measured in several ways, most of which are justified statistically as well as through theoretical axiomatic approaches. In this context, size-correct statistical inference is an enduring challenge. One reason is that the underlying distributions often have thick tails, which contaminate standard asymptotic and bootstrap-based procedures (Davidson and Flachaire, 2007; Cowell and Flachaire, 2007). Another reason is that two different distributions can yield equal measures, which complicates comparisons (Dufour et al., 2019).

An important additional difficulty is that common inequality measures – such as the generalized entropy (GE) and Gini indices – involve transformations of parameters (Cowell and Flachaire, 2015). Formally, denote by X the random variable with a typical realization representing say the income of a randomly chosen individual in the population, and let F_X refer to the CDF of X . Given a predetermined parameter – denoted γ – that characterizes the sensitivity to changes over different parts of the income distribution, the GE measure – denoted GE_γ – can be defined as a function of the ratio of two particular moments of F_X : the mean $\mu_X = \mathbb{E}_F(X)$ and $v_X(\gamma) = \mathbb{E}_F(X^\gamma)$.¹ Such nonlinear forms may easily be ill-conditioned or poorly identified. The *first* objective of this paper is to underscore and address resulting inference problems. Identification broadly refers to our ability to recover objects of interest from available models and data (Dufour and Hsiao, 2008). Despite a sizeable literature on inequality, methods that take into account the irregularities underscored in the weak identification literature appear to be missing in this context.²

More to the point from the index comparison perspective, most available approaches for this purpose focus on *significance* tests. The *second* objective of this paper is to document the fragility of decisions relying on traditional interpretations of – significant or insignificant – test results, when the difference under test can be weakly identified. In particular, when a zero difference cannot be rejected, we show that because of the definition of conventional inequality indices, one may also not be able to refute a large spectrum of possible values of this difference. From a policy perspective, this indicates that available samples are uninformative on inequality changes, which stands in sharp contrast to a no-change conclusion.

The *third* objective is to propose tractable identification-robust confidence sets for inequality indices – in particular, for differences between such indices – which require the same basic inputs

¹This definition implies that GE_γ is more sensitive to differences in the top (bottom) tail with more positive (negative) γ .

²See *e.g.* Dufour (1997), Andrews and Cheng (2013), Kleibergen (2005), Andrews and Mikusheva (2015), Beaulieu et al. (2013), Bertanha and Moreira (2016), and references therein; see also Bahadur and Savage (1956) and Gleser and Hwang (1987).

as their standard counterparts. Whereas usual companion variances and covariances as well as critical values need to be computed, the alternative test statistics are formed and inverted analytically into confidence sets that will reflect the underlying identification status.

The *fourth* objective is to discuss challenges for empirical researchers and policy-makers in light of the above observations. We study evidence on economic convergence; see *e.g.* Romer (1994) for a historical critical perspective. We show that conflict in test decisions and uninformative confidence sets cannot be ruled out with standard measures and data sets. To the best of our knowledge, this problem and our proposed solution have escaped formal notice.

Indeed, the literature on statistical inference for inequality measures is relatively recent; see Cowell and Flachaire (2015) for a comprehensive survey. In particular, the standard bootstrap is known to fail, and alternative methods remain scarce. For testing the equality of two inequality measures from independent samples, Dufour et al. (2019) suggest a permutational approach for the two-sample problem which outperforms other asymptotic and bootstrap methods available in the literature. However, these results are limited to testing the equality of two inequality measures and do not provide a way of making inference on a *possibly non-zero* difference nor building a confidence interval on the difference.

In the present paper, we propose Fieller-type methods for set inference on the GE family of inequality indices. This family satisfies a set of key axiomatic principles and is widely used in practice.³ We study the general comparison problem of testing any *possibly non-zero* difference between measures, with either independent or dependent samples. Moving from testing a zero difference to assessing the size of the difference is much more informative from both statistical and economic viewpoints, including potential policy recommendations.

The fact that inequality measures in general, and those considered in this paper in particular, can be expressed as ratios of moments or ratios of functions of moments, provides a strong motivation for our work since Fieller-type methods are typically used for inference on ratios. Fieller's original solution for the means of two independent normal random variables was extended to multivariate normals (Bennett, 1959), general exponential (Cox, 1967) and linear (Zerbe, 1978; Dufour, 1997) regression models, dynamic models with possibly persistent covariates (Bernard et al., 2007, 2019) and for simultaneous inference on multiple ratios (Bolduc et al., 2010). For a good review of inference on ratios, see Franz (2007).

On the GE class of inequality indices, this paper makes the following contributions. *First*, we provide analytical and tractable solutions for proposed confidence sets. *Second*, we show in a simulation study that the proposed solutions are more reliable than Delta counterparts. *Third*, we show that our approach outperforms most simulation-based alternatives including the permutation

³These include scale invariance, the Pigou-Dalton transfer, the symmetry and the Dalton population principle. It is also additively decomposable. See Cowell (2000) for a detailed discussion on these and other properties of indices.

test of Dufour et al. (2019). *Fourth*, our solution covers tests for any given value of the difference [i.e. not just zero, in contrast with Dufour et al. (2019)], allowing the construction of confidence sets through test inversion. *Fifth*, we provide useful empirical evidence supporting the seemingly counter-intuitive bounds that Fieller-type methods can produce.

Key simulation results illustrate the superiority of Fieller-type methods across the board: (1) the improved level control (over the Delta method) is especially notable for indices that put more weight on the right tail of the distribution *i.e.* as γ increases; (2) size improvements preserve power; (3) results are robust to different assumptions on the shape of the null distributions; (4) tests based on the Fieller-type method outperform available permutation tests when the distributions under the null hypothesis are different. A permutational approach is not available (to date) for the general problem we consider here. Overall, while irregularities arising from the right tail have long been documented, we find that left-tail irregularities are equally important in explaining the failure of standard inference methods for inequality measures.

Our empirical study on growth demonstrates the practical relevance of these theoretical results. Using per-capita income data for 48 U.S. states, we analyze the convergence hypothesis by comparing the inequality levels between 1946 and 2016. In contrast to the bulk of this literature, we depart from just testing and build confidence sets to document the economic and policy significance of statistical decisions. The empirical literature on growth relies on the variance of log incomes as a measure of dispersion in per-capita income distributions (Blundell et al., 2008). But this measure violates the Pigou-Dalton principle (Araar and Duclos, 2006). We use GE indices instead, since these satisfy the axioms suggested in the inequality measurement literature. We document specific cases where the variance of log incomes decrease while the GE_2 measure indicates the opposite.

We find that inter-state inequality has declined over the 1946-2016 period indicating convergence across the states. For the GE_2 index, the Fieller-type and Delta methods lead to contradictory conclusions: in contrast to the former, the latter suggests that inequality declines are insignificant at usual levels. Results with non-OECD countries stress the severe consequences of ignoring identification problems: with the GE_2 index, the Fieller-type method produces an unbounded set, which casts serious doubts on the reliability of the no-change results using the Delta method.

The rest of the paper is organized as follows. Section 2 derives Fieller-type confidence sets. Section 3 reports the results of the simulation study. Section 4 contains the inter-state convergence application, and Section 5 concludes.

2 Fieller-type confidence sets for Generalized Entropy inequality measures

An inequality measure is a measure of dispersion for the distribution of a random variable. Throughout this paper, it will be convenient to focus on income distributions, though our results also apply to other variables relevant to inequality studies, such as wage, health, and consumption distributions. Many inequality measures, including the GE_γ class, solely depend on the underlying distribution and can typically be written as a functional which maps the space of the cumulative distribution function (**CDF**) to the nonnegative real line \mathbb{R}_+^0 .

Our aim is to make inference on the GE_γ measure for any given $\gamma \in (0, 2)$. In particular, we wish to build an asymptotic Fieller-type confidence set (**FCS**) for the difference between two measures. We call this problem the *two-sample problem*, as opposed to the *one-sample problem* where the objective consists in testing and building a confidence interval for a single index. The crucial difference between a FCS and its Wald-type counterpart based on using an approximate standard error derived by the Delta method (**DCS**) is that FCS start by reformulating the null hypothesis in a linear form. One then proceeds by inverting the square of the t-test associated with the reformulated linear hypothesis. This avoids the irregularities which affect the validity of the Delta method (*e.g.*, as the denominator approaches zero).

A consequence of rewriting the null hypothesis in linear form is that the variance used by the Fieller-type statistic depends on the true value of the tested parameter. This leads to a quadratic inequality problem. The resulting confidence regions are not standard, in the sense that they may be asymmetric, consisting of two disjoint unbounded confidence intervals or the whole real line \mathbb{R} . Nevertheless, unbounded intervals are an attractive feature of the method which addresses coverage problems (Koschat et al., 1987; Gleser and Hwang, 1987; Dufour, 1997; Dufour and Jasiak, 2001; Dufour and Taamouti, 2005, 2007; Bertanha and Moreira, 2016). For a geometric comparison of the Fieller and Delta methods, see Hirschberg and Lye (2010).

Let X be a positive random variable such that both moments $\mathbb{E}_F(X)$ and $\mathbb{E}_F(X^\gamma)$ are finite, *i.e.*

$$\mathbb{P}[X > 0] > 0, \quad 0 < \mu_X := \mathbb{E}_F(X) < \infty, \quad 0 < \nu_X(\gamma) := \mathbb{E}_F(X^\gamma) < \infty. \quad (2.1)$$

Then the $GE_\gamma(X)$ measure can be expressed as in Shorrocks (1980):

$$\begin{aligned} GE_\gamma(X) &= \frac{1}{\gamma(\gamma-1)} \left[\frac{\mathbb{E}_F(X^\gamma)}{[\mathbb{E}_F(X)]^\gamma} - 1 \right] \quad \text{for } \gamma \neq 0, 1, \\ GE_0(X) &= \mathbb{E}_F[\log(X)] - \log[\mathbb{E}_F(X)] \\ GE_1(X) &= \frac{\mathbb{E}_F[X \log(X)]}{\mathbb{E}_F(X)} - \log[\mathbb{E}_F(X)]. \end{aligned} \quad (2.2)$$

This class of measures includes several common indices, including two well-known ones introduced by Theil (1967): the Mean Logarithmic Deviation (**MLD**), which is the limiting value of the

$GE_\gamma(X)$ as γ approaches zero, and the Theil index, which is the limiting value of the $GE_\gamma(X)$ as γ approaches 1. When $\gamma = 2$, the index is equal to half the squared coefficient of variation and is related to the Hirschman-Herfindahl (HH) index, used in industrial organization (Schluter, 2012). The Atkinson index can be obtained from the $GE_\gamma(X)$ index using an appropriate transformation.

Denote by X the random variable representing incomes of individuals from the first population with CDF F_X , and by Y the incomes of individuals from the second population with CDF F_Y , both satisfying (2.1). We assume we have i.i.d. samples X_1, \dots, X_n and Y_1, \dots, Y_m from each population. The empirical distribution functions (**EDFs**) associated with these samples are:

$$\hat{F}_X(x) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}(X_i \leq x), \quad \hat{F}_Y(y) = \frac{1}{m} \sum_{j=1}^m \mathbf{1}(Y_j \leq y), \quad (2.3)$$

where $\mathbf{1}(\cdot)$ is the indicator function that takes the value 1 if the argument is true, and 0 otherwise. The following presentation sets $\gamma \neq 1, 0$. The Theil index can be treated along the same lines, beginning from the expressions in (2.2). The MLD measure eschews the inference problem associated with ratios; for further insights, see Cowell and Flachaire (2018) and the references therein. Under standard laws of large numbers, we can consistently estimate the index $GE_\gamma(X)$ and $GE_\gamma(Y)$ by

$$\widehat{GE}_\gamma(X) := \frac{1}{\gamma(\gamma-1)} \left[\frac{\hat{v}_X(\gamma)}{\hat{\mu}_X^\gamma} - 1 \right], \quad \widehat{GE}_\gamma(Y) := \frac{1}{\gamma(\gamma-1)} \left[\frac{\hat{v}_Y(\gamma)}{\hat{\mu}_Y^\gamma} - 1 \right], \quad (2.4)$$

where

$$\hat{\mu}_X := \int x d\hat{F}_X = \frac{1}{n} \sum_{i=1}^n X_i, \quad \hat{v}_X(\gamma) := \int x^\gamma d\hat{F}_X = \frac{1}{n} \sum_{i=1}^n X_i^\gamma, \quad (2.5)$$

$$\hat{\mu}_Y := \int y d\hat{F}_Y = \frac{1}{m} \sum_{j=1}^m Y_j, \quad \hat{v}_Y(\gamma) := \int y^\gamma d\hat{F}_Y = \frac{1}{m} \sum_{j=1}^m Y_j^\gamma. \quad (2.6)$$

We will now propose level $1 - \alpha$ confidence sets for the difference under consideration, which can be written as a ratio of the four moments μ_x , v_x , μ_y , and v_y :

$$\Delta GE_\gamma := GE_\gamma(X) - GE_\gamma(Y) = \frac{v_X(\gamma)\mu_Y^\gamma - v_Y(\gamma)\mu_X^\gamma}{\gamma(\gamma-1)\mu_Y^\gamma\mu_X^\gamma}. \quad (2.7)$$

ΔGE_γ can be estimated by substituting estimates of the relevant moments [see (2.5) - (2.6)]:

$$\Delta \widehat{GE}_\gamma := \widehat{GE}_\gamma(X) - \widehat{GE}_\gamma(Y) = \frac{\hat{v}_X(\gamma)\hat{\mu}_Y^\gamma - \hat{v}_Y(\gamma)\hat{\mu}_X^\gamma}{\gamma(\gamma-1)\hat{\mu}_Y^\gamma\hat{\mu}_X^\gamma}. \quad (2.8)$$

Our analysis covers three cases defined by the following assumptions.

Assumption 2.1. *Samples are of equal sizes and independent.*

Assumption 2.2. *Samples are of unequal sizes and independent.*

Assumption 2.3. *Samples are of equal sizes and dependent.*

Denote by λ the vector of the moments μ_X , v_X , μ_Y , and v_Y :

$$\lambda = (\lambda_1, \lambda_2, \lambda_3, \lambda_4)' = (\mu_X, v_X, \mu_Y, v_Y)', \quad (2.9)$$

and by $\hat{\lambda}$ the vector of the estimates of the moments in λ defined in (2.5) - (2.6). Furthermore, on setting

$$\mathbf{T} := \tau \otimes I_2, \quad \tau := \begin{bmatrix} n & 0 \\ 0 & m \end{bmatrix} \quad (2.10)$$

where I_2 is the 2×2 identity matrix. Under the standard assumption that the estimated moments defined in (2.9) are asymptotically normal, we can write:

$$\mathbf{T}^{-1/2}(\hat{\lambda} - \lambda) \xrightarrow{D} N(\mathbf{0}, \Sigma), \quad \Sigma := [\sigma_{ij}]_{i,j=1,\dots,4}. \quad (2.11)$$

The standard DCS is obtained by inverting the square (or the absolute value) of the t-test associated with

$$H_D(\Delta_0) : \Delta GE_\gamma = \Delta_0 \quad (2.12)$$

where Δ_0 is any known admissible value of ΔGE_γ , including possibly $\Delta_0 = 0$, for equality. We derive the DCS and FCS for each of Assumptions 2.1 - 2.3. These cases will actually differ only by the expression of the variance. Thus to avoid redundancy, we will derive the method in its most general form and state the restrictions required to obtain the relevant formulae in each case.

By inverting a test statistic with respect to the parameter tested (Δ_0 in this case), we mean collecting the values of the parameter for which the test cannot be rejected at a given significance level α . Assuming that the estimator is asymptotically normal, this can be carried out by solving the following inequality for Δ_0 :

$$(\Delta \widehat{GE}_\gamma - \Delta_0)^2 \leq z_{\alpha/2}^2 \widehat{V}[\Delta \widehat{GE}_\gamma] \quad (2.13)$$

where $\Delta \widehat{GE}_\gamma = \widehat{GE}_\gamma(X) - \widehat{GE}_\gamma(Y)$ and $z_{\alpha/2}$ is the asymptotic two-tailed critical value at the significance level α (i.e., $\mathbb{P}[Z \geq z_{\alpha/2}] = \alpha/2$ for $Z \sim N[0, 1]$) and is the estimate of the asymptotic variance. The solution of (2.13) yields the Delta-method confidence set:

$$\text{DCS}(\Delta GE_\gamma; 1 - \alpha) = \left[\Delta \widehat{GE}_\gamma \pm z_{\alpha/2} [\widehat{V}(\Delta \widehat{GE}_\gamma)]^{1/2} \right]. \quad (2.14)$$

The formula for the asymptotic variance $V(\Delta \widehat{GE}_\gamma)$ in (2.14) depends on the assumptions made

on the observations. For the assumptions 2.1 - 2.3, we get by using the Delta method:

$$\text{under Ass. 2.1 : } V(\Delta\widehat{GE}_\gamma) = \frac{1}{n} \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial \Delta GE_\gamma}{\partial \lambda_i} \frac{\partial \Delta GE_\gamma}{\partial \lambda_j} \sigma_{ij} + \frac{1}{n} \sum_{i=3}^4 \sum_{j=3}^4 \frac{\partial \Delta GE_\gamma}{\partial \lambda_i} \frac{\partial \Delta GE_\gamma}{\partial \lambda_j} \sigma_{ij}, \quad (2.15)$$

$$\text{under Ass. 2.2 : } V(\Delta\widehat{GE}_\gamma) = \frac{1}{n} \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial \Delta GE_\gamma}{\partial \lambda_i} \frac{\partial \Delta GE_\gamma}{\partial \lambda_j} \sigma_{ij} + \frac{1}{m} \sum_{i=3}^4 \sum_{j=3}^4 \frac{\partial \Delta GE_\gamma}{\partial \lambda_i} \frac{\partial \Delta GE_\gamma}{\partial \lambda_j} \sigma_{ij}, \quad (2.16)$$

$$\text{under Ass. 2.3 : } V(\Delta\widehat{GE}_\gamma) = \frac{1}{n} \sum_{i=1}^4 \sum_{j=1}^4 \frac{\partial \Delta GE_\gamma}{\partial \lambda_i} \frac{\partial \Delta GE_\gamma}{\partial \lambda_j} \sigma_{ij}. \quad (2.17)$$

$\widehat{V}(\Delta\widehat{GE}_\gamma)$ may then be computed by replacing σ_{ij} with $\hat{\sigma}_{ij}$, and λ with $\hat{\lambda}$.

In contrast, we propose a Fieller-type set by inverting the square of the t-test associated with a linearized counterpart of $H_D(\Delta_0)$. We thus first reformulate the null hypothesis in a linear form (without the ratio transformation). This can be obtained through the multiplication of both sides of (2.12) by the common denominator $\gamma(\gamma-1)\mu_X^\gamma\mu_Y^\gamma$. For presentation clarity, we denote by θ_1 and θ_2 the numerator and the denominator in (2.7):

$$\theta_1 = v_X(\gamma)\mu_Y^\gamma - v_Y(\gamma)\mu_X^\gamma, \quad \theta_2 = \gamma(\gamma-1)\mu_Y^\gamma\mu_X^\gamma. \quad (2.18)$$

The linear form of the null hypothesis is the following:

$$H_F(\Delta_0) : \Theta(\Delta_0) = 0 \quad \text{where} \quad \Theta(\Delta_0) := \theta_1 - \theta_2\Delta_0. \quad (2.19)$$

We then consider the acceptance region associated with the t-test of this linear hypothesis:

$$\widehat{\Theta}(\Delta_0)^2 \leq z_{\alpha/2}^2 \widehat{V}[\widehat{\Theta}(\Delta_0)] \quad (2.20)$$

where we use the moment-type estimators based on (2.5) - (2.6), *i.e.*

$$\widehat{\Theta}(\Delta_0) := \hat{\theta}_1 - \hat{\theta}_2\Delta_0, \quad \hat{\theta}_1 := \hat{v}_X(\gamma)\hat{\mu}_Y^\gamma - \hat{v}_Y(\gamma)\hat{\mu}_X^\gamma, \quad \hat{\theta}_2 := \gamma(\gamma-1)\hat{\mu}_Y^\gamma\hat{\mu}_X^\gamma, \quad (2.21)$$

and $\widehat{V}[\widehat{\Theta}(\Delta_0)]$ is a consistent estimator of $V[\widehat{\Theta}(\Delta_0)]$, the asymptotic variance of $\widehat{\Theta}(\Delta_0)$ under $H_F(\Delta_0)$. Note the latter consistency needs to hold only under the null hypothesis $H_F(\Delta_0)$. On using the asymptotic normality assumption (2.11), the acceptance region (2.20) yields a confidence for ΔGE_γ with level $1 - \alpha$ (asymptotically):

$$\text{FCS}[\Delta GE_\gamma; 1 - \alpha] = \{\Delta_0 : \widehat{\Theta}(\Delta_0)^2 \leq z_{\alpha/2}^2 \widehat{V}[\widehat{\Theta}(\Delta_0)]\} \quad (2.22)$$

We call $\text{FCS}[\Delta GE_\gamma; 1 - \alpha]$ the level- $(1 - \alpha)$ *Fieller-type confidence set* for ΔGE_γ . Estimating $\mathbb{V}[\hat{\Theta}(\Delta_0)]$ will require estimating the asymptotic covariance of $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)'$. For future reference, we denote the latter and the corresponding estimator as follows:

$$\mathbb{V}(\hat{\theta}) = \begin{bmatrix} \mathbb{V}(\hat{\theta}_1) & \mathbb{C}(\hat{\theta}_1, \hat{\theta}_2) \\ \mathbb{C}(\hat{\theta}_1, \hat{\theta}_2) & \mathbb{V}(\hat{\theta}_2) \end{bmatrix}, \quad \hat{\mathbb{V}}(\hat{\theta}) = \begin{bmatrix} \hat{\mathbb{V}}(\hat{\theta}_1) & \hat{\mathbb{C}}(\hat{\theta}_1, \hat{\theta}_2) \\ \hat{\mathbb{C}}(\hat{\theta}_1, \hat{\theta}_2) & \hat{\mathbb{V}}(\hat{\theta}_2) \end{bmatrix}. \quad (2.23)$$

The form of the Fieller-type confidence may not be clear from (2.22). The following theorem characterizes $\text{FCS}[\Delta GE_\gamma; 1 - \alpha]$ in an explicit way.

Theorem 2.1. *Let $\hat{\mathbb{V}}(\hat{\theta})$ be an estimate of $\mathbb{V}(\hat{\theta})$ in (2.23). Then the confidence set $\text{FCS}[\Delta GE_\gamma; 1 - \alpha]$ defined in (2.22) can be computed as follows:*

$$\begin{aligned} \text{FCS}[\Delta GE_\gamma; 1 - \alpha] &= \{\Delta_0 : A\Delta_0^2 + B\Delta_0 + C \leq 0\} \\ &= \begin{cases} \left[\frac{-B-\sqrt{D}}{2A}, \frac{-B+\sqrt{D}}{2A} \right] & \text{if } D \geq 0 \text{ and } A > 0 \\ \left[-\infty, \frac{-B+\sqrt{D}}{2A} \right] \cup \left[\frac{-B-\sqrt{D}}{2A}, +\infty \right] & \text{if } D \geq 0 \text{ and } A < 0 \\ \left[-\infty, -\frac{C}{B} \right] & \text{if } A = 0 \text{ and } B > 0 \\ \left[-\frac{C}{B}, \infty \right] & \text{if } A = 0 \text{ and } B < 0 \\ \mathbb{R} & \text{if } [A = B = 0 \text{ and } C \leq 0] \text{ or } [D < 0 \text{ and } A \leq 0] \\ \emptyset & \text{if } [A = B = 0 \text{ and } C > 0] \text{ or } [D < 0 \text{ and } A > 0] \end{cases} \end{aligned} \quad (2.24)$$

where

$$A := \hat{\theta}_2^2 - z_{\alpha/2}^2 \hat{\mathbb{V}}(\hat{\theta}_2), \quad B := -2[\hat{\theta}_1 \hat{\theta}_2 - z_{\alpha/2}^2 \hat{\mathbb{C}}(\hat{\theta}_1, \hat{\theta}_2)], \quad C := \hat{\theta}_1^2 - z_{\alpha/2}^2 \hat{\mathbb{V}}(\hat{\theta}_1), \quad (2.25)$$

$$\begin{aligned} D := B^2 - 4AC &= 4z_{\alpha/2}^2 \{ [\hat{\theta}_1^2 \hat{\mathbb{V}}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{\mathbb{V}}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{\mathbb{C}}(\hat{\theta}_1, \hat{\theta}_2)] \\ &\quad + z_{\alpha/2}^2 [\hat{\mathbb{C}}(\hat{\theta}_1, \hat{\theta}_2)^2 - \hat{\mathbb{V}}(\hat{\theta}_1) \hat{\mathbb{V}}(\hat{\theta}_2)] \}. \end{aligned} \quad (2.26)$$

If furthermore $\hat{\mathbb{V}}(\hat{\theta})$ is positive definite, then

$$D < 0 \implies [A < 0 \text{ and } C < 0]. \quad (2.27)$$

The proof is available in Appendix A. Unlike the Delta method, the Fieller-type method satisfies the theoretical result which states that, for a confidence interval of a locally almost unidentified (LAU) parameter, or a parametric function, to attain correct coverage, it should allow for a non-zero probability of being unbounded (Koschat et al., 1987; Gleser and Hwang, 1987; Dufour, 1997;

Dufour and Taamouti, 2005, 2007; Bertanha and Moreira, 2016).

Theorem 2.1 allows for non-positive definite matrix $\hat{V}(\hat{\theta})$ [at least, for the specific sample considered]. When $\hat{V}(\hat{\theta})$ is positive definite, (2.27) implies that $\text{FCS}[\Delta GE_\gamma; 1 - \alpha]$ may be empty only when $A = B = 0$ and $C > 0$, *i.e.* $A\Delta_0^2 + B\Delta_0 + C = C > 0$ [an event with zero probability when $(\hat{\theta}_1, \hat{\theta}_2)'$ has a Gaussian distribution]. Note that the condition $A > 0$ means that θ_2 is significantly different from zero [according to the criterion $\hat{\theta}_2^2/\hat{V}(\hat{\theta}_2) > z_{\alpha/2}^2$], while $C > 0$ means that θ_1 is significantly different from zero [according to the criterion $\hat{\theta}_1^2/\hat{V}(\hat{\theta}_1) > z_{\alpha/2}^2$].

Consistent estimation of these depends on the assumptions made on the observations $[X_1, \dots, X_n$ and $Y_1, \dots, Y_m]$. For the assumptions 2.1, 2.2 and 2.3, we get (using the delta method):

$$\text{under Ass. 2.1: } \quad \mathbf{V}(\hat{\theta}_1) = \frac{1}{n}S_{11}, \quad \mathbf{V}(\hat{\theta}_2) = \frac{1}{n}S_{22}, \quad \mathbf{C}(\hat{\theta}_1, \hat{\theta}_2) = \frac{1}{n}(S_{12} + S_{21}), \quad (2.28)$$

$$\text{under Ass. 2.2: } \quad \mathbf{V}(\hat{\theta}_1) = \frac{1}{n}S_{11}, \quad \mathbf{V}(\hat{\theta}_2) = \frac{1}{m}S_{22}, \quad \mathbf{C}(\hat{\theta}_1, \hat{\theta}_2) = \frac{1}{n}S_{12} + \frac{1}{m}S_{21}, \quad (2.29)$$

$$\text{under Ass. 2.3: } \quad \mathbf{C}(\hat{\theta}_k, \hat{\theta}_l) = \frac{1}{n} \sum_{i=1}^4 \sum_{j=1}^4 \frac{\partial \theta_k}{\partial \lambda_i} \frac{\partial \theta_l}{\partial \lambda_j} \sigma_{ij} \quad \text{for } k = 1, 2, l = 1, 2, \quad (2.30)$$

where

$$S_{11} := \sum_{i=1}^4 \sum_{j=1}^4 \frac{\partial \theta_1}{\partial \lambda_i} \frac{\partial \theta_1}{\partial \lambda_j} \sigma_{ij}, \quad S_{22} := \sum_{i=1}^4 \sum_{j=1}^4 \frac{\partial \theta_2}{\partial \lambda_i} \frac{\partial \theta_2}{\partial \lambda_j} \sigma_{ij}, \quad (2.31)$$

$$S_{12} := \sum_{i=1}^2 \sum_{j=1}^2 \frac{\partial \theta_1}{\partial \lambda_i} \frac{\partial \theta_2}{\partial \lambda_j} \sigma_{ij}, \quad S_{21} := \sum_{i=3}^4 \sum_{j=3}^4 \frac{\partial \theta_1}{\partial \lambda_i} \frac{\partial \theta_2}{\partial \lambda_j} \sigma_{ij}. \quad (2.32)$$

The above presumes asymptotic normality of the underlying criteria. In fact, the considered measures are known transformations of two moments the estimators of which are asymptotically normal under standard regularity assumptions; see Davidson and Flachaire (2007) and Cowell and Flachaire (2007). These typically require that the first two moments exists and are finite. Asymptotic normality of the statistics in (2.13) and (2.20) thus follows straightforwardly. Nevertheless, convergence in this context is known to be slow, especially when the distribution of the data is heavy-tailed and with indices that are sensitive to the upper tail. Our simulations confirm these issues, yet the Fieller-based criteria perform better than the Delta method in finite samples because these eschew problems arising from the ratio.

3 Simulation evidence

This section describes a simulation study designed to compare the finite-sample properties of FCS to the standard DCS. This will be done for the two popular inequality measures nested in the

general entropy class of inequality measures: the Theil Index (GE_1), and half of the coefficient of variation squared (GE_2) which is related to the Hirschman-Herfindahl (HH) index. The tables and figures are in the appendix.

We report the *rejection frequencies* of the tests underlying the proposed confidence sets, under both the null hypothesis (level control) and the alternative (power). Under the null hypothesis, these can also be interpreted as 1 minus the corresponding *coverage probability* for the associated confidence set. So we are studying here both the operating characteristics of tests used and the coverage probabilities of the confidence sets defined above. For further insight on confidence set properties, we also study the frequency of unbounded outcomes and the width of the bounded ones.

Since available inference methods perform poorly when the underlying distributions are heavy-tailed, we designed our simulation experiments to cover such distributions by simulating the data from the Singh-Maddala distribution, which was found to successfully mimic observed income distributions for developed countries such as Germany (Brachmann et al., 1995). Another reason to use the Singh-Maddala distribution is that it was widely used in the literature which makes our results directly comparable to previously proposed inference methods. The CDF of the Singh-Maddala distribution can be written as

$$F_X(x) = 1 - \left[1 + \left(\frac{x}{b_X} \right)^{a_X} \right]^{-q_X} \quad (3.1)$$

where a_X , q_X and b_X are the three parameters defining the distribution. a_X influences both tails, while q_X only affects the right tail. b_X is a scale parameter to which we give little attention as the inequality measures considered in this paper are scale invariant. This distribution is a member of the five-parameter generalized beta distribution and its upper tail behaves like a Pareto distribution with a tail index equal to the product of the two shape parameters a_X and q_X ($\xi_X = a_X q_X$). The k -th moment exists for $-a_X < k < \xi_X$ which implies that a sufficient condition for the mean and the variance to exist is $-a_X < 2 < \xi_X$.

The moment of order γ of Singh-Maddala distribution have the following closed form:

$$v_X(\gamma) := \mathbb{E}(X^\gamma) = \frac{b_X^\gamma \Gamma(\gamma a_X^{-1} + 1) \Gamma(q_X - \gamma a_X^{-1})}{\Gamma(q_X)} \quad (3.2)$$

where $\Gamma(\cdot)$ is the gamma function. For $\gamma = 1$, this yields the mean of X [$\mu_X = v_X(1) = \mathbb{E}(X)$] and, for $\gamma = 2$, the second moment of X [$v_X(2) = \mathbb{E}(X^2)$]. Similarly, replacing X by Y in the above expressions, we can compute μ_Y and $v_Y(2)$. Using the values of these moments, we compute analytical expressions for $GE_\gamma(X)$ and $GE_\gamma(Y)$. Each experiment involves 10000 replications and sample sizes of $n = 50, 100, 250, 500, 1000, 2000$. The nominal level α is set at 5%.

The hypotheses of interest take the form $H_0(\gamma) : GE_\gamma(X) - GE_\gamma(Y) = \Delta_0$, for $\gamma = 1$ or 2. Even

though we emphasize the important problem of testing equality ($\Delta_0 = 0$), we also consider the problem of testing nonzero differences ($\Delta_0 \neq 0$). Our simulation experiments cover the following designs.

1. Experiment I – Independent samples of equal sizes ($m = n$):
 - (a) $\Delta_0 = 0$ with $F_X = F_Y$; (b) $\Delta_0 = 0$ with $F_X \neq F_Y$; and (c) $\Delta_0 \neq 0$ (hence $F_X \neq F_Y$).
2. Experiment II – Independent samples of unequal sizes ($m \neq n$):
 - (a) $\Delta_0 = 0$ with $F_X = F_Y$; (b) $\Delta_0 = 0$ with $F_X \neq F_Y$; and (c) $\Delta_0 \neq 0$ (hence $F_X \neq F_Y$).
3. Experiment III – Dependent samples of equal sizes ($m = n$):
 - (a) $\Delta_0 = 0$ with $F_X = F_Y$; (b) $\Delta_0 = 0$ with $F_X \neq F_Y$; and (c) $\Delta_0 \neq 0$ (hence $F_X \neq F_Y$).

The simulation results are presented graphically through plotting the rejection frequencies against the number of observations. When the number of observations is different between the two samples, we plot the rejection frequencies against the number of observations of the smallest sample.

For the Delta method, we use the critical region $[\Delta\widehat{GE}_\gamma - \Delta_0]^2 > z_{\alpha/2}^2 \widehat{V}[\Delta\widehat{GE}_\gamma]$, based on (2.13); for the Fieller method, we use the critical region $\widehat{\Theta}(\Delta_0)^2 > z_{\alpha/2}^2 \widehat{V}[\widehat{\Theta}(\Delta_0)]$, as described in (2.20). Power is investigated by assuming distributions with heavier left and right tails to draw the first sample, and distributions with less heavy left and right tails to draw the second sample. We do so by considering DGPs with a lower value of the shape parameter a_X and a higher value of the shape parameter a_Y . The rejection frequencies under the alternative are not size-controlled, yet we compare power when both methods have similar sizes.

Our extensive simulation study reveals several important results. First, the Fieller-type method outperforms the Delta method under most specifications, and when it does not, it performs as well as the Delta method. Put differently, the Fieller-type method was never dominated by Delta method. Second, the Fieller-type method is more robust to irregularities arising from both the left and right tails. Third, the Fieller-type method gains become more sizeable as the sensitivity parameter γ increases. Fourth, the performance of the Fieller-type method matches, and for some cases exceeds, the permutation method which is considered one of the best performing methods proposed in the literature so far for the two-sample problem. In the remainder of this section we take a closer look at the simulation evidence supporting the above findings.

Experiment I: Independent samples of equal sizes – The left panels of Figures B.1 and B.2 depict the rejection frequencies against the sample size for GE_1 and GE_2 respectively. Here the distributions are assumed identical [$F_X = F_Y$]. Comparing the two panels, we notice that better size control with the Fieller-type method is more noticeable for GE_2 : the size gains are larger when the

index used is more sensitive to the changes in the right tail of the underlying distributions. As the sample size increases the rejection probabilities of the two methods converge to the same level.

In the second specification, the indices are identical, but the underlying distributions are not $[\Delta_0 = 0 \text{ with } F_X \neq F_Y]$. The left panel of Figure B.3 plots the FCS and DCS rejection frequencies for this scenario. Again, the results suggest that the Fieller-type method outperforms the Delta method in small samples in terms of size, and the gains are most prominent for GE_2 . The gains are smaller in this scenario compared to the previous one. As we will show later, the Fieller-type method will not solve the over-rejection problem under all scenarios, but it will reduce size distortions in many cases, and when it does not, it performs as well as the Delta method.

We now move to the third scenario, where we consider different distributions under the null hypothesis and unequal inequality indices $[\Delta_0 \neq 0]$. In this scenario, the difference under the null hypothesis can take any admissible value (possibly different from zero). Testing a zero value, although informative, does not always translate into a confidence interval. Hence, one of our contributions lies in considering the non-zero null hypothesis which allows us to rely for inference on the more-informative confidence sets approach rather than testing the equality of the difference between the two indices to one specific value.

The results, as shown in the left panels of Figures B.5 and B.6, suggest a considerable improvement. In both panels, the Fieller-type method leads to size gains and almost achieves correct size. The improvements are more pronounced for the GE_2 index. The right panels of Figures B.1 to B.6 illustrate the power of FCS and DCS for both GE_1 and GE_2 under the three scenarios considered: $[\Delta_0 = 0 \text{ with } F_X = F_Y]$, $[\Delta_0 = 0 \text{ with } F_X \neq F_Y]$ and $[\Delta_0 \neq 0]$ respectively. The results show that the Fieller-type method is as powerful as the Delta method when compared at sample sizes where both FCS and DCS have similar empirical rejection frequencies.

Experiment II: Independent samples of unequal sizes – Empirically, when comparing inequality levels spatially or over time, it is unlikely one encounters samples with the same size. Thus, it is useful to assess the performance of our proposed method when the sample sizes are unequal. To do so, we adjust our simulation design by setting the number of observations of the second sample to be as twice as large as the first sample. If we denote the size of the first sample by n and that of the second by m , then $n = 2m$.⁴ The results are analogous to those obtained in the first experiment, under which sample sizes were equal, in the sense that the Fieller-type method improves level control for both GE_1 and GE_2 , with a larger improvement for GE_2 . The size and power simulation results for the three scenarios considered here are available in the online appendix.

Experiment III: Dependent samples of equal sizes – Another interesting case is the one where the samples are dependent. This occurs mostly when comparing inequality levels before and after a policy change, such as comparing pre-tax and post-tax income inequality levels, or comparing

⁴The results presented here are not sensitive to choice of the ratio between n and m

the distributional impact of a macroeconomic shock. To accommodate for such dependencies, we modify the simulation design as follows: the samples are drawn in pairs from the joint distribution, which we denote F_{XY} , where the correlation between the two marginal distributions is generated using a Gumbel copula with a high Kendall's correlation coefficient of 0.8. For this case, results are in line with the independent cases, in small samples and when larger γ is used. Size and power plots are available in the online appendix.

Comparing the Fieller-type method with the permutation method – As outlined in the introduction, the permutation-based Monte-Carlo test approach proposed in Dufour et al. (2019) stands out as one of the best performing nonparametric inference method for testing the equality of two inequality indices. The authors focus on the Theil and the Gini indices. The permutation testing approach provides exact inference when the null distributions are identical ($F_X = F_Y$) and it leads to a sizeable size distortion reduction when the null distributions are sufficiently close ($F_X \approx F_Y$). However, as the null distributions differ, the performance of the method deteriorates.

Figures B.7 and B.8 plot size and power of the permutation Fieller-type methods against the tail index of F_Y . As in Dufour et al. (2019), we fix the tail index of the null distribution F_X to 4.76. When the distributions under the null hypothesis are identical, the permutation method is exact and thus it is important to compare methods when exactness does not hold. Our results point to two main advantages of the Fieller-type method over the permutation method: for the Theil index, the Fieller-type method is more powerful and these power gains are magnified as the difference between the indices becomes larger. On the other hand, when considering the GE_2 , there are size gains mainly when the tail index is relatively small (*i.e.*, when the right tail is heavier). These size gains are not associated with power loss as the right panel of the same figure illustrates.

The attraction of the Fieller-type method with respect to the permutation approach goes beyond the superior performance highlighted above. Unlike the Fieller-type method, its applicability is restricted to the null hypothesis of equality ($\Delta_0 = 0$), and further theoretical developments would be needed to test more general hypotheses. Building confidence intervals using a permutation-based or another simulation-based method (such as the bootstrap) would also require a computationally intensive numerical inversion (*e.g.*, through a grid search). So another appealing feature of the Fieller-type approach comes from the fact that it is computationally easy to implement.

Behavior with respect to the tails – To better understand under what circumstances does the Fieller-type method improves level control, we assess the performance of the proposed method to different tail shapes. The literature has focused on the role of heavy right tails in the deterioration of the Delta method confidence sets. However, as our results indicate, heavy left tails also contribute to the under-performance of the standard inference procedures. The Fieller-type method is less prone to such irregularities arising from both ends of the distributions and thus it reduces size distortions whether the cause of the under-performance is arising from the left tail or the right

tail. This is supported by the results reported below in Tables C.3 and C.4. The results in these tables rely on samples of 50 observations. Table C.3 reports the percentage difference of the rejection frequencies as the right tails of the two distributions become thicker. The right-tail shape is determined by the tail index ($\xi_X = a_X q_X$). The smaller the tail index, the thicker is the right tail of the distribution under consideration. The reliability advantage of the Fieller-type method (over the Delta method) increases as the right tail of the distributions gets thicker.

To study the impact of the left tail, the parameters of the first distribution are fixed at $a_X = 2.8$ and $q_X = 1.7$, while a_Y and q_Y are varied such that the left tail becomes thicker and the right tail is left unchanged. This is done by decreasing a_Y , and increasing q_Y enough to keep the tail index fixed ($\xi_X = \xi_Y = 4.76$). The last column of Table C.4 shows the percentage difference of the rejection frequencies between the Fieller-type and Delta methods. As the left tail thickens, the performance of the Delta method deteriorates relative to the Fieller-type method, and thus the Fieller method better captures irregularities in the left tail. This conclusion holds regardless of whether the left tail of the second distribution is lighter or thicker than the left tail of the first distribution.

Fieller-type method and the sensitivity parameter γ – A consistent conclusion from our results is that the Fieller’s-induced size gains are more prominent for GE_2 compared to GE_1 , that is, when the sensitivity parameter γ increases from 1 to 2. This might suggest that as γ increases, size gains from the Fieller-type method increase. Such generalization is indeed supported by simulation evidence illustrated by Figure B.9. The left panel plots rejection frequencies of DCS and FCS for $\gamma \in [0, 3.5]$ for independent samples. The right panel considers dependent samples. As γ becomes larger, FCS outperforms DCS at an increasing rate. The superiority of the Fieller-type method in this context is unaffected by the independence assumption as shown in the right panel where the rejection frequencies are plotted against γ for dependent samples with Kendall’s correlation of 0.8.

Recall that the parameter γ characterizes the sensitivity of the index to changes at the tails of the distribution. For instance, the index becomes more sensitive to changes at the upper tails as γ increases (assuming positive γ). Thus, relative to the Delta method, the performance of the Fieller-type method in the two-sample problem improves as the right tail of the underlying distributions becomes heavier. This conclusion, as we saw from the results above, is robust to the assumptions about the independence of the samples and to the distance between the two null distributions.

The identical performance of the Fieller-type method and Delta method at $\gamma = 0$ is expected as the underlying t-tests inverted in the process of building FCS and DCS are identical since the null hypothesis is no longer a ratio. To see that, recall that the limiting solution for $GE_\gamma(\cdot)$ at $\gamma = 0$ is equal to $\mathbb{E}_F[\log(X)] - \log[\mathbb{E}_F(X)]$. Graphically, we can see that both methods start off at the same rejection frequencies when $\gamma = 0$, and then diverge as γ increases.

Robustness to the shapes of the null distributions – So far, our simulation experiments have focused on comparing the finite-sample performance of FCS and DCS by studying their behavior

as the number of observation increases, holding the parameters of the two underlying null distributions constant. Here we try to check the robustness of our results by fixing the number of observations at 50 and allowing the parameters (a_X, q_X, a_Y and q_Y) to vary. This type of analysis highlights the (in)sensitivity of our conclusions regarding the Fieller-type method to the shape of the null distributions. In left panel of Figure B.10, we plot the rejection frequencies of both methods against the sensitivity parameter ξ_X for the Theil index. We set ξ_X equal to 4.76 and allow ξ_Y to vary between 3.05 and 6.255. In the right panel, we focus on the GE_2 index. Here ξ_X is fixed at 4.76 again and the parameter ξ_Y ranges between 3.293 and 5.7107.

For small samples, the gains of the Fieller-type method are maintained regardless the shape of the distribution. The gains are more pronounced for GE_2 compared to GE_1 . These two graphs show that the gains attained by the Fieller-type method are not arbitrary and that they hold for various parametric assumptions of the underlying distributions.

Slow convergence – Inequality estimates are characterized by slow convergence when underlying distributions are heavy-tailed. This problem has in fact motivated most of the proposed asymptotic refinements in this literature [see Davidson and Flachaire (2007); Cowell and Flachaire (2007)]. Our results in Table C.6 corroborate this fact, as over-rejections remain even with samples as large as 200000, particularly with the GE_2 which puts more weight on the upper tail of the distribution. On balance, our main finding is the superiority of the Fieller method in finite samples.

Widths of the confidence sets – The last two columns of Table C.6 show the average widths of the FCS and the DCS for the two sample problem. Since the Fieller’s method can produce unbounded confidence sets, we take the average of the widths based on the bounded confidence sets. In general, compared to the FCS widths, the DCS widths are shorter with small samples, *i.e.* they are shorter when the Delta method rejection frequencies are higher than those of Fieller. This suggest that the DCS are too short and thus they tend to undercover the true difference between the indices. As the sample size increases, the two methods exhibit similar performance and the widths coincides.

4 Application: Regional economic convergence

In this section, we present empirical evidence on the relevance of our theoretical results to applied economic work. We assess economic convergence across the U.S. states between 1946 and 2016. In what follows, unless stated otherwise, tests and confidence sets are at the 5% level.

The late 1980’s witnessed a new wave of interest in economic convergence that was spurred by the revival of growth models. The convergence hypothesis, first theorized by the popular Solow growth model, postulates that in the long run, economies will converge to similar per-capita income levels. The convergence question is important from theoretical and policy perspectives.

Theoretically, Romer (1994) and Rebelo (1991) argue that the rejection of the convergence hypothesis provides empirical support for the endogenous growth model and evidence against the neoclassical growth model. In the latter models, per-capita income convergence results from the diminishing return to capital assumption. This assumption implies that the return to capital increases in economies with low level of capital and decreases in capital-abundant economies. Moreover, since the rate of return on capital is higher in poorer economies, investments will migrate from rich economies to poorer ones, further enhancing growth and reducing the gap between them. On the other hand, in endogenous growth models as in Romer (1994) and Rebelo (1991), the diminishing rate of return on capital is considered implausible once knowledge is assumed to be one of the production factors. Thus, the model does not predict convergence, but on the contrary predicts that divergence might occur.

Empirically, policy-makers are interested in learning about the dynamics of income dispersion across regions/states so they can engage in redistributive policies when needed or to assess the distributional impact of a specific policy. Among the various definitions of convergence provided in the literature, two definitions appear to dominate the work on this topic: β -convergence and σ -convergence (Barro, 2012; Barro and Sala-i Martin, 1992; Quah, 1996; Sala-i Martin, 1996; Higgins et al., 2006). Although related, these two measures might lead to different conclusions as they capture different dimensions of economic convergence. For an analytical treatment of the relationship between the two measures, see Higgins et al. (2006).

β -convergence occurs when there is a negative relationship between the growth rate and the initial level of per-capita income, that is, when poor economies grow at a faster rate than the rich ones. The σ -convergence concept focuses on the dispersion of the income distribution which is typically measured in this literature by the variance of the logs. The variance of logs is scale-independent and thus multiplying the per-capita incomes by a scale k has no impact on the dispersion level. Alternative scale-independent measures of dispersion such as inequality measures have generally not been utilized in convergence analysis. The only exception is Young et al. (2008) which reported the Gini coefficient for comparison purposes with reference to the variance of logs.

One feature of inequality measures such as the Gini coefficient and the GE measures is that they respect the Pigou-Dalton principle, which states that a rank preserving transfer from a richer individual/state to a poorer individual/state should make the distribution at least as equitable. In the context of economic convergence, this principle is particularly relevant. For instance, if the US government makes a transfer from a richer state to a poorer one, one would expect dispersion between states to decline. The Gini and GE measures would capture this decline, whereas the variance of logs might indicate no change or even an increase in dispersion. The fact that the variance of logs violates the Pigou-Dalton principle is usually neglected in the literature on the grounds that the problem occurs only at the extreme right tail of the distribution. However, Foster

and Ok (1999) show that disagreement between the variance of logs and inequality measures can result from changes in incomes in other parts of the distribution including the left tail. The following example (Foster and Ok, 1999) underscores the importance of the Pigou-Dalton principle and its implications for convergence. Consider two income distributions defined by the following incomes (2, 5, 10, 28, 40) and (2, 5, 10, 34, 34) where the latter is associated with a transfer from the richest [40 to 34] incomes to poorer ones [28 to 34]. The resulting change in the variance of logs, from 1.5125 to 1.5154, suggests an increase of inequality. In contrast, the GE_2 index declines from 0.3696 to 0.3446, thereby capturing the expected distributional impact of such a transfer.

Our empirical analysis of per-capita income dispersion across the US is motivated by comparably peculiar statistics. Consider the publicly available per-capita income at the state level for 48 out of the 50 states (as the data for Alaska and Hawaii is not available). The variance of logs between the years 2000 and 2016 indicates a 3% increase in dispersion, whereas GE_2 indicates a decline in dispersion by 0.3%. This provides a compelling basis for the more comprehensive inferential analysis reported next.

Using the same data source, we first compute the Theil index for the per-capita income distributions of 1946 and 2016. Then we construct the Delta and Fieller confidence sets for the difference between the two indices. A standard interpretation of differences between the two confidence intervals (at the considered level) implies that one will reject the null hypothesis $\Delta GE_\gamma = \Delta_0$ for a given Δ_0 while the other fails to reject it. Special attention should be paid to the $\Delta_0 = 0$ case, as decisions might reverse the conclusion on whether convergence holds or not.

Using the Theil index, our results in the first column of Table 4.1 indicate that per-capita income inequality across states has declined between 1946 and 2016. The decline in inequality implies convergence. This is compatible with the general convergence trend reported in the literature (Barro and Sala-i Martin, 1992; Bernat Jr, 2001; Higgins et al., 2006). Although the Fieller and Delta-method confidence sets are not identical, they still lead to the same conclusion which is that the decline of inequality is statistically different from zero at the level used.

In the second column of Table 4.1, we consider the same problem using GE_2 index rather than the Theil one. This index puts more weight on the right tail of the distribution. In this case, the results also indicate a decline of inequality across states. Inequality in 1946 was 0.02679 and declined by -0.01163 by 2016. The confidence sets based on the Delta and Fieller-type methods lead to opposite conclusions about the statistical significance of the decline in inequality. DCS fails to reject the null hypothesis of no change in inequality, thus the decline in inequality based on DCS is not statistically different from zero. On the other hand, the Fieller-type methods rejects the hypothesis of no inequality change, which entails that the decline is significant.

In addition to DCS and FCS, we report the permutational p -values. For the GE_2 , the p -value is less than 5% and thus we reject the null hypothesis of no change in inequality contradicting

Table 4.1: Estimates and confidence intervals of the change in inequality across U.S. states between 1946 and 2016.

	Theil Index / GE_1	GE_2
First sample - 1946	0.02743	0.02679
Second sample -2016	0.0144	0.01516
$GE_\gamma(2016) - GE_\gamma(1946)$	-0.01303	-0.01163
Delta C.I.	[-0.02486, -0.001204] Inequality decreases	[-0.02349, 0.00024] No change in Inequality
Fieller's C.I.	[-0.02531, -0.00155] Inequality decreases	[-0.02456, -0.00043] Inequality decreases
Permutation test p – Value	0.014 Inequality decreases	0.014 Inequality decreases
Number of states	48	48

Table 4.2: Estimates and confidence intervals of the change in inequality across non-OECD countries

	Theil Index / GE_1	GE_2
First sample - 1960	0.717621	1.46631
Second sample -2013	0.78726	1.45076
$GE_\gamma(2013) - GE_\gamma(1960)$	0.06964	-0.01554
Delta C.I.	[-0.35694, 0.49623]	[-1.15143, 1.120337]
Fieller's C.I.	[-0.40436, 0.63075]	\mathbb{R}
Permutation test p – value	0.886	0.992
Number of countries	72	72

the conclusion based on the Delta method. This constitutes an empirical evidence supporting the findings of Dufour et al. (2019).

Two conclusions can be drawn from our findings. First, the Fieller-type and the Delta methods can lead to different confidence sets in practice which documents the empirical relevance of our theoretical findings. Second, disparities between both sets can lead to spurious conclusions about inequality changes if one set includes zero while the other does not. From a policy point of view, this disparity is crucial, especially if important policy actions are motivated by underlying analysis.

We next turn to non-OECD countries between 1960 and 2013. Table 4.2 presents estimates and confidence sets for the difference of inequality measures between the two periods. The main finding here is that the Fieller-type confidence set based on the GE_2 index is the whole real line \mathbb{R} . These results confirm that decisions based on Delta-method are spurious, and that a no-change conclusion is flawed: data and measure are, instead, uninformative.

The permutational method leads results similar to Delta and the Fieller-type methods for non-OECD countries. Available permutation tests although preferable size-wise to their standard counterparts, are difficult to invert to build confidence sets. Instead, the confidence sets proposed here can be unbounded and thus avoid misleading statistical inferences and policy decisions, in particular from seemingly insignificant tests. The econometric literature on inequality has long emphasized the need to avoid over-sized tests. Rightfully, spurious rejections are misleading. Our results document a different although related problem: even with adequately sized no-change tests, weak identification can undercut the reliability of policy advice resulting from insignificant no-change test outcomes. Far more attention needs to be paid to confidence sets. Moreover, sets that can be unbounded make empirical and policy work far more credible than it can be using bounded alternatives or no-change tests that cannot be inverted.

5 Conclusion

This paper introduces a Fieller-type method for two-sample inference problem on the GE class of inequality indices. Simulation results confirm that the Fieller-type method outperforms standard counterparts including the permutation test. Size gains are most prominent when using indices that put more weight on the right tail of the distribution and results are robust to different assumptions about the shape of the null distributions. While irregularities arising from the right tail have long been documented, we find that left tail irregularities are equally important in explaining the failure of standard inference methods. On recalling that permutation tests are difficult to invert, our results underscore the usefulness of the Fieller-type method for evidence-based policy. An empirical analysis of economic convergence reinforces this result, and casts a new light on traditional controversies in the growth literature.

Fieller's approach is frequently applied in medical research and to a lesser extent in applied economics despite its solid theoretical foundations (Srivastava, 1986; Willan and O'Brien, 1996; Johannesson et al., 1996; Laska et al., 1997). This could be due to the seemingly counter-intuitive non-standard confidence sets it produces which economists often find hard to interpret. Consequently, many applied researchers encountering the estimation of ratios avoid using it and opt to use methods that yield closed intervals regardless of theoretical validity. This paper illustrates serious empirical and policy flaws that may result from such practices in inequality analysis.

References

- Andrews, D. W. and Cheng, X. (2013). Maximum likelihood estimation and uniform inference with sporadic identification failure. *Journal of Econometrics*, 173(1):36–56.
- Andrews, I. and Mikusheva, A. (2015). Maximum likelihood inference in weakly identified dynamic stochastic general equilibrium models. *Quantitative Economics*, 6(1):123–152.
- Araar, A. and Duclos, J.-Y. (2006). *Poverty and equity: Measurement, policy, and estimation with DAD*. Springer, New York.
- Bahadur, R. R. and Savage, L. J. (1956). The nonexistence of certain statistical procedures in nonparametric problems. *The Annals of Mathematical Statistics*, 27(4):1115–1122.
- Barro, R. J. (2012). Convergence and modernization revisited. Technical report, National Bureau of Economic Research.
- Barro, R. J. and Sala-i Martin, X. (1992). Convergence. *Journal of Political Economy*, 100(2):223–251.
- Beaulieu, M.-C., Dufour, J.-M., and Khalaf, L. (2013). Identification-robust estimation and testing of the zero-beta CAPM. *Review of Economic Studies*, 80(3):892–924.
- Bennett, B. (1959). On a multivariate version of Fieller’s theorem. *Journal of the Royal Statistical Society: Series B (Methodological)*, 21(1):59–62.
- Bernard, J.-T., Chu, B., Khalaf, L., Voia, M., et al. (2019). Non-standard confidence sets for ratios and tipping points with applications to dynamic panel data. *Annals of Economics and Statistics*, (134):79–108.
- Bernard, J.-T., Idoudi, N., Khalaf, L., and Yélou, C. (2007). Finite sample inference methods for dynamic energy demand models. *Journal of Applied Econometrics*, 22(7):1211–1226.
- Bernat Jr, G. A. (2001). Convergence in state per capita personal income, 1950-99. *Survey of Current Business*, 81(6):36–48.
- Bertanha, M. and Moreira, M. J. (2016). Impossible inference in econometrics: Theory and applications. *arXiv preprint arXiv:1612.02024*.
- Blundell, R., Pistaferri, L., and Preston, I. (2008). Consumption inequality and partial insurance. *American Economic Review*, 98(5):1887–1921.
- Bolduc, D., Khalaf, L., and Yélou, C. (2010). Identification robust confidence set methods for inference on parameter ratios with application to discrete choice models. *Journal of Econometrics*, 157(2):317–327.
- Brachmann, K., Stich, A., and Trede, M. (1995). Evaluating parametric income distribution models. *Allgemeines Statistisches Archiv*, 80:285–298.
- Cowell, F. A. (2000). Measurement of inequality. In *Handbook of Income Distribution*, volume 1, pages 87–166. Elsevier, Amsterdam.
- Cowell, F. A. and Flachaire, E. (2007). Income distribution and inequality measurement: The

- problem of extreme values. *Journal of Econometrics*, 141(2):1044–1072.
- Cowell, F. A. and Flachaire, E. (2015). Statistical methods for distributional analysis. In *Handbook of Income Distribution*, volume 2, pages 359–465. Elsevier, Amsterdam.
- Cowell, F. A. and Flachaire, E. (2018). Inequality measurement and the rich: why inequality increased more than we thought. *Suntory and Toyota International Centres for Economics and Related Disciplines, LSE*.
- Cox, D. R. (1967). Fieller’s theorem and a generalization. *Biometrika*, 54(3-4):567–572.
- Davidson, R. and Flachaire, E. (2007). Asymptotic and bootstrap inference for inequality and poverty measures. *Journal of Econometrics*, 141(1):141–166.
- Dufour, J.-M. (1997). Some impossibility theorems in econometrics with applications to structural and dynamic models. *Econometrica*, 65(6):1365–1387.
- Dufour, J.-M., Flachaire, E., and Khalaf, L. (2019). Permutation tests for comparing inequality measures. *Journal of Business & Economic Statistics*, 37(3):457–470.
- Dufour, J.-M. and Hsiao, C. (2008). Identification. In *The New Palgrave Dictionary of Economics*, The New Palgrave Dictionary of Econometrics. Palgrave MacMillan, 2nd edition.
- Dufour, J.-M. and Jasiak, J. (2001). Finite sample limited information inference methods for structural equations and models with generated regressors. *International Economic Review*, 42(3):815–844.
- Dufour, J.-M. and Taamouti, M. (2005). Projection-based statistical inference in linear structural models with possibly weak instruments. *Econometrica*, 73(4):1351–1365.
- Dufour, J.-M. and Taamouti, M. (2007). Further results on projection-based inference in IV regressions with weak, collinear or missing instruments. *Journal of Econometrics*, 139(1):133–153.
- Foster, J. E. and Ok, E. A. (1999). Lorenz dominance and the variance of logarithms. *Econometrica*, 67(4):901–907.
- Franz, V. H. (2007). Ratios: A short guide to confidence limits and proper use. *arXiv preprint arXiv:0710.2024*.
- Gleser, L. J. and Hwang, J. T. (1987). The nonexistence of 100 (1- α)% confidence sets of finite expected diameter in errors-in-variables and related models. *The Annals of Statistics*, 15(4):1351–1362.
- Higgins, M. J., Levy, D., and Young, A. T. (2006). Growth and convergence across the United States: Evidence from county-level data. *The Review of Economics and Statistics*, 88(4):671–681.
- Hirschberg, J. and Lye, J. (2010). A geometric comparison of the delta and Fieller confidence intervals. *The American Statistician*, 64(3):234–241.
- Johannesson, M., Jönsson, B., and Karlsson, G. (1996). Outcome measurement in economic evaluation. *Health economics*, 5(4):279–296.

- Kleibergen, F. (2005). Testing parameters in GMM without assuming that they are identified. *Econometrica*, 73(4):1103–1123.
- Koschat, M. A. et al. (1987). A characterization of the Fieller solution. *The Annals of Statistics*, 15(1):462–468.
- Laska, E. M., Meisner, M., and Siegel, C. (1997). Statistical inference for cost–effectiveness ratios. *Health economics*, 6(3):229–242.
- Quah, D. T. (1996). Empirics for economic growth and convergence. *European Economic Review*, 40(6):1353–1375.
- Rebelo, S. (1991). Long-run policy analysis and long-run growth. *Journal of Political Economy*, 99(3):500–521.
- Romer, P. M. (1994). The origins of endogenous growth. *Journal of Economic Perspectives*, 8(1):3–22.
- Sala-i Martin, X. X. (1996). The classical approach to convergence analysis. *The Economic Journal*, 106(437):1019–1036.
- Schluter, C. (2012). On the problem of inference for inequality measures for heavy-tailed distributions. *The Econometrics Journal*, 15(1):125–153.
- Shorrocks, A. F. (1980). The class of additively decomposable inequality measures. *Econometrica*, 48(3):613–625.
- Srivastava, M. (1986). Multivariate bioassay, combination of bioassays, and Fieller’s theorem. *Biometrics*, 42(1):131–141.
- Theil, H. (1967). *Economics and information theory*. North-Holland, Amsterdam.
- Willan, A. R. and O’Brien, B. J. (1996). Confidence intervals for cost-effectiveness ratios: An application of Fieller’s theorem. *Health Economics*, 5(4):297–305.
- Young, A. T., Higgins, M. J., and Levy, D. (2008). Sigma convergence versus beta convergence: Evidence from US county-level data. *Journal of Money, Credit and Banking*, 40(5):1083–1093.
- Zerbe, G. O. (1978). On Fieller’s theorem and the general linear model. *The American Statistician*, 32(3):103–105.

Appendices

A Proof of Theorem 2.1

From (2.22) and (2.21), the Fieller-type confidence region for Δ can be rewritten as follows:

$$\begin{aligned}
 \text{FCS}[\Delta GE\gamma; 1 - \alpha] &= \{\Delta_0 : \hat{\Theta}(\Delta_0)^2 \leq z_{\alpha/2}^2 \hat{V}[\hat{\Theta}(\Delta_0)]\} \\
 &= \{\Delta_0 : [\hat{\theta}_1 - \hat{\theta}_2 \Delta_0]^2 \leq z_{\alpha/2}^2 \hat{V}[\hat{\Theta}(\Delta_0)]\} \\
 &= \{\Delta_0 : (\hat{\theta}_2^2 \Delta_0^2 - 2\hat{\theta}_1 \hat{\theta}_2 \Delta_0 + \hat{\theta}_1^2) \leq z_{\alpha/2}^2 [\hat{V}(\hat{\theta}_1) - 2\hat{C}(\hat{\theta}_1, \hat{\theta}_2) \Delta_0 + \hat{V}(\hat{\theta}_2) \Delta_0^2]\} \\
 &= \{\Delta_0 : [\hat{\theta}_2^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_2)] \Delta_0^2 + 2[z_{\alpha/2}^2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) - \hat{\theta}_1 \hat{\theta}_2] \Delta_0 + [\hat{\theta}_1^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_1)] \leq 0\} \\
 &= \{\Delta_0 : A \Delta_0^2 + B \Delta_0 + C \leq 0\}, \tag{A.1}
 \end{aligned}$$

where

$$A := \hat{\theta}_2^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_2), \quad B := 2[z_{\alpha/2}^2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) - \hat{\theta}_1 \hat{\theta}_2], \quad C := \hat{\theta}_1^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_1). \tag{A.2}$$

The expression (2.24) then follows on solving the quadratic inequality in (A.1) for Δ_0 , where

$$\begin{aligned}
 D &= B^2 - 4AC \\
 &= 4[z_{\alpha/2}^2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) - \hat{\theta}_1 \hat{\theta}_2]^2 - 4[\hat{\theta}_2^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_2)][\hat{\theta}_1^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_1)] \\
 &= 4\{[z_{\alpha/2}^4 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2 + \hat{\theta}_1^2 \hat{\theta}_2^2 - 2z_{\alpha/2}^2 \hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)] \\
 &\quad - [\hat{\theta}_1^2 \hat{\theta}_2^2 + z_{\alpha/2}^4 \hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - z_{\alpha/2}^2 \hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) - z_{\alpha/2}^2 \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1)]\} \\
 &= 4z_{\alpha/2}^2 \{[z_{\alpha/2}^2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2 - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)] - [z_{\alpha/2}^2 \hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) - \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1)]\} \\
 &= 4z_{\alpha/2}^2 \{[\hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)] \\
 &\quad + z_{\alpha/2}^2 [\hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2 - \hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2)]\}; \tag{A.3}
 \end{aligned}$$

for similar arguments, see Dufour and Jasiak (2001) or Bolduc et al. (2010).

Suppose now that $\hat{V}(\hat{\theta})$ is positive definite with $z_{\alpha/2}^2 > 0$. Then,

$$\det[\hat{V}(\hat{\theta})] = \hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2 > 0 \tag{A.4}$$

and

$$\begin{aligned}
\hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) &= (\hat{\theta}_1, -\hat{\theta}_2)' \begin{bmatrix} \hat{V}(\hat{\theta}_1) & \hat{C}(\hat{\theta}_1, \hat{\theta}_2) \\ \hat{C}(\hat{\theta}_1, \hat{\theta}_2) & \hat{V}(\hat{\theta}_2) \end{bmatrix} \begin{bmatrix} \hat{\theta}_1 \\ -\hat{\theta}_2 \end{bmatrix} \\
&= (\hat{\theta}_1, -\hat{\theta}_2)' \hat{V}(\hat{\theta}) \begin{bmatrix} \hat{\theta}_1 \\ -\hat{\theta}_2 \end{bmatrix} \geq 0.
\end{aligned} \tag{A.5}$$

From (A.3), we see that

$$\begin{aligned}
D < 0 &\iff \{ \hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) < z_{\alpha/2}^2 [\hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2] \} \\
&\iff z^* := \frac{\hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)}{\hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2} < z_{\alpha/2}^2.
\end{aligned} \tag{A.6}$$

Further,

$$\begin{aligned}
\frac{\hat{\theta}_2^2}{\hat{V}(\hat{\theta}_2)} - z^* &= \frac{\hat{\theta}_2^2}{\hat{V}(\hat{\theta}_2)} - \frac{\hat{\theta}_1^2 \hat{V}(\hat{\theta}_2) + \hat{\theta}_2^2 \hat{V}(\hat{\theta}_1) - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)}{\hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2} \\
&= \frac{-[\hat{\theta}_2^2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2 + \hat{\theta}_1^2 \hat{V}(\hat{\theta}_2)^2 - 2\hat{\theta}_1 \hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) \hat{V}(\hat{\theta}_2)]}{\hat{V}(\hat{\theta}_2) [\hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2]} \\
&= \frac{-[\hat{\theta}_2 \hat{C}(\hat{\theta}_1, \hat{\theta}_2) - \hat{\theta}_1 \hat{V}(\hat{\theta}_2)]^2}{\hat{V}(\hat{\theta}_2) [\hat{V}(\hat{\theta}_1) \hat{V}(\hat{\theta}_2) - \hat{C}(\hat{\theta}_1, \hat{\theta}_2)^2]} < 0
\end{aligned} \tag{A.7}$$

so that $D < 0$ implies

$$\frac{\hat{\theta}_2^2}{\hat{V}(\hat{\theta}_2)} < z^* < z_{\alpha/2}^2 \tag{A.8}$$

and

$$A = \hat{\theta}_2^2 - z_{\alpha/2}^2 \hat{V}(\hat{\theta}_2) < 0. \tag{A.9}$$

Finally, to see that $C < 0$, we simply observe that $D < 0$ implies $B^2 < 4AC$, hence on using $A < 0$,

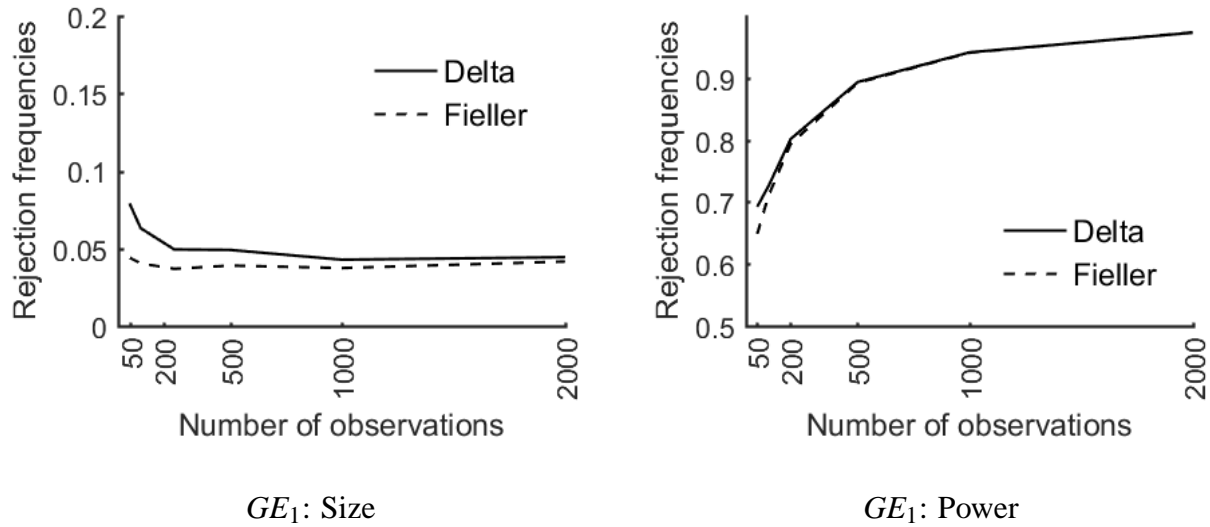
$$C < \frac{B^2}{4A} < 0. \tag{A.10}$$

This establishes (2.27).

B Figures

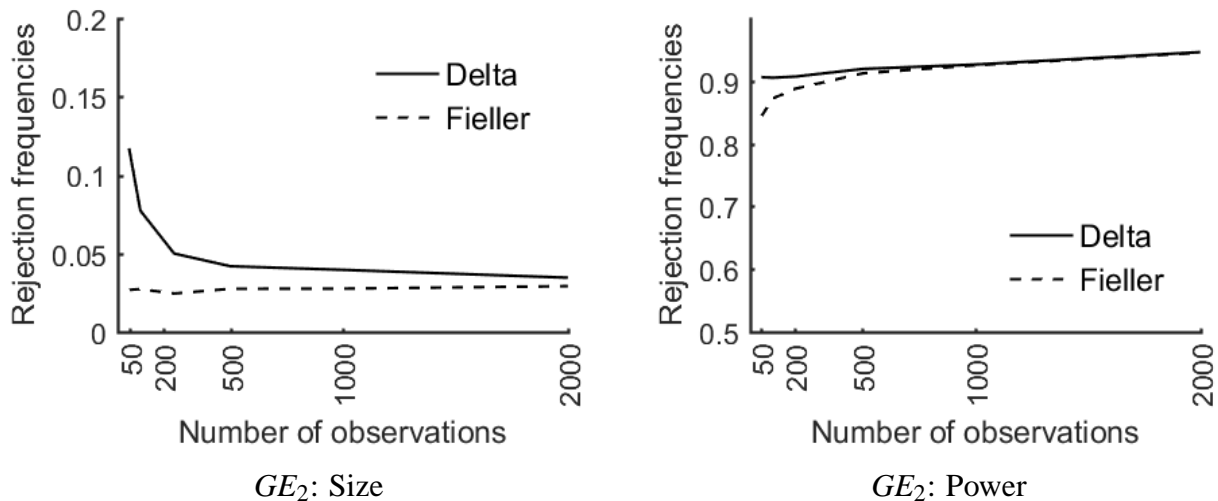
B.1 Experiment I; Design (I-a) – Independent samples: $n = m, F_X = F_Y, \Delta_0 = 0$

Figure B.1: Design (I-a) – Size and power of Delta and Fieller-type tests for GE_1 comparisons
 $H_0: GE_1(X) = GE_1(Y)$, Nominal size = 0.05



Left panel: $SM_X(a_X = 5.8, q_X = 0.499616)$, $SM_Y(a_Y = 5.8, q_Y = 0.499616)$. $GE_1(X) = GE_1(Y) = 0.14011$
 Right panel: $SM_X(a_X = 4.8, q_X = 0.499616)$, $SM_Y(a_Y = 6.8, q_Y = 0.499616)$. $GE_1(X) = 0.22857$,
 $GE_1(Y) = 0.09514$

Figure B.2: Design (I-a) – Size and power of Delta and Fieller-type tests for GE_2 comparisons
 $H_0: GE_2(X) = GE_2(Y)$, Nominal size = 0.05

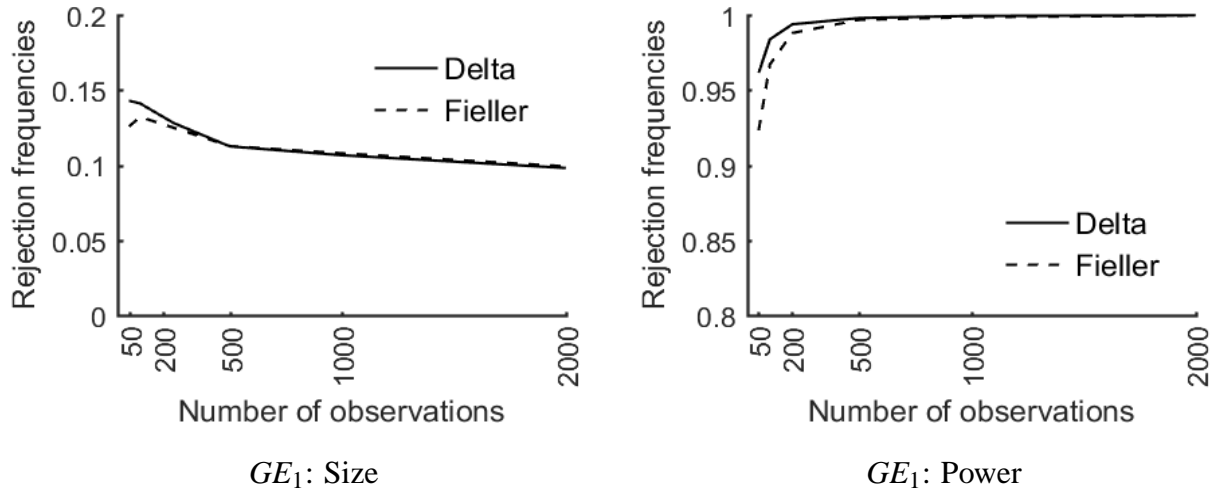


Left panel: $SM_X(a_X = 5.8, q_X = 0.499616)$, $SM_Y(a_Y = 5.8, q_Y = 0.499616)$. $GE_2(X) = GE_2(Y) = 0.24396$
 Right panel: $SM_X(a_X = 4.8, q_X = 0.499616)$, $SM_Y(a_Y = 6.8, q_Y = 0.499616)$. $GE_2(X) = 0.63705$,
 $GE_2(Y) = 0.13806$.

B.2 Experiment I; Design (I-b) – Independent samples: $n = m, F_X \neq F_Y, \Delta_0 = 0$

Figure B.3: Design (I-b) – Size and power of Delta and Fieller-type tests for GE_1 comparisons.

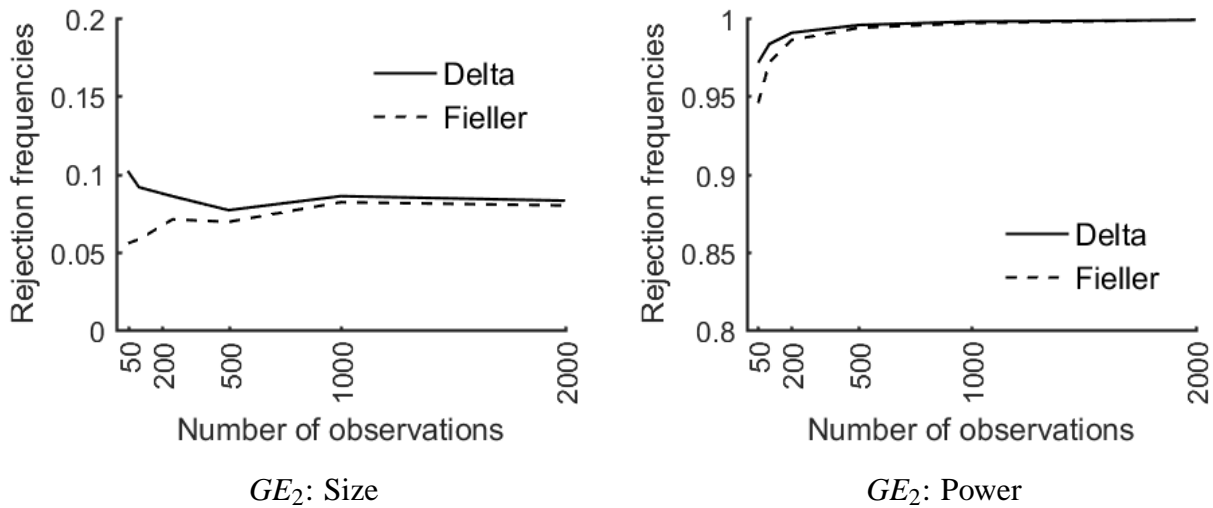
$$H_0: GE_1(X) = GE_1(Y), \text{Nominal size} = 0.05$$



Left panel: $SM_X(a_X = 2.8, q_X = 1.7), SM_Y(a_Y = 5.8, q_Y = 0.499616)$. $GE_1(X) = GE_1(Y) = 0.14011$
 Right panel: $SM_X(a_X = 1.8, q_X = 1.7), SM_Y(a_Y = 6.8, q_Y = 0.499616)$. $GE_1(X) = 0.33830, GE_1(Y) = 0.09514$.

Figure B.4: Design (I-b) – Size and power of Delta and Fieller-type tests for GE_2 comparisons

$$H_0: GE_2(X) = GE_2(Y), \text{Nominal size} = 0.05$$

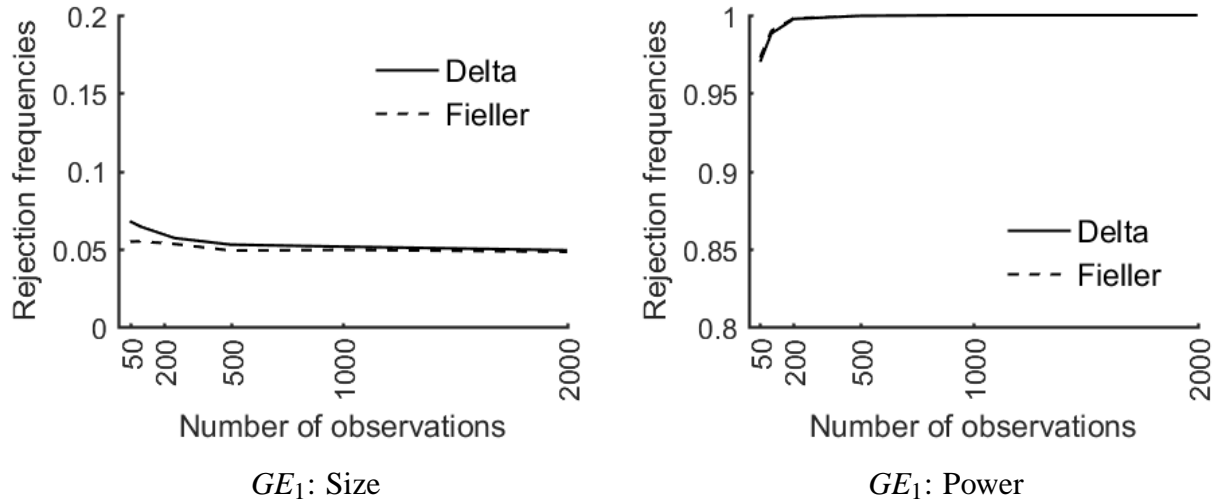


Left panel: $SM_X(a_X = 2.8, q_X = 1.7), SM_Y(a_Y = 3.8, q_Y = 0.9831)$. $GE_2(X) = GE_2(Y) = 0.16204$
 Right panel: $SM_X(a_X = 1.8, q_X = 1.7), SM_Y(a_Y = 4.8, q_Y = 0.9831)$. $GE_2(X) = 0.5479, GE_2(Y) = 0.08835$.

B.3 Experiment I; Design (I-c) – Independent samples: $n = m, F_X \neq F_Y, \Delta_0 \neq 0$

Figure B.5: Design (I-c) – Size and power of Delta and Fieller-type tests for GE_1 comparisons

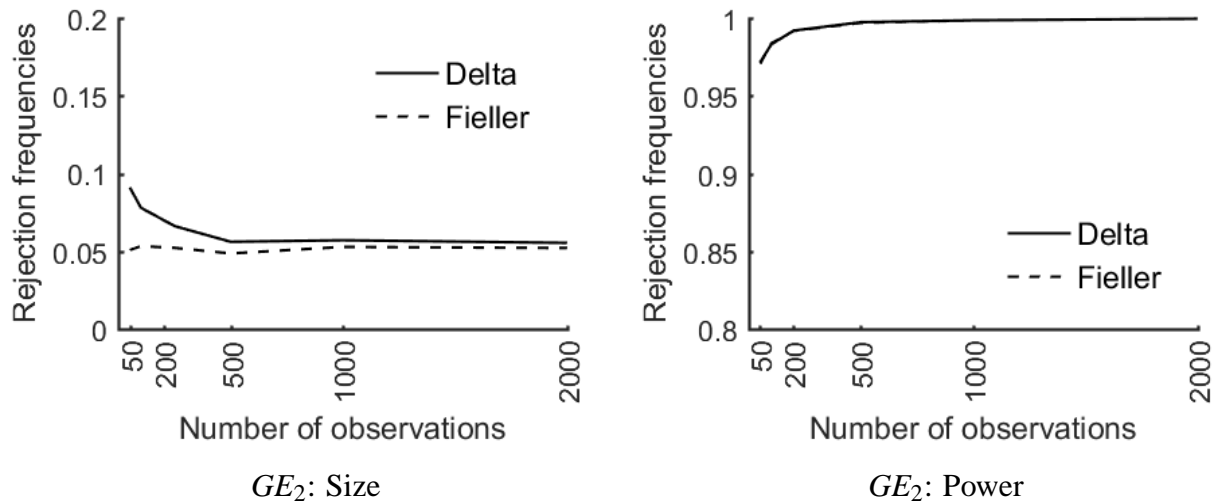
$$H_0: GE_1(X) - GE_1(Y) = 0.04670, \text{Nominal size} = 0.05$$



Left panel: $SM_X(a_X = 2.8, q_X = 1.7), SM_Y(a_Y = 3.8, q_Y = 1.3061)$. $GE_1(X) = 0.14011, GE_1(Y) = 0.09340$
 Right panel: $SM_X(a_X = 1.8, q_X = 1.7), SM_Y(a_Y = 4.8, q_Y = 1.3061)$. $GE_1(X) = 0.33829, GE_1(Y) = 0.05839$

Figure B.6: Design (I-c) – Size and power of Delta and Fieller-type tests for GE_2 comparisons

$$H_0: GE_2(X) - GE_2(Y) = 0.05401, \text{Nominal size} = 0.05$$

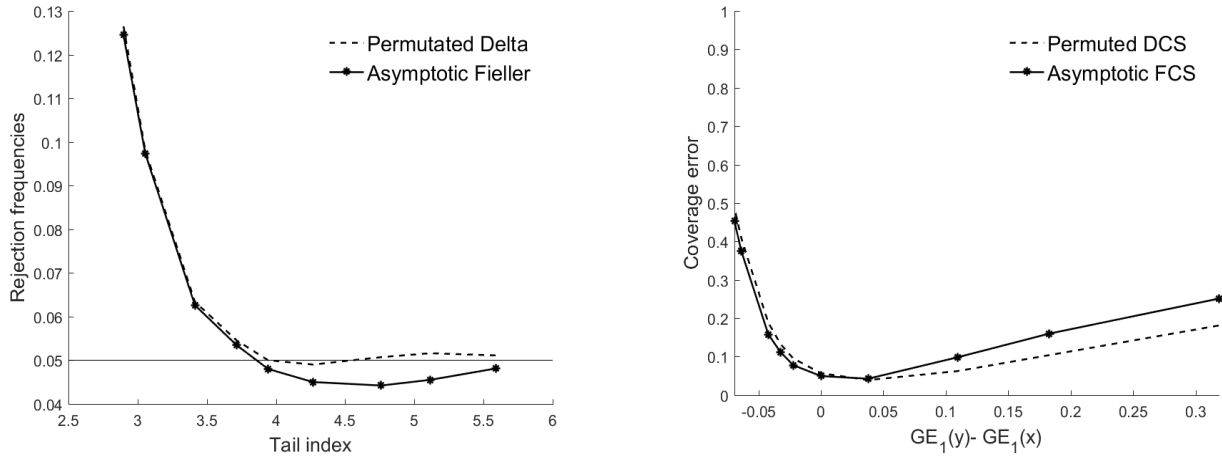


Left panel: $SM_X(a_X = 2.8, q_X = 1.7), SM_Y(a_Y = 3.8, q_Y = 1.2855)$. $GE_2(X) = 0.16203, GE_2(Y) = 0.10802$
 Right panel: $SM_X(a_X = 1.8, q_X = 1.7), SM_Y(a_Y = 4.8, q_Y = 1.2855)$. $GE_2(X) = 0.54790, GE_2(Y) = 0.06367$

B.4 Comparing Fieller's method and the permutation method

Figure B.7: Size and Power of two-sample tests

Independent samples: $n = m$. $F_X = F_Y$, $GE_1(X) = GE_1(Y)$. Nominal size = 0.05



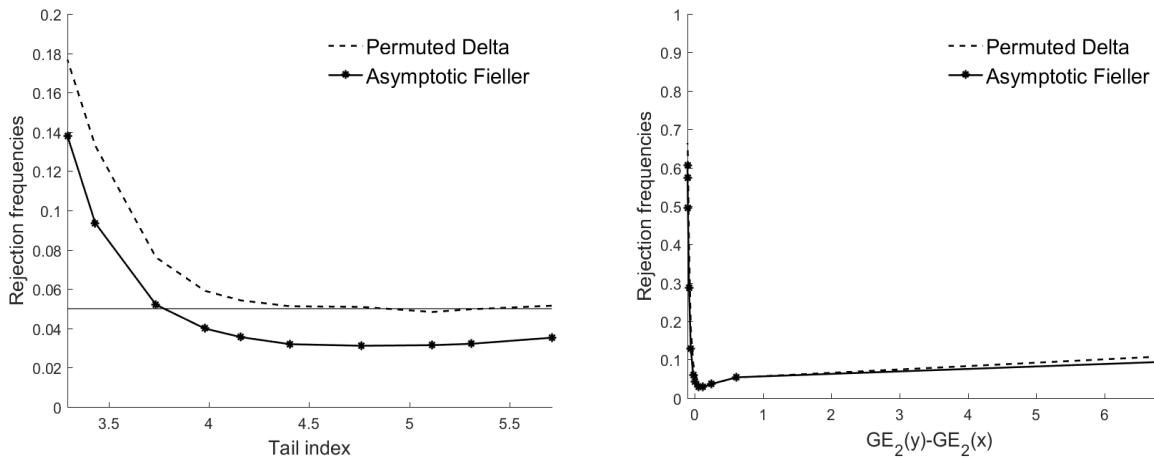
GE_1 : Size

GE_1 : Power

Note – The left panel pertains to test levels: the rejection frequencies of asymptotic the Fieller-type and Permutated Delta method are plotted against the tail index: $\xi = [2.897, 6.256]$. Power analysis is presented in the right panel: rejection frequencies are plotted against the difference between the two indices $GE_1(Y) - GE_1(X)$, with $q_Y = 10$.

Figure B.8: Size and Power of two-sample tests

Independent samples: $n = m$. $F_X = F_Y$ and $GE_\gamma(X) = GE_\gamma(Y)$. Nominal size = 0.05



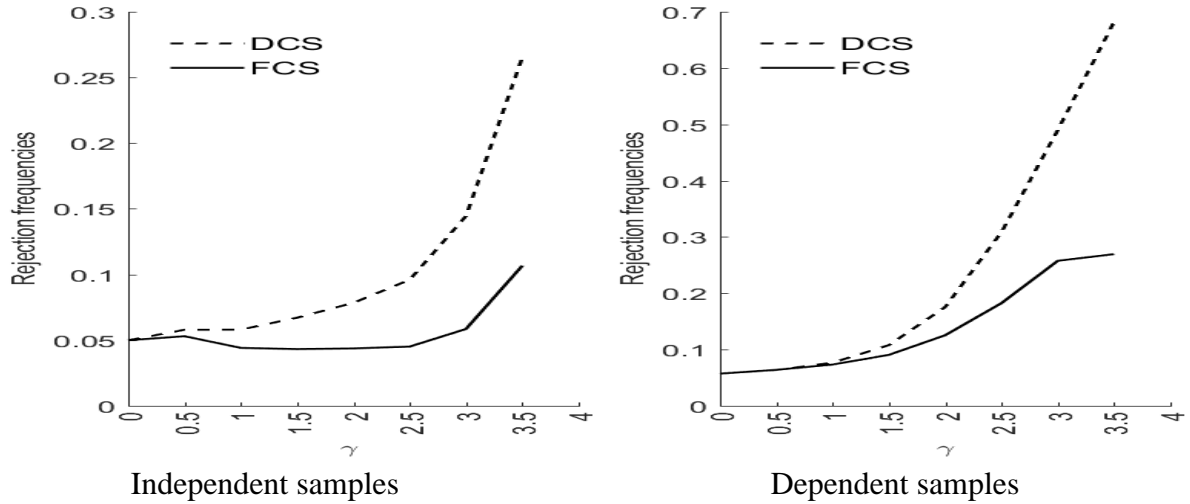
GE_2 : Size

GE_2 : Power

Note – The left panel pertains to test levels: the rejection frequencies of asymptotic Fieller-type and Permutated delta methods are plotted against the tail index: $\xi = [2.897, 6.256]$. Power analysis is presented in the right panel: the rejection frequencies are plotted against the difference between the two indices $GE_2(Y) - GE_2(X)$, with $q_2 = 10$.

B.5 Behavior with respect to the sensitivity parameter γ

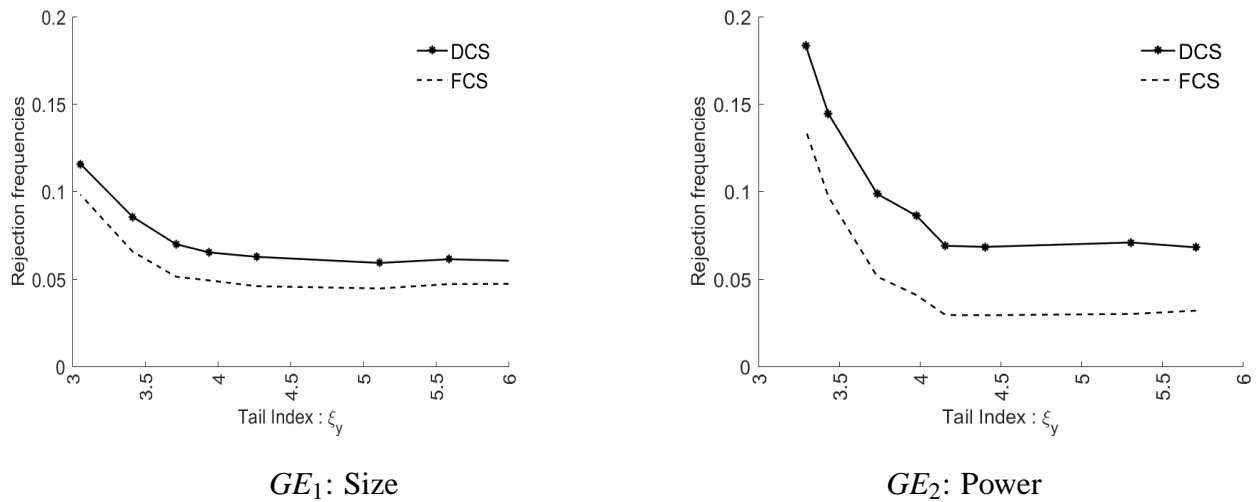
Figure B.9: Rejection frequencies of the tests inverted to derive the Delta method and Fieller's confidence sets over the sensitivity parameter γ . *Nominal size* = 0.05



Note – The distributions under the null hypothesis are identical and defined by: $SM_X(a_X = 2.8, q_X = 1.7)$ and $SM_Y(a_Y = 2.8, q_Y = 2)$. $n = m = 50$.

B.6 Robustness to the shape of the null distributions

Figure B.10: Rejection frequencies of the tests inverted to derive the Delta method and Fieller's confidence sets over the tail index ξ_Y . *Nominal size* = 0.05



Note – In the left panel, we consider the Theil index where ξ_X is fixed at 4.76 and $\xi_Y = [3.055, 6.255]$. In the right panel, we consider GE_2 with ξ_X is fixed at 4.76 and $\xi_Y = [3.293, 5.7107]$. $n = m = 50$.

C Tables

C.1 Effect of right tail thickness

Table C.3: Rejection frequencies of Delta and Fieller methods:
effect of right-tail thickness; n=50.

a_X	q_X	q_Y	a_Y	$\xi_X = \xi_Y$	$GE_1(X) = GE_1(Y)$	$GE_2(X) = GE_2(Y)$	PDL - GE_1	PDL - GE_2
5	2.1	5	2.1	10.5	0.04075	0.04096	2.84	10.58
5	1.9	5	1.9	9.5	0.04268	0.04326	4.47	13.99
5	1.7	5	1.7	8.5	0.04524	0.04639	5.19	20.59
5	1.5	5	1.5	7.5	0.0488	0.05084	8.81	24.2
5	1.3	5	1.3	6.5	0.05401	0.05763	13.60	31.96
5	1.1	5	1.1	5.5	0.06230	0.06906	16.88	42.72
5	0.9	5	0.9	4.5	0.07708	0.09155	29.70	56.74
5	0.7	5	0.7	3.5	0.10877	0.15046	36.87	66.55
5	0.5	5	0.5	2.5	0.20464	0.49151	53.75	84.86

Note – PDL stands for the percentage difference of the levels of the Delta and the Fieller-type method. The results in this table pertain to the percentage difference of the DCS and FCS levels as the right tails of both distributions gets thicker. The left tails of both distributions are fixed (a_X and a_Y are fixed) and the right tails gets thicker (with smaller ξ_X and ξ_Y). Column 8 reports the percentage difference associated with the null hypothesis $H_{01}: GE_1(X) = GE_1(Y)$ and column 9 reports the percentage difference associated with the null hypothesis $H_{02}: GE_2(X) = GE_2(Y)$.

C.2 Effect of left tail thickness

Table C.4: Rejection frequencies of Delta and Fieller methods: effect of left-tail thickness; $n=50$.

a_X	q_X	a_Y	q_Y	$\xi_X = \xi_Y$	$GE_1(X)$	$GE_2(X)$	$GE_1(Y)$	$GE_2(Y)$	$\Delta_{0,1}$	$\Delta_{0,2}$	PDL- GE_1	PDL- GE_2
2.8	1.7	5.8	0.821	4.76	0.14012	0.16204	0.0628	0.07347	0.07732	0.08857	18.74	38.6
2.8	1.7	5.2	0.915	4.76	0.14012	0.16204	0.06957	0.08095	0.07055	0.08109	20.80	39.95
2.8	1.7	4.8	0.992	4.76	0.14012	0.16204	0.07524	0.0872	0.06488	0.07484	19.78	41.72
2.8	1.7	4.2	1.133	4.76	0.14012	0.16204	0.08666	0.09998	0.05346	0.06206	21.35	45.9
2.8	1.7	3.8	1.253	4.76	0.14012	0.16204	0.09685	0.11148	0.04327	0.05056	23.41	48.41
2.8	1.7	3.2	1.488	4.76	0.14012	0.16204	0.11866	0.13661	0.02146	0.02543	23.82	52.19
2.8	1.7	3	1.587	4.76	0.14012	0.16204	0.12848	0.14816	0.01164	0.01388	26.22	56.03
2.8	1.7	2.6	1.831	4.76	0.14012	0.16204	0.15401	0.17888	-0.01389	-0.01684	25.37	56.94
2.8	1.7	2.4	1.983	4.76	0.14012	0.16204	0.17092	0.19982	-0.0308	-0.03778	27.01	58.79
2.8	1.3	1.3	2.8	3.64	0.17535	0.23133	0.44414	0.682	-0.26879	-0.45067	30.08	49.28
2.8	1.3	1.5	2.426	3.64	0.17535	0.23133	0.36978	0.53789	-0.19443	-0.30656	32.65	52.58
2.8	1.3	1.7	2.141	3.64	0.17535	0.23133	0.31577	0.44332	-0.14042	-0.21199	32.77	54.64
2.8	1.3	1.9	1.915	3.64	0.17535	0.23133	0.27516	0.37727	-0.09981	-0.14594	32.98	59.81
2.8	1.3	2.1	1.733	3.64	0.17535	0.23133	0.24375	0.32895	-0.0684	-0.09762	37.37	62.30
2.8	1.3	2.3	1.582	3.64	0.17535	0.23133	0.21891	0.29233	-0.04356	-0.061	40.25	66.83
2.8	1.3	2.5	1.456	3.64	0.17535	0.23133	0.19888	0.26379	-0.02353	-0.03246	41.10	71.34

Note- PDL stands for the percentage difference of the levels of the Delta and the Fieller's method. The results in this table pertain to the percentage difference of the DCS and FCS levels as the left tails of both distributions gets thicker. The right tails of both distributions are fixed ($\xi_X = \xi_Y = 4.76$) while the left tail of the second distribution gets thicker (with smaller a_Y). Column 12 reports the percentage difference associated with the null hypothesis $H01: GE_1(X) - GE_1(Y) = \Delta_{0,1}$ and column 13 reports the percentage difference associated with the null hypothesis $H02: GE_1(X) - GE_1(Y) = \Delta_{0,2}$. The values of $\Delta_{0,1}$ and $\Delta_{0,2}$ are given in columns 10 and 11 respectively.

C.3 Boundedness and width of the confidence intervals

Table C.5: Rejection probabilities and widths of confidence sets based on the Delta and Fieller-type methods: One-sample problem

n	Rejection Delta	Rejection Fieller	Bounded	Union of two disjoint sets	Unbounded	Width Fieller	Width Delta
50	0.3758	0.2616	9841	105	54	1.4339	0.6316
100	0.3211	0.2773	9983	16	1	0.7616	0.6026
200	0.2707	0.258	9998	2	0	0.6324	0.5462
500	0.2219	0.2244	10000	0	0	0.4482	0.4325
1000	0.1764	0.1796	10000	0	0	0.3635	0.3575
2000	0.1626	0.167	10000	0	0	0.2746	0.2726
10000	0.1077	0.1095	10000	0	0	0.1474	0.1472
20000	0.0990	0.1006	10000	0	0	0.1098	0.1097
100000	0.0753	0.0756	10000	0	0	0.0544	0.0544
200000	0.0686	0.0698	10000	0	0	0.0395	0.0395

Note – The coverage rate of the confidence set is equal to $1 - (\text{Rejection probability})$. The results in this table pertains to the GE_2 index with $SM_X(a_X = 1.1, q_X = 4.327273)$. $H_0: GE_2 = 0.71577$.

Table C.6: Rejection probabilities and widths of confidence sets based on the Delta and Fieller-type methods

n	Rejection Delta	Rejection Fieller	Bounded	Union of two disjoint sets	Unbounded	Width Fieller	Width Delta
50	0.1843	0.1161	9955	35	10	0.1031	0.0655
100	0.1666	0.1293	9997	3	0	0.0642	0.0548
200	0.1468	0.1297	9999	1	0	0.0461	0.0436
500	0.1316	0.125	10000	0	0	0.032	0.0313
1000	0.1187	0.1168	10000	0	0	0.0239	0.0237
2000	0.1049	0.1047	10000	0	0	0.0179	0.0179
10000	0.0790	0.0787	10000	0	0	0.0090	0.0090
20000	0.0761	0.0766	10000	0			0.0066
100000	0.0663	0.0663	10000	0	0	0.0032	0.0032
200000	0.0616	0.0617	10000	0	0	0.0023	0.0023

Note – The coverage rate of the confidence set is equal to $1 - (\text{Rejection probability})$. The results in this table pertains to GE_2 index with $SM_X(a_X = 2.8, q_X = 1.7)$ and $SM_Y(a_Y = 3.8, q_Y = 1.2855)$. $H_0: GE_2(X) - GE_2(Y) = 0.05401$.