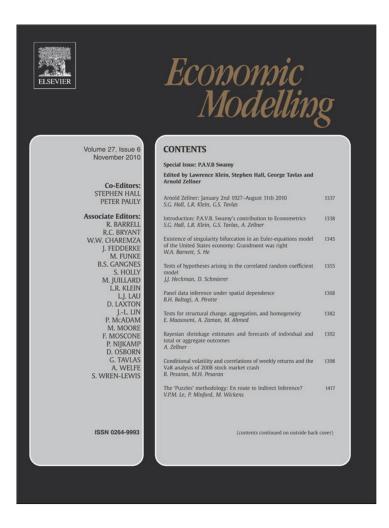
Provided for non-commercial research and education use. Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

http://www.elsevier.com/copyright

Economic Modelling 27 (2010) 1429-1435

Contents lists available at ScienceDirect



Economic Modelling

journal homepage: www.elsevier.com/locate/ecmod

Empirical likelihood confidence intervals for the Gini measure of income inequality

Yongsong Qin^a, J.N.K. Rao^{b,*}, Changbao Wu^c

^a Department of Mathematics, Guangxi Normal University, Guilin, Guanxi, 541004, China

^b School of Mathematics and Statistics, Carleton University, Ottawa, Ontario, Canada, K1S 5B6

^c Department of Statistics and Actuarial Science, University of Waterloo, Waterloo, Ontario, Canada, N2L 3G1

ARTICLE INFO

Keywords: Bootstrap percentile Bootstrap-t Confidence interval Coverage probability Gini coefficient Normal approximation

ABSTRACT

Gini coefficient is among the most popular and widely used measures of income inequality in economic studies, with various extensions and applications in finance and other related areas. This paper studies confidence intervals on the Gini coefficient for simple random samples, using normal approximation, bootstrap percentile, bootstrap-t and the empirical likelihood method. Through both theory and simulation studies it is shown that the intervals based on normal or bootstrap approximation are less satisfactory for samples of small or moderate size than the bootstrap-calibrated empirical likelihood ratio confidence intervals which perform well for all sample sizes. Results for stratified random sampling are also presented. © 2010 Elsevier B.V. All rights reserved.

1. Introduction

Income inequality has long been an active research area in economic studies. Among various measures of income inequality proposed in the statistical and economic literature, the Gini coefficient, *G*, is probably the most popular and widely used measure. It was originated from Gini's mean difference (Gini 1912, 1936), and is closely related to the Lorenz curve, the popular measure for the size distribution of income and wealth. Lorenz curves are also widely used in economic analysis (Kakwani, 1977).

Let $F(y) = P(Y \le y)$ be the cumulative distribution function of a nonnegative continuous random variable *Y*. We will refer to *Y* as the income variable. Let *X* and *Y* be two independent random variables following the same distribution F(y). The Gini mean difference is then defined as

$$D = E|X-Y| = \int_0^{+\infty} \int_0^{+\infty} |x-y| \, dF(x) \, dF(y).$$

The value of *D* is the average absolute difference of incomes of two randomly selected individuals and hence reflects the income inequality in the population. Noting that $0 \le D \le 2\mu$, where $\mu = E(Y) = \int_0^{+\infty} y dF(y)$ is the population mean income, the Gini coefficient, *G*, is defined as the normalized mean difference, i.e., $G = D/(2\mu) \in [0, 1]$, which can be equivalently written as (David, 1968)

$$G = \frac{1}{\mu} \int_0^{+\infty} \{2F(y) - 1\} y dF(y).$$
(1)

The Gini coefficient is also closely related to another popular measure of income inequality, the Lorenz curve (Lorenz, 1905;

E-mail address: jrao@math.carleton.ca (J.N.K. Rao).

Sendler, 1979). Let $F^{-1}(t) = \inf\{\xi: F(\xi) \ge t\}$ for $t \in [0, 1]$. The Lorenz curve based on the income distribution $F(\cdot)$ is then defined as

$$L(\alpha, F) = \frac{1}{\mu} \int_0^{\alpha} F^{-1}(t) dt = \frac{1}{\mu} \int_0^{F^{-1}(\alpha)} x dF(x)$$

for $\alpha \in [0, 1]$. The Gini coefficient *G* is equal to twice the area between a 45-degree line and the Lorenz curve, i.e., $G = 2\left\{0.5 - \int_0^1 L(\alpha, F) d\alpha\right\}$.

There exists an extensive literature on the Gini measure of income inequality. In addition to various applications and extensions in economic studies, statistical investigations focused largely on variance estimation; see, for instance, Glasser (1962), Sandström et al. (1985, 1988), Yitzhaki (1991), Karagiannis and Kovacevic (2000), among others. In particular, Yitzhaki (1991) calculated jackknife variance estimators of the plug-in moment estimator, \hat{G} , of G, under simple random sampling and stratified random sampling. However, confidence intervals for the Gini coefficient have not been studied by previous authors, with the exception of Sandström et al. (1988) where 95% normal approximation confidence intervals based on three variance estimators were briefly mentioned.

This paper presents confidence intervals on the Gini coefficient, *G*, using normal and bootstrap approximations and empirical likelihood (EL) based methods. We first consider the case of independent and identically distributed (*iid*) samples (or simple random samples when the sampling fraction is negligible), and then extend the results to stratified random sampling. In Section 2, we establish the asymptotic normality of the point estimator, \hat{G} , of *G* and construct confidence intervals on *G* based on the normal approximation. Confidence intervals on *G* based on the bootstrap percentile and the bootstrap-t methods are also given. In Section 3, the limiting distribution of the EL ratio statistic is established and the EL ratio confidence intervals are presented. A bootstrap-

^{*} Corresponding author. Tel.: +1 613 520 2600x2167.

^{0264-9993/\$ –} see front matter 0 2010 Elsevier B.V. All rights reserved. doi:10.1016/j.econmod.2010.07.015

Y. Qin et al. / Economic Modelling 27 (2010) 1429-1435

calibrated EL confidence interval on *G* is also presented. Results of a limited simulation study on the finite sample performance of the proposed confidence intervals are reported in Section 4. Extensions to stratified random sampling are outlined in Section 5. Proofs of theorems are relegated to Appendix A.

2. Normal and bootstrap approximation confidence intervals

Let $\{y_1, \dots, y_n\}$ be an *iid* sample from F(y) and $F_n(u) = n^{-1} \sum_{j=1}^n I(y_j \le u)$ be the empirical distribution function based on the sample, where $I(\cdot)$ denotes the indicator function. Noting that $G = E[\{2F(Y) - 1\}Y]/E(Y)$, a simple plug-in moment estimator of *G* is given by

$$\hat{G} = \frac{1}{\hat{\mu}} \cdot \frac{1}{n} \sum_{i=1}^{n} [\{2F_n(y_i) - 1\}y_i],$$
(2)

where $\hat{\mu} = \overline{y}$ is the sample mean. Let $h(u_1, u_2) = I(u_2 \le u_1)u_1 + I(u_1 \le u_2)u_2$. For $u_1 \ge 0$, let

$$h_1(u_1) = Eh(u_1, Y) = u_1 F(u_1) + \int_{u_1}^{\infty} y dF(y).$$
(3)

We have the following result on the asymptotic normality of \hat{G} .

Theorem 1. Suppose that $0 \le E(Y^2) \le \infty$. Then, as $n \to \infty$,

 $\sqrt{n}(\hat{G}-G) \xrightarrow{d} N(0,\sigma_1^2),$

where $\sigma_1^2 = \mu^{-2} Var\{2h_1(Y) - (G+1)Y\}$ and \xrightarrow{d} denotes convergence in distribution.

Using this result, a $(1 - \alpha)$ -level normal approximation confidence interval on *G* is given by

$$\left(\hat{G} - z_{\alpha/2} \frac{\hat{\sigma}_1}{\sqrt{n}}, \quad \hat{G} + z_{\alpha/2} \frac{\hat{\sigma}_1}{\sqrt{n}}\right),\tag{4}$$

where $z_{\alpha/2}$ is the upper $\alpha/2$ quantile from the standard normal distribution and

$$\hat{\sigma}_{1}^{2} = \frac{1}{\hat{\mu}^{2}} \cdot \frac{1}{n-1} \sum_{i=1}^{n} \left(u_{1i} - \overline{u}_{1} \right)^{2}$$
(5)

with

$$u_{1i} = 2\hat{h}_1(y_i) - (\hat{G} + 1)y_i, \quad \overline{u}_1 = \frac{1}{n}\sum_{i=1}^n u_{1i}$$
(6)

and

$$\hat{h}_1(u) = uF_n(u) + \frac{1}{n} \sum_{j=1}^n y_j I(y_j \ge u).$$
(7)

The symmetric interval (Eq. (4)) has asymptotically correct coverage probability for large samples. For small samples, however, the normal interval (Eq. (4)) tends to have under-coverage problems, as observed from the simulation results reported in Section 4. In addition, the tail error rates of this interval also tend to be unbalanced, due to skewness of income distributions.

The normal approximation can be replaced by bootstrap procedures. A $(1-\alpha)$ -level confidence interval based on the bootstrap percentile of $\hat{G}-G$ is given by

$$\left(\hat{G} - P_{1-\alpha/2}, \quad \hat{G} - P_{\alpha/2}\right),\tag{8}$$

where P_{α} is the 100 α th percentile of the sampling distribution of $\hat{G}^* - \hat{G}$, and \hat{G}^* is the estimator of *G* calculated based on a bootstrap sample $\{y_1^*, \dots, y_n^*\}$ taken from the original sample $\{y_1, \dots, y_n\}$ by simple random sampling with replacement. The percentile P_{α} can be obtained through Monte Carlo approximations by drawing a large number of bootstrap samples. Let $\hat{G}^*(b)$ be the estimate of *G* computed from the *b*th bootstrap sample $\{y_1^*(b), \dots, y_n^*(b)\}, b = 1, \dots, B$ and let $\hat{G}^*[1] \leq \dots \leq \hat{G}^*[B]$ be the ordered sequence of the $\hat{G}^*(b)$ s. Then $P_{\alpha} \doteq \hat{G}^*[\alpha B] - \hat{G}$.

The bootstrap-t confidence interval on G is constructed as

$$\left(\hat{G} - T_{1-\alpha/2}\frac{\hat{\sigma}_1}{\sqrt{n}}, \quad \hat{G} - T_{\alpha/2}\frac{\hat{\sigma}_1}{\sqrt{n}}\right),\tag{9}$$

where T_{α} is the 100 α th percentile of the sampling distribution of $(\hat{G}^* - \hat{G}) / (\hat{\sigma}_1^* / \sqrt{n})$, and \hat{G}^* and $\hat{\sigma}_1^* / \sqrt{n}$ are the estimator of *G* and the associated standard error based on a bootstrap sample $\{y_1^*, ..., y_n^*\}$. Once again, T_{α} can be obtained through Monte Carlo approximations.

3. Empirical likelihood ratio confidence intervals

The empirical likelihood (EL) method is a nonparametric approach and is particularly suitable to handle inferential problems involving skewed distributions. EL confidence intervals, obtained from profiling the empirical likelihood ratio statistic, are range respecting and transformation invariant. The shape and orientation of the EL intervals are determined by the data (Owen, 2001), unlike the normal approximation and bootstrap intervals. The log-EL ratio statistic for $\theta = G$ is given by

$$R(\theta) = \sum_{i=1}^{n} \log\{n\tilde{p}_i(\theta)\},\tag{10}$$

where $\tilde{p}_1(\theta), \cdots, \tilde{p}_n(\theta)$ maximize the log-EL function $l(\mathbf{p}) = \sum_{i=1}^n log(p_i)$ subject to the following set of constraints:

$$p_i > 0$$
, $\sum_{i=1}^{n} p_i = 1$ and $\sum_{i=1}^{n} p_i [\{2F_n(y_i) - 1\}y_i - \theta y_i] = 0.$ (11)

The last constraint in Eq. (11) is induced by the estimating equation $E[\{2F(Y) - 1\}Y - \theta Y] = 0$ which defines the parameter $\theta = G$, with the unknown distribution function $F(\cdot)$ replaced by the empirical distribution function $F_n(\cdot)$.

Let $Z(y_i, \theta) = \{2F_n(y_i) - 1\}y_i - \theta y_i, i = 1, \dots, n$. It can be shown, by using the Lagrange multiplier method, that

$$R(\theta) = -\sum_{i=1}^{n} \log\{1 + \lambda Z(y_i, \theta)\},\$$

where λ is the solution to the equation

$$\frac{1}{n}\sum_{i=1}^{n} \frac{Z(y_i,\theta)}{1+\lambda Z(y_i,\theta)} = 0.$$

Theorem 2 establishes the asymptotic distribution of the log-EL ratio statistic $R(\theta)$.

Theorem 2. Suppose that $0 < E(Y^3) < \infty$. Then, as $n \to \infty$,

$$-2R(\theta) \xrightarrow{d} \frac{\sigma_3^2}{\sigma_2^2} \chi^2(1),$$

where $\sigma_2^2 = Var\{2YF(Y) - (\theta + 1)Y\}$, $\sigma_3^2 = Var\{2h_1(Y) - (\theta + 1)Y\}$, and $h_1(\cdot)$ is defined in Eq. (3).

Using this result, a $(1-\alpha)$ -level EL ratio confidence interval on G can be constructed as

$$\Big\{\theta|-2R(\theta) \leq \hat{k}^{-1}\chi_{\alpha}^{2}(1)\Big\},\tag{12}$$

where $\chi^2_{\alpha}(1)$ is the upper α quantile of the χ^2 distribution with one degree of freedom. The scaling factor, \hat{k} , in Eq. (12) is given by $\hat{k} = \hat{\sigma}_2^2 / \hat{\sigma}_3^2$, where

$$\hat{\sigma}_{2}^{2} = \frac{1}{n-1} \sum_{i=1}^{n} \left(u_{2i} - \overline{u}_{2} \right)^{2}$$
(13)

with

$$u_{2i} = 2y_i F_n(y_i) - (\hat{G} + 1)y_i, \quad \overline{u}_2 = \frac{1}{n} \sum_{j=1}^n u_{2j}$$
(14)

and

$$\hat{\sigma}_3^2 = \frac{1}{n-1} \sum_{i=1}^n (u_{1i} - \overline{u}_1)^2, \tag{15}$$

where u_{1i} and \overline{u}_1 are defined in Eq. (6). Note that \hat{k} is a consistent estimator of $k = \sigma_2^2/\sigma_3^2$ as $n \to \infty$.

The EL ratio confidence interval (Eq. (12)) has asymptotically correct $(1 - \alpha)$ -level coverage probability for large samples. The EL-based interval, however, often has under-coverage problems when the sample size *n* is small or moderate. This has been observed in many other applications of the EL method (Owen, 2001). The following bootstrap-calibrated EL interval is an attractive alternative under such scenarios.

Let $\{y_1^*, \dots, y_n^*\}$ be a bootstrap sample selected from the original sample $\{y_1, \dots, y_n\}$ using simple random sampling with replacement. Also, let $R_1^*(\hat{G})$ be the value of $R(\theta)$ calculated from the bootstrap sample, using $\theta = \hat{G}$. Repeat the process independently for a large number of times, B, to get $R_1^*(\hat{G}), \dots, R_B^*(\hat{G})$. Let C_α be the upper 1000% sample quantile of the B values $-R_1^*(\hat{G}), \dots, -R_B^*(\hat{G})$. The $(1 - \alpha)$ -level bootstrap-calibrated EL ratio interval on G can then be constructed as

$$\{\theta \mid -R(\theta) \le C_{\alpha}\}. \tag{16}$$

Another major advantage of using the bootstrap-calibrated EL interval (Eq. (16)) is that the scale factor $\hat{k} = \hat{\sigma}_2^2 / \hat{\sigma}_3^2$, which is required in the EL ratio interval (Eq. (12)), is not needed here in the construction of the interval. It can be shown, by following the same lines of the proof of Theorem 2, that the bootstrap version of the EL ratio function $R_1^*(\hat{G})$ converges to the same scaled $\chi^2(1)$ distribution as $n \to \infty$, and hence the two intervals (Eqs. (12) and (16)) have the same asymptotic coverage probabilities under large samples. The bootstrap-calibrated interval (Eq. (16)), however, performs better when sample sizes are small or moderate, as shown in the simulation study reported in Section 4.

4. A simulation study

We examined the finite sample performances of five confidence intervals for the Gini coefficient *G* through a simulation study: (i) the normal approximation interval (Eq. (4)), denoted by NA; (ii) the bootstrap percentile interval (Eq. (8)), denoted by BTp; (iii) the bootstrap-t interval (Eq. (9)), denoted by BT; (iv) the EL ratio interval (Eq. (12)), denoted by EL1, based on the scaled χ^2 approximation; and (v) the EL ratio interval (Eq. (16)), denoted by EL2, using the bootstrap calibration method. Four different population distributions were considered: (i) the χ^2 distribution with one degree of freedom ($\chi^2(1)$); (ii) the standard exponential distribution (*Exp*(1)); and (iv) the standard lognormal distribution (*LN*(0, 1)). The population distributions considered here represent potential income distributions one might encounter in real-world situations. The true value of the Gini coefficient *G* for *Exp*(1) is 0.5 and the true values of *G* for $\chi^2(1), \chi^2(3)$ and LN(0,1) are approximately 0.6366, 0.4244 and 0.5205, respectively, obtained through Monte Carlo simulations.

Confidence intervals on *G* were evaluated in terms of the simulated coverage probability (CP), lower (L) and upper (U) tail error rates and the average length (AL). For each simulated sample, we constructed confidence intervals on *G* using NA, BTp, BTt, EL1 and EL2, with B = 2000 for the bootstrap procedures. The simulation process was repeated R = 2000 times for each of the four population distributions and selected sample sizes ranging from n = 20 to n = 80.

Tables 1 and 2 report the simulation results for the 95% confidence intervals on G. Major observations from the simulation can be summarized as follows: (i) The normal approximation interval (NA), the bootstrap percentile interval (BTp) and the EL interval based on the χ^2 approximation (EL1) have very similar performances and none of them seems to be satisfactory when $n \leq 60$. (ii) The bootstrap-calibrated EL interval (EL2) has coverage probabilities very close to the nominal value for all sample sizes considered when the population distribution is $\chi^2(1)$, $\chi^2(3)$ or Exp(1). (iii) The EL2 intervals also demonstrate balanced tail error rates for the $\chi^2(1)$ distribution and to a lesser extent for the $\chi^2(3)$ and Exp(1)distributions. (iv) The bootstrap-t interval (BTt) has coverage probabilities comparable to EL2 when the population distribution is $\chi^2(1)$ or Exp(1) but the BTt intervals are wider in those cases. For instance, the lengths of the BTt and EL2 intervals are respectively 0.297 and 0.269 for n = 20 and the $\chi^2(1)$ distribution. The BTt and EL2 intervals are similar in length for the $\chi^2(3)$ distribution but the coverage probability of the BTt interval is not as good as the EL2 interval when n = 20. (v) None of the methods provides very good results for the lognormal distribution, but the bootstrap-calibrated EL2 intervals have marginally acceptable results regardless of the sample size: coverage probabilities around 92% for EL2 compared to 86-90% for NA, BTp, BTt and EL1. Moreover, in the lognormal case, EL2 outperforms BTt in terms of AL and yet gives coverage probabilities closer to nominal 95% than BTt. For example, for n = 20 Table 2 gives 92.8% and 0.281 as CP and AL for EL2 compared to 86.2% and 0.350 for BTt.

In a recent paper Giorgi et al. (2006) reported results from a simulation study that the bootstrap-t confidence intervals for the so-called S-Gini and E-Gini indices have superb coverage probabilities for all cases considered in their paper. The superiority of the bootstrap-t interval, however, does not seem to show up for the scenarios examined here on the Gini coefficient *G*. The bootstrap-calibrated

Table 1 Simulation results for 95% confidence intervals on *G*: $\chi^2(1)$ and $\chi^2(3)$.

п	CI	L	СР	U	AL	L	СР	U	AL	
		$\chi^{2}(1)$	$\chi^{2}(1)$				$\chi^{2}(3)$			
20	NA	8.4	88.9	2.7	0.244	9.9	87.9	2.2	0.208	
	ВТр	7.4	86.5	6.1	0.243	4.7	88.4	6.9	0.204	
	BTt	2.3	93.5	4.2	0.297	3.3	91.9	4.8	0.235	
	EL1	6.9	89.9	3.2	0.233	9.7	87.8	2.5	0.201	
	EL2	3.3	94.8	1.9	0.269	4.9	94.2	0.9	0.240	
40	NA	5.7	91.5	2.8	0.183	6.2	91.5	2.3	0.154	
	ВТр	5.0	90.7	4.3	0.181	3.8	90.6	5.6	0.152	
	BTt	2.4	94.0	3.6	0.207	2.9	93.5	3.6	0.167	
	EL1	4.8	92.1	3.1	0.179	6.3	91.3	2.4	0.152	
	EL2	3.2	94.4	2.4	0.194	4.7	93.6	1.7	0.166	
60	NA	5.2	92.4	2.4	0.151	4.9	92.5	2.6	0.128	
	ВТр	4.8	91.5	3.7	0.150	2.9	92.3	4.8	0.126	
	BTt	2.6	94.3	3.1	0.166	2.5	94.1	3.4	0.135	
	EL1	4.7	92.6	2.7	0.149	4.8	92.5	2.7	0.127	
	EL2	3.7	94.1	2.2	0.157	3.9	94.1	2.0	0.134	
80	NA	3.5	93.8	2.7	0.132	4.0	94.5	1.5	0.112	
	ВТр	3.4	93.2	3.4	0.132	2.7	93.9	3.4	0.110	
	BTt	1.6	95.3	3.1	0.143	2.5	95.4	2.1	0.117	
	EL1	3.4	93.6	3.0	0.130	3.8	94.7	1.5	0.111	
	EL2	3.0	94.4	2.6	0.136	3.2	95.5	1.3	0.115	

1432

Y. Qin et al. / Economic Modelling 27 (2010) 1429-1435

 Table 2

 Simulation results for 95% confidence intervals on G: Exp(1) and LN (0,1).

п	CI	L	СР	U	AL	L	СР	U	AL	
		Exp(1	Exp(1)				LN (0,1)			
20	NA	7.4	89.8	2.8	0.228	6.4	85.6	8.0	0.248	
	BTp	4.4	88.6	7.0	0.224	6.8	78.9	14.3	0.243	
	BTt	2.5	93.2	4.3	0.265	2.7	86.2	11.1	0.350	
	EL1	6.6	90.3	3.1	0.220	5.3	86.0	8.7	0.237	
	EL2	3.6	95.2	1.2	0.256	1.9	92.8	5.3	0.281	
40	NA	6.0	91.0	3.0	0.170	2.9	88.3	8.8	0.199	
	ВТр	4.5	89.7	5.8	0.167	4.8	82.5	12.7	0.193	
	BTt	2.9	93.5	3.6	0.186	1.0	89.4	9.6	0.271	
	EL1	5.7	91.1	3.2	0.167	2.7	87.8	9.5	0.193	
	EL2	3.9	94.0	2.1	0.179	1.3	91.2	7.5	0.216	
60	NA	5.4	92.7	1.9	0.141	2.4	90.7	6.9	0.171	
	BTp	3.9	92.3	3.8	0.139	4.1	86.6	9.3	0.167	
	BTt	2.8	94.4	2.8	0.151	1.2	92.0	6.8	0.223	
	EL1	5.0	93.0	2.0	0.139	2.2	90.6	7.2	0.168	
	EL2	3.8	94.4	1.8	0.146	1.3	92.2	6.5	0.183	
80	NA	3.9	93.9	2.2	0.123	1.8	91.4	6.8	0.155	
	ВТр	3.3	92.8	3.9	0.122	3.2	87.9	8.9	0.151	
	BTt	2.5	94.2	3.3	0.130	1.0	92.1	6.9	0.197	
	EL1	3.8	93.6	2.6	0.122	1.8	91.4	6.8	0.152	
	EL2	3.4	94.6	2.0	0.126	1.2	92.2	6.6	0.164	

empirical likelihood ratio confidence interval EL2 seems to be the best among the five intervals included in our study.

5. Stratified random sampling

Suppose that the population is divided into *S* strata with known stratum sizes N_1, \dots, N_S . Let $W_s = N_s/N$, $s = 1, \dots, S$, where $N = \sum_{s=1}^{S} N_s$ is the overall population size. Independent simple random samples of sizes n_s , $s = 1, \dots, S$ are drawn from the strata, and the strata sampling fractions, n_s/N_s , are assumed to be negligible. Hence, we regard the sample in stratum *s*, $\{y_{si}, i = 1, \dots, n_s\}$, as an *iid* sample generated by the continuous random variable Y_s with distribution function $F_s(y) = P(Y_s \le y)$. The overall population mean and distribution function are respectively given by

$$\mu = \sum_{s=1}^{S} W_{s}\mu_{s}$$
 and $F(y) = \sum_{s=1}^{S} W_{s}F_{s}(y)$

where $\mu_s = E(Y_s)$. The estimators of μ and F(y) are respectively given by

$$\hat{\mu} = \sum_{s=1}^{S} W_s \hat{\mu}_s$$
 and $F_n(y) = \sum_{s=1}^{S} W_s F_{ns}(y)$, (17)

where $\hat{\mu}_s = n_s^{-1} \sum_{i=1}^{n_s} y_{si}$ and $F_{ns}(y) = n_s^{-1} \sum_{i=1}^{n_s} I(y_{si} \le y)$. Finally, the plug-in moment estimator of the Gini coefficient *G* based on the stratified sample is given by

$$\hat{G}_{st} = \frac{1}{\hat{\mu}} \sum_{s=1}^{S} W_s \frac{1}{n_s} \sum_{i=1}^{n_s} \{2F_n(y_{si}) - 1\} y_{si}.$$

5.1. Normal approximation confidence intervals

Let $n = \sum_{s=1}^{S} n_s$ be the overall sample size. We have the following result on the asymptotic normality of \hat{G}_{st} .

Theorem 3. Suppose that $0 \le E(Y_s^2) \le \infty$, $n_s/n \to \tau_s$ ($0 \le \tau_s \le 1$), s = 1 - S. Then, as $n \to \infty$,

 $\sqrt{n}\left(\hat{G}_{st}\!-\!G\right) \xrightarrow{d} N\!\left(0,\sigma_{a1}^2\right),$

where

$$\sigma_{a1}^{2} = \frac{1}{\mu^{2}} \sum_{s=1}^{S} \tau_{s}^{-1} W_{s}^{2} Var \Big\{ 2Y_{s} F(Y_{s}) + 2 \int_{Y_{s}}^{\infty} y dF(y) - (G+1)Y_{s} \Big\}.$$

Using this result, a normal approximation confidence interval for *G* with asymptotically correct coverage probability $1 - \alpha$ is given by

$$\left(\hat{G}_{st} - z_{\alpha/2} \frac{\hat{\sigma}_{a1}}{\sqrt{n}}, \quad \hat{G}_{st} + z_{\alpha/2} \frac{\hat{\sigma}_{a1}}{\sqrt{n}}\right), \tag{18}$$

where $z_{\alpha/2}$ is the upper $\alpha/2$ quantile from the standard normal distribution and

$$\hat{\sigma}_{a1}^{2} = \frac{1}{\hat{\mu}^{2}} \sum_{s=1}^{S} \frac{n}{n_{s}} W_{s}^{2} \cdot \frac{1}{n_{s}-1} \sum_{i=1}^{n_{s}} \left(u_{1si} - \overline{u}_{1s} \right)^{2}$$
(19)

with

$$u_{1si} = 2\hat{h}_{a1}(y_{si}) - (\hat{G}_{st} + 1)y_{si}, \quad \overline{u}_{1s} = \frac{1}{n_s} \sum_{i=1}^{n_s} u_{1si}$$
(20)

and

$$\hat{h}_{a1}(u) = uF_n(u) + \sum_{s=1}^{S} W_s \frac{1}{n_s} \sum_{j=1}^{n_s} y_{sj} I(y_{sj} \ge u).$$
(21)

5.2. Empirical likelihood ratio confidence intervals

Under stratified random sampling, the log-EL ratio statistic for $\theta = G$ is given by

$$R(\theta) = \sum_{s=1}^{S} \sum_{i=1}^{n_s} \log\{n_s \tilde{p}_{si}(\theta)\},$$
(22)

where $\{\tilde{p}_{si}(\theta), i = 1, \dots, n_s, s = 1, \dots, S\}$ maximize the log-EL function $l(\mathbf{p}) = \sum_{s=1}^{S} \sum_{i=1}^{n_s} log(n_s p_{si})$ subject to the following set of constraints:

$$p_{si} > 0, \sum_{i=1}^{n_s} p_{si} = 1, s = 1, \dots, S \text{ and } \sum_{s=1}^{S} W_s \sum_{i=1}^{n_s} p_{si} [\{2F_n(y_{si}) - 1\}y_{si} - \theta y_{si}] = 0.$$

(23)

Let $Z(y_{si}, \theta) = \{2F_n(y_{si})-1\}y_{si}-\theta y_{si}, i = 1, \cdots, n_s, s = 1, \cdots, S$. It can be shown, by using the Lagrange multiplier method, that

$$p_{si} = \frac{1}{n_s} \frac{1}{1 + m_s \lambda Z^*(y_{si}, \theta)},$$

where $m_s = nW_s n_s^{-1}$, $Z^*(y_{si}, \theta) = Z(y_{si}, \theta) - \sum_{i=1}^{n_s} p_{si}Z(y_{si}, \theta)$, and $\lambda = \lambda(\theta)$ is the solution of the equation

$$\sum_{s=1}^{S} W_{s} \frac{1}{n_{s}} \sum_{i=1}^{n_{s}} \frac{Z(y_{si}, \theta)}{1 + m_{s} \lambda Z^{*}(y_{si}, \theta)} = 0.$$

Thus,

$$R(\theta) = -\sum_{s=1}^{S} \sum_{i=1}^{n_s} \log\{1 + m_s \lambda Z^*(y_{si}, \theta)\}.$$

Algorithms for evaluating the log-EL ratio function under stratified random sampling and other sampling designs are given by Zhong and Rao (2000) and Wu (2004).

Theorem 4 below establishes the asymptotic distribution of the log-EL ratio statistic $R(\theta)$.

Theorem 4. Suppose that $E(Y_s^3) < \infty$, $n_s/n \rightarrow \tau_s(0 < \tau_s < 1)$, $s = 1, \dots, S$. Then, as $n \to \infty$,

$$-2R(\theta) \xrightarrow{d} \frac{\mu^2 \sigma_{a1}^2}{\sigma_{a2}^2} \chi^2(1),$$

where σ_{a1}^2 is defined in Theorem 3 and

$$\sigma_{a2}^{2} = \sum_{s=1}^{S} W_{s}^{2} \tau_{s}^{-1} E\{2Y_{s}F(Y_{s}) - (\theta + 1)Y_{s}\}^{2}.$$

Using the result in Theorem 4, a $(1 - \alpha)$ -level confidence interval on G with asymptotically correct coverage probability can be constructed as

$$\left\{\theta \mid -2R(\theta) \le \hat{k}_a^{-1} \chi_\alpha^2(1)\right\},\tag{24}$$

where $\chi^2_{\alpha}(1)$ is the upper α quantile of the χ^2 distribution with one degree of freedom. The scaling factor \hat{k}_a in Eq. (24) is given by $\hat{k}_a = \hat{\sigma}_{a2}^2 / (\hat{\mu}^2 \hat{\sigma}_{a1}^2)$, where $\hat{\mu}$ is given by Eq. (17), $\hat{\sigma}_{a1}^2$ is given by Eq. (19) and

$$\hat{\sigma}_{a2}^{2} = \sum_{s=1}^{S} W_{s}^{2} \frac{n}{n_{s}^{2}} \sum_{i=1}^{n_{s}} \{2y_{si}F_{n}(y_{si}) - (\hat{G}_{st} + 1)y_{si}\}^{2}$$

Note that \hat{k}_a is a consistent estimator of $k_a = \sigma_{a2}^2 / (\mu^2 \sigma_{a1}^2)$. A bootstrap-calibrated EL ratio confidence interval can also be constructed along the lines of Eq. (16) by drawing independent bootstrap samples $\{y_{si}, i = 1, \dots, n_s\}$ from $\{y_{si}, i = 1, \dots, n_s\}$, $s = 1, \dots, S$ by simple random sampling with replacement.

6. Additional remarks

The Gini coefficient is defined as $G = E[\{2F(Y) - 1\}Y]/E(Y)$, which is equivalently defined by the estimating equation

$$E[\{2F(Y)-1\}Y-GY] = 0.$$
(25)

This is not an easy-to-handle estimating equation since it also involves the unknown distribution function, $F(\cdot)$.

The maximum EL estimator of G may be defined as the maximizer of the log-EL function $l(\theta) = \sum_{i=1}^{n} \log\{p_i(\theta)\}$, where for a fixed θ the $p_i(\theta)$ maximize $\sum_{i=1}^{n} \log(p_i)$ subject to $\sum_{i=1}^{n} p_i = 1$ and

$$\sum_{i=1}^{n} p_i \Big[\{ 2\hat{F}(y_i) - 1 \} y_i - \theta y_i \Big] = 0,$$
(26)

$$\hat{F}(y_i) = \sum_{j=1}^{n} p_i I(y_j \le y_i), \ i = 1, \cdots, n.$$
(27)

A computational difficulty arises since Eqs. (26) and (27) together define a non-linear constraint on the p_i . For any fixed θ , however, the solution $p_i(\theta)$ can be obtained through an iterative procedure as follows:

- (i) Start with $p_i^{(0)} = 1/n$ and let $\hat{F}(y_i) = \sum_{j=1}^n p_i^{(0)} I(y_j \le y_i)$; (ii) Find $p_i^{(1)}$ which maximize $\sum_{i=1}^n log(p_i)$ subject to $\sum_{i=1}^n p_i = 1$ and Eq. (26) only;
- (iii) Let $p_i^{(0)} = p_i^{(1)}$ and iterate between (i) and (ii) until convergence.

While theoretical properties of the maximum EL estimator $\hat{\theta}$ of *G*, which is the maximizer of $l(\theta)$, are not immediately clear, simulation results (not reported here) showed that $\hat{\theta}$ is virtually identical to the plug-in moment estimator, \hat{G} . The direct use of $F_n(y_i)$ in the constraint (11) has a major advantage of computational simplicity compared to the use of $\hat{F}(y_i)$ in the constraint (26).

Empirical likelihood confidence intervals on the Gini mean difference D = E[X - Y] were discussed by Wood et al. (1996) and Jing et al. (2008) using results from U-statistics. Those results, however, do not apply to the Gini coefficient G which is a ratio of D and 2μ . The sample version of *D* is a *U*-statistic with the kernel h(x, y) = |x - y| but the sample version of $\mu = E(X)$ is a *U*-statistic with kernel h(x) = x. Therefore, it seems that none of the existing approaches on empirical likelihood methods for U-statistics can simultaneously handle two Ustatistics with kernels of different degrees. Moreover, the U-statistic approach does not readily extend to stratified sampling, even for the Gini mean difference *D*, unlike the approach proposed in this paper.

Most income distributions are heavily skewed, and the bootstrapcalibrated empirical likelihood method is a very attractive approach to interval estimation under such scenarios, as demonstrated by the simulation results reported in Section 4 for the Gini coefficient G. Empirical likelihood-based interval estimation from complex survey samples, such as data from stratified multi-stage sampling, is currently under investigation, using the pseudo empirical likelihood approach of Wu and Rao (2006). Empirical likelihood-based interval estimation for other measures of income distributions, such as low income proportions, Lorenz curve ordinate and quantile share, is also under investigation.

Acknowledgements

This research was supported by grants from the Natural Sciences and Engineering Research Council of Canada awarded to Rao and Wu, and grant 10971038 from the National Natural Science Foundation of China awarded to Qin.

Appendix A. Proofs

Proof of Theorem 1. The proof given here is based on standard results for von Mises statistics, the associated U-statistics and the projections of U-statistics (Serfling, 1980, Chapter 5). It is necessary, however, to spell out the details since we will need to extend the approach to prove Theorem 3 for stratified random sampling.

Let
$$\mu_1 = E \{ 2h_1(Y) - Y \}, \quad X_{n1} = n^{-1} \sum_{i=1}^n \{ 2h_1(y_i) - y_i - \mu_1 \},$$

 $X_{n2} = n^{-1} \sum_{i=1}^n (y_i - \mu),$ and
 $V_n = n^{-2} \sum_{i=1}^n \sum_{j=1}^n 2I (y_j \le y_i) y_i$
 $= n^{-2} \sum_{i=1}^n \sum_{j=1}^n \{ I (y_j \le y_i) y_i + I (y_i \le y_j) y_j \}.$

Noting that V_n is a von Mises statistic, we denote its corresponding *U*-statistic as U_n . When $E(Y^2) < \infty$, it is well-known (Serfling, 1980, p. 206) that $\sqrt{n}(V_n - U_n) = o_p(1)$, where $o_p(1)$ denotes a term that goes to zero in probability as $n \rightarrow \infty$. The projection of U_n is

$$\hat{U}_n = Eh_1(Y) + \frac{2}{n} \sum_{i=1}^n \{h_1(y_i) - Eh_1(Y)\},\$$

where $h_1(\cdot)$ is defined in Eq. (3) and $Eh_1(Y) = 2E\{YF(Y)\}$. It follows from Serfling (1980, p. 190) that $\sqrt{n}(\hat{U}_n - U_n) = o_p(1)$. Thus

$$\begin{split} \hat{G} &= \frac{V_n - \hat{\mu}}{\hat{\mu}} = \frac{\hat{U}_n - \hat{\mu}}{\hat{\mu}} + o_p \left(n^{-1/2} \right) \\ &= \frac{n^{-1} \sum_{i=1}^{n} \{2h_1(y_i) - y_i - Eh_1(Y)\}}{\hat{\mu}} + o_p \left(n^{-1/2} \right) \\ &= \frac{X_{n1} + \mu_1 - Eh_1(Y)}{X_{n2} + \mu} + o_p \left(n^{-1/2} \right) \\ &= \frac{X_{n1} + 2E\{YF(Y)\} - \mu}{X_{n2} + \mu} + o_p \left(n^{-1/2} \right) \\ &= \frac{X_{n1} + G\mu}{X_{n2} + \mu} + o_p \left(n^{-1/2} \right), \end{split}$$
(28)

Y. Qin et al. / Economic Modelling 27 (2010) 1429-1435

where $o_p(n^{-1/2})$ denotes a term of lower order than $n^{-1/2}$ in probability. Here we have used the facts that $Eh_1(Y) = 2E\{YF(Y)\}$ and $E[\{2F(Y) - 1\}Y - GY] = 0$. Note that $\sqrt{n}X_{nj} = O_p(1), j = 1, 2$, where $O_p(1)$ denotes a term bounded in probability. A standard Taylor series expansion leads to $1/(X_{n2} + \mu) = \mu^{-1} - \mu^{-2}X_{n2} + O_p(n^{-1})$, where $O_p(n^{-1})$ denotes a term of order of n^{-1} in probability. Substituting this into Eq. (28), we have

$$\hat{G} = G + \mu^{-1}(X_{n1} - GX_{n2}) + o_p \left(n^{-1/2}\right).$$
⁽²⁹⁾

The result of Theorem 1 then follows from Eq. (29) and the central limit theorem (CLT) applied to $X_{n1} - GX_{n2}$.

Proof of Theorem 2. It can be shown that

$$\max_{1 \le i \le n} |Z(y_i, \theta)| = o_p\left(n^{1/2}\right) \tag{30}$$

and that

$$\frac{1}{n}\sum_{i=1}^{n} Z^{2}(y_{i},\theta) = \frac{1}{n}\sum_{i=1}^{n} \left[\{2F_{n}(y_{i})-1\}y_{i}-\theta y_{i}\right]^{2} + o_{p}(1) = \sigma_{2}^{2} + o_{p}(1),$$
(31)

where $\sigma_2^2 = E\{2YF(Y) - (\theta + 1)Y\}^2 = Var\{2YF(Y) - (\theta + 1)Y\}$. From the proof of Theorem 1, we have

$$n^{-1/2} \sum_{i=1}^{n} Z(y_i, \theta) = n^{-1/2} \sum_{i=1}^{n} \{2h_1(y_i) - (\theta + 1)y_i - Eh_1(Y)\} + o_p(1).$$

On the other hand, $Eh_1(Y) = E\{2YF(Y)\}$ and $E[\{2F(Y) - 1\}Y - \theta Y] = 0$ imply that $E\{2h_1(Y) - (\theta + 1)Y - Eh_1(Y)\} = 0$. Therefore, by the CLT, we have

$$n^{-1/2} \sum_{i=1}^{n} Z(y_i, \theta) \xrightarrow{d} N(0, \sigma_3^2),$$
(32)

where $\sigma_3^2 = Var\{2h_1(Y) - (\theta + 1)Y\}$. From Eqs. (30), (31), (32), and the proof of Theorem 1 in Owen (1990), it can be shown that

$$-2R(\theta) = \left\{\frac{1}{n}\sum_{i=1}^{n} Z^{2}(y_{i},\theta)\right\}^{-1} \left\{n^{-1/2}\sum_{i=1}^{n} Z(y_{i},\theta)\right\}^{2} + o_{p}(1) \stackrel{d}{\to} \frac{\sigma_{3}^{2}}{\sigma_{2}^{2}}\chi^{2}(1).$$

Proof of Theorem 3. Let

$$\begin{split} H(u,v) &= uI(v \le u), \\ \theta_{st} &= EH(Y_s,Y_t), \\ h_{st10}(u) &= EH(u,Y_t) - \theta_{st}, \\ h_{st01}(u) &= EH(Y_s,u) - \theta_{st}, \\ h_{s1}(u) &= E\{H(u,Y_s) + H(Y_s,u)\} = uF_s(u) + \int_u^{\infty} y dF_s(y), \end{split}$$

and

$$\begin{split} V_n &= 2\sum_{s=1}^{S} \sum_{i=1}^{n_s} W_s n_s^{-1} y_{si} F_n(y_{si}) \\ &= 2\sum_{s=1}^{S} \sum_{t=1}^{S} W_s W_t n_s^{-1} n_t^{-1} \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} y_{si} I(y_{tj} \leq y_{si}) \\ &= 2\sum_{s=1}^{S} W_s^2 n_s^{-2} \sum_{i=1}^{n_s} \sum_{y=1}^{n_s} y_{si} I(y_{sj} \leq y_{si}) \\ &+ 2\sum_{s=1}^{S} \sum_{t \neq s} W_s W_t n_s^{-1} n_t^{-1} \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} y_{si} I(y_{tj} \leq y_{si}) \\ &= 2\sum_{s=1}^{S} W_s^2 V_{ns} + 2\sum_{s=1}^{S} \sum_{t \neq s} W_s W_t V_{nst}, \end{split}$$

where $V_{ns} = n_s^{-2} \sum_{i=1}^{n_s} \sum_{j=1}^{n_s} y_{si} I(y_{sj} \le y_{si})$ and $V_{nst} = n_s^{-1} n_t^{-1} \sum_{i=1}^{n_s} \sum_{j=1}^{n_t} y_{si} I(y_{tj} \le y_{si})$.

Noting that V_{ns} is a von Mises statistic and $n_s/n \rightarrow \tau_s$ (0< τ_s <1), it follows from the proof of Theorem 1 that

$$2V_{ns} = Eh_{s1}(Y_s) + \frac{2}{n_s} \sum_{i=1}^{n_s} \{h_{s1}(y_{si}) - Eh_{s1}(Y_s)\} + o_p(n^{-1/2}).$$

When $s \neq t$, V_{nst} is a two-sample *U*-statistic. From Theorem 12.6 in van der Vaart (1998), we have

$$\begin{split} V_{nst} &= \theta_{st} + \frac{1}{n_s} \sum_{i=1}^{n_s} h_{st10}(y_{si}) + \frac{1}{n_t} \sum_{i=1}^{n_t} h_{st01}(y_{ti}) + o_p(n^{-1/2}) \\ &= \theta_{st} + \frac{1}{n_s} \sum_{i=1}^{n_s} \left[y_{si} F_t(y_{si}) - E\left\{ Y_s F_t(Y_s) \right\} \right] \\ &+ \frac{1}{n_t} \sum_{i=1}^{n_t} \left[\int_{y_{ti}}^{\infty} y dF_s(y) - E\left\{ \int_{Y_t}^{\infty} y dF_s(y) \right\} \right] + o_p(n^{-1/2}). \end{split}$$

It follows that

$$V_{n} = 2\sum_{s=1}^{S} \sum_{t=1}^{S} W_{s}W_{t}\theta_{st} + 2\sum_{s=1}^{S} W_{s}\frac{1}{n_{s}}\sum_{i=1}^{n_{s}} \left[y_{si}F(y_{si}) - E\{Y_{s}F(Y_{s})\} + \int_{y_{si}}^{\infty} ydF(y) - E\{\int_{Y_{s}}^{\infty} ydF(y)\} \right] + o_{p}(n^{-1/2}).$$

Further, from $F(y) = \sum_{s=1}^{S} W_s F_s(y)$ and the definition of *G* in Eq. (1), we have

$$2\sum_{s=1}^{S}\sum_{t=1}^{S}W_{s}W_{t}\theta_{st}-\mu=G\mu$$

Thus

$$\hat{G}_{st} = \frac{V_n - \hat{\mu}}{\hat{\mu}} = \frac{X_{n1} + G\mu}{X_{n2} + \mu} + o_p \left(n^{-1/2} \right),$$
(33)

where

$$X_{n1} = 2 \sum_{s=1}^{S} W_s \frac{1}{n_s} \sum_{i=1}^{n_s} \left[y_{si} F(y_{si}) - E\{Y_s F(Y_s)\} + \int_{y_{si}}^{\infty} y dF(y) - E\{\int_{Y_s}^{\infty} y dF(y)\} \right] \\ - (\hat{\mu} - \mu)$$

and $X_{n2} = \hat{\mu} - \mu$. Note that, by the CLT, $\sqrt{n}X_{nj} = O_p(1), j = 1, 2$. Taylor series expansion leads to $1/(X_{n2} + \mu) = \mu^{-1} - \mu^{-2}X_{n2} + O_p(n^{-1})$. Substituting this into Eq. (33), we have

$$\hat{G}_{st} = G + \mu^{-1}(X_{n1} - GX_{n2}) + o_p(n^{-1/2}) = G + \mu^{-1}X_{n3} + o_p(n^{-1/2}),$$
(34)

where

$$\begin{split} X_{n3} &= 2\sum_{s=1}^{S} W_s \frac{1}{n_s} \sum_{i=1}^{n_s} \left[y_{si} F(y_{si}) - E\{Y_s F(Y_s)\} + \int_{y_{si}}^{\infty} y dF(y) - E\{\int_{Y_s}^{\infty} y dF(y)\} \right] \\ &- (G+1) (\hat{\mu} - \mu). \end{split}$$

The result of Theorem 3 then follows from Eq. (34), using the CLT for stratified random sampling. $\hfill \Box$

Proof of Theorem 4. First we can show that

$$\sum_{s=1}^{S} W_{s} n_{s}^{-1} m_{s} \sum_{i=1}^{n_{s}} \{ Z^{*}(y_{si}, \theta) \}^{2} = \sum_{s=1}^{S} W_{s} n_{s}^{-1} m_{s} \sum_{i=1}^{n_{s}} Z^{2}(y_{si}, \theta) + o_{p}(1)$$
(35)
= $o_{a2}^{2} + o_{p}(1).$

Y. Qin et al. / Economic Modelling 27 (2010) 1429-1435

Note that

$$\sum_{s=1}^{S} W_s n_s^{-1} \sum_{i=1}^{n_s} Z^*(y_{si}, \theta) = \sum_{s=1}^{S} W_s n_s^{-1} \sum_{i=1}^{n_s} Z(y_{si}, \theta)$$

Thus, from the proof of Theorem 3, we have

$$\sqrt{n} \sum_{s=1}^{n} W_s n_s^{-1} \sum_{i=1}^{n_s} Z^*(y_{si}, \theta) \xrightarrow{d} N\left(0, \mu^2 \sigma_{a1}^2\right).$$
(36)

From Eqs. (35), (36) and the proof of Theorem 1 in Owen (1990), it can be shown that

$$-2R(\theta) = \left\{ \sum_{s=1}^{S} W_s n_s^{-1} m_s \sum_{i=1}^{n_s} \left\{ Z^*(y_{si}, \theta) \right\}^2 \right\}^{-1} \\ \times \left\{ \sqrt{n} \sum_{s=1}^{S} W_s n_s^{-1} \sum_{i=1}^{n_s} Z^*(y_{si}, \theta) \right\}^2 + o_p(1) \stackrel{d}{\to} \frac{\mu^2 \sigma_{a1}^2}{\sigma_{a2}^2} \chi^2(1)$$

References

- David, H.A., 1968. Gini's mean difference rediscovered. Biometrika 55, 573–575. Gini, C., 1912. Variabilità e mutabilità, contributo allo studio delle distribuzioni e relazioni statistiche. Studi Economico-Giuridici della R. Università di Cagliari.
- Gini, C., 1936. On the measure of concentration with special reference to income and wealth. Abstract of Paper Presented at the Cowles Commission Research Conference on Economics and Statistics. Colorado College Press, Colorado Springs. Giorgi, G.M., Palmitesta, P., Provasi, C., 2006. Asymptotic and bootstrap inference for the generalized Gini Indices. Metron LXIV, 107–124.

- Glasser, G.J., 1962. Variance formulas for the mean difference and coefficient of concentration. Journal of the American Statistical Association 57, 648–654.
- Jing, B.Y., Yuan, J.Q., Zhou, W., 2008. Empirical likelihood for non-degenerate *U*statistics. Statistics and Probability Letters 78, 599–607.
- Kakwani, N.C., 1977. Applications of Lorenz curves in economic analysis. Econometrica 45, 719–727.
- Karagiannis, E., Kovacevic, M., 2000. A method to calculate the jackknife variance estimator for the Gini coefficient. Oxford Bulletin of Economics and Statistics 62, 119–122.
- Lorenz, M.O., 1905. Methods for measuring concentration of wealth. Journal of the American Statistical Association 9, 209–219.
- Owen, A.B., 1990. Empirical likelihood ratio confidence regions. Annals of Statistics 18, 90–120.
- Owen, A.B., 2001. Empirical Likelihood. Chapman & Hall/CRC.
- Sandström, A., Wretman, J.H., Waldén, B., 1985. Variance estimators of the Gini coefficient: simple random sampling. Metron 43, 41–70.
- Sandström, A., Wretman, J.H., Waldén, B., 1988. Variance estimators of the Gini coefficient –probability sampling. Journal of Business & Economic Statistics 6, 113–119. Sendler, W., 1979. On statistical inference in concentration measurement. Metrika 26,
- 109–122. Serfling, R.J., 1980. Approximation Theorems of Mathematical Statistics. John Wiley &
- Sons, New York. van der Vaart, A.W., 1998. Asymptotic Statistics. Cambridge University Press, New York.
- Van der Vaart, A.W., 1990. Asymptotic statistics. Cambridge Oniversity Press, New York. Wood, A.T.A., Do, K.A., Broom, N.M., 1996. Sequential linearization of empirical likelihood constraints with application to U-statistics. Journal of Computational and Graphical Statistics 5, 365–385.
- Wu, C., 2004. Some algorithmic aspects of the empirical likelihood method in survey sampling. Statistica Sinica 14, 1057–1067.
- Wu, C., Rao, J.N.K., 2006. Pseudo-empirical likelihood ratio confidence intervals for complex surveys. The Canadian Journal of Statistics 34, 359–375.
- Yitzhaki, S., 1991. Calculating Jackknife variance estimators for parameters of the Gini method. Journal of Business & Economic Statistics 9, 235–239.
- Zhong, B., Rao, J.N.K., 2000. Empirical likelihood inference under stratified random sampling using auxiliary population information. Biometrika 87, 929–938.