

## Convergence of the Sample Mean Difference to the normal distribution: simulation results

Francesca Greselin <sup>‡</sup>

Michele Zenga <sup>‡</sup>

*Summary:* The present work aims to obtain the value of minimum sample size required by a good approximation by the normal curve for the sample mean difference. Particular care is given to what happens in the tails of the curves, with the aim of deriving confidence intervals for Gini's mean difference. This goal is obtained by empirical methods and the presented results have an explorative nature. Simulation data have been obtained sampling from different distributions, considering symmetry versus asymmetry and the existence of the moments as main aspects in the underlying distribution. These remarks lead to the choice of the normal, the rectangular, the exponential and the Pareto distributions. All the obtained results indicate that the shape of the distribution from which the samples are generated is critically related to the minimum sample sizes required for a good approximation of the tails of the sample mean difference to the normal curve.

*Keywords:* Gini Mean Difference, asymptotic distribution, convergence, U-statistic.

### 1. Introduction

The estimator  $\hat{\Delta}_n$  of Gini's mean difference  $\Delta$ , in a remarkable work due to Hoeffding (1948), was classified as a U-statistics. The whole class of U-statistics owns some optimal properties (Halmos, 1946): among them, they are unbiased estimators for the related functional on the population and they are asymptotically normal; moreover, their variance is a function of the sample size  $n$  and of some main population characteristics. A systematic

---

<sup>‡</sup> Quantitative Methods for Economics and Business Sciences, University of Milano-Bicocca, P.za dell'Ateneo Nuovo 1, 20126 Milano  
(e-mail: [francesca.greselin@unimib.it](mailto:francesca.greselin@unimib.it); [michele.zenga@unimib.it](mailto:michele.zenga@unimib.it))

Although it is the result of a close collaboration, this paper was specifically elaborated as follows: Section 1, 3 and 5 are due to M. Zenga, while Sections 2 and 4 are due to F. Greselin.

presentation of the theory of U-statistics can be found in Koroljuk and Borovskich (1994), although in this work only their central limit theorem is required.

An unbiased estimator  $\hat{V}ar(\hat{\Delta}_n)$  for  $Var(\hat{\Delta}_n)$  was proposed by Cowell (1989) and recently Zenga *et al.* (2004) have derived it through a different methodology. This result allows to make inference about Gini's mean difference  $\Delta$ : for instance, to derive confidence intervals, the convergence of  $\hat{\Delta}_n$  to the normal distribution has to be investigated. In general, the sample size  $n$  that assures a good approximation will depend on the population distribution from which the sample is drawn.

The aim of this work is to explore, by simulation methods, the minimum sample sizes that guarantee a good approximation of the distribution of the statistic  $\hat{\Delta}_n$  by the normal distribution, considering different underlying distributions. Symmetric and asymmetric distributions will be considered, giving particular care in choosing and varying their parameters, mainly when they are related to the existence of moments.

The present paper is organized as follows. Section 2 is devoted to some needed definitions and notations. The methodology is presented and discussed in section 3. The details concerning the simulation experiments are given, along with the discussion of the results, within section 4. Finally, section 5 concludes and points out some possible developments.

## 2. Notations and definitions

Let  $X$  be a continuous random variable (c.r.v.) with probability density function  $f(x)$ , for  $x \in P$ . Gini's mean difference  $\Delta$  is defined by:

$$\Delta = \int \int_{-\infty-\infty}^{+\infty+\infty} |x - y| f(x) f(y) dx dy. \quad (1)$$

Let  $\mu$  and  $\sigma^2$  denote, respectively, the mean and the variance of the c.r.v.  $X$ . In this paper it is assumed that  $\sigma^2$  is finite: this assures the asymptotic convergence to normality of  $\hat{\Delta}$ . Let  $(X_1, \dots, X_i, \dots, X_n)$  denote a random sample (s.w.r.) of size  $n$  ( $n > 3$ ) from the c.r.v.  $X$ , where the r.v.  $X_i$  ( $i = 1, 2, \dots, n$ ) are i.i.d., so that the sample mean difference (hereafter, s.m.d.) without repetition  $\hat{\Delta}_n$  is given by:

$$\hat{\Delta}_n = \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{\substack{j=1 \\ i \neq j}}^n |X_i - X_j|. \quad (2)$$

The expected value of  $\hat{\Delta}_n$  is:

$$E(\hat{\Delta}_n) = \Delta, \quad \text{for every } \Delta. \quad (3)$$

The general formula for the variance of  $\hat{\Delta}_n$  has been derived first by Nair (1936), then, in a simpler form, by Lomnicki (1952) (see also Michetti and Dall'Aglio, 1957), as:

$$\text{Var}(\hat{\Delta}) = \frac{4}{n(n-1)} \left[ \sigma^2 + (n-2)\mathfrak{F} - \frac{(2n-3)}{2}\Delta^2 \right]. \quad (4)$$

where the functional  $\mathfrak{F}$  is given by:

$$\mathfrak{F} = \iiint |x-y||x-z|f(x)f(y)f(z) dx dy dz. \quad (5)$$

The sample mean difference without repetition is hence a mean squared error consistent estimator for  $\Delta$ .

### 3. Measuring the convergence of $\hat{\Delta}_n$ to the normal distribution

In order to evaluate the convergence to normality of  $\hat{\Delta}_n$  a  $X^2$  test may be used, or, more generally, a goodness of fit test. Actually, in order to derive confidence intervals for  $\Delta$ , we are not interested in measuring how  $\hat{\Delta}_n$  behaves in all the real line, but rather how its tails behave. To measure the departure from normality of the distribution of the sample mean difference, the nominal probability  $(1-\alpha)$  assigned to the asymptotic interval:

$$1-\alpha = \lim_{n \rightarrow \infty} P \left\{ \Delta - z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\Delta}_n)} \leq \hat{\Delta}_n \leq \Delta + z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\Delta}_n)} \right\} \quad (6)$$

has to be compared<sup>1</sup> with the probability  $p_{\hat{\Delta}_n}$  obtained for fixed and finite  $n$ :

$$p_{\hat{\Delta}_n} = P \left\{ \Delta - z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\Delta}_n)} \leq \hat{\Delta}_n \leq \Delta + z_{1-\alpha/2} \sqrt{\text{Var}(\hat{\Delta}_n)} \right\} \quad (7)$$

It is worth noting that the probability  $p_{\hat{\Delta}_n}$  defined in (7) depends on the nominal risk  $\alpha$ , the sample size  $n$  and the underlying distribution function  $F(x)$  from which the samples are drawn:  $p_{\hat{\Delta}_n} = p_{\hat{\Delta}_n}(\alpha, n, F(x))$ . If the exact distribution of the sample mean difference  $\hat{\Delta}_n$  was known for some  $F(x)$ , the probability  $p_{\hat{\Delta}_n}$  could be evaluated; unfortunately this has been done only

---

<sup>1</sup> As usual,  $z_{1-\alpha/2}$  is the  $(1-\alpha/2)$ -quantile of the standard normal variable.

for a very few distributions (see Ali 1969, Crocetta and Loperfido 2004), moreover its computational complexity grows as  $n$  increases.

An estimation of the unknown probability  $p_{\hat{\Delta}_n}$  can be obtained by simulations. A high number  $B$  of pseudo samples of size  $n$  are drawn from a given continuous distribution  $F(x)$ ; from each sample a simulated estimation, say  $\delta_i$ , of the sample mean difference  $\hat{\Delta}_n$  is provided, so that  $B$  values  $\delta_i$  ( $i = 1, \dots, B$ ) are obtained. The set of these values can be considered as an empirical variable, say  $\hat{\Delta}^{sim}$ . The ratio between the number of values of  $\delta_i$  falling in the  $(1 - \alpha)$ -confidence interval centred on  $\Delta$  defined in (7) and the total number  $B$  of the simulated estimations  $\delta_i$  is a sample relative frequency:

$$\hat{p}_{\hat{\Delta}_n} = \frac{\#\left\{ \delta_i \mid \delta_i \in \Delta \pm z_{1-\alpha/2} \sqrt{Var(\hat{\Delta}_n)}; i = 1, \dots, B \right\}}{B}. \quad (8)$$

Being  $\hat{p}_{\hat{\Delta}_n}$  a relative frequency, its expected value and variance are given by:

$$E(\hat{p}_{\hat{\Delta}_n}) = p_{\hat{\Delta}_n} \quad \text{and} \quad Var(\hat{p}_{\hat{\Delta}_n}) = \frac{1}{B} p_{\hat{\Delta}_n} (1 - p_{\hat{\Delta}_n}). \quad (9)$$

The values obtained by  $\hat{p}_{\hat{\Delta}_n}$ , for large values of  $B$ , are good estimations of the unknown probability  $p_{\hat{\Delta}_n}$ .

In order to appreciate and discuss the simulated results, the same approach is carried out for confidence intervals for  $E(X) = \mu$ , so giving the possibility of a useful comparison with an analogous - but well known - situation. The sample relative frequencies, in this case, are given by:

$$\hat{p}_{\bar{X}_n} = \frac{\#\left\{ \bar{x}_i \mid \bar{x}_i \in \mu \pm z_{1-\alpha/2} \sqrt{Var(\bar{X}_n)}; i = 1, \dots, B \right\}}{B} \quad (10)$$

where  $\bar{x}_i$  is the simulated value of the sample mean  $\bar{X}_n = n^{-1} \sum_{j=1}^n X_j$ , obtained for each of the  $B$  samples. For high values of  $B$  the values of  $\hat{p}_{\bar{X}_n}$  are good approximations of the unknown probabilities  $p_{\bar{X}_n}$ :

$$p_{\bar{X}_n} = P\left\{ \mu - z_{1-\alpha/2} \sqrt{Var(\bar{X}_n)} \leq \bar{X}_n \leq \mu + z_{1-\alpha/2} \sqrt{Var(\bar{X}_n)} \right\}, \quad (11)$$

and  $\lim_{n \rightarrow \infty} p_{\bar{X}_n} = 1 - \alpha$ .

In the following section, a series of simulation data will be presented, for

different values of  $\alpha$  and the sample size  $n$ , and for different distributions of the continuous random variable  $X$ .

#### 4. Minimum sample sizes assuring a good approximation

Four different continuous distributions, widely employed in modelling real data, will be considered throughout this section: the Normal, the Rectangular, the Exponential and the Pareto distributions. By these fairly simple models we can anyway explore and compare symmetric and asymmetric cases, and the influence of the existence of the moments (related to tails heaviness) in the underlying distribution.

The simulations are obtained by a software program written in C++, using pseudorandom numbers generated by the IMSL statistical library.

For each continuous population,  $B = 10,000$  samples of size  $n$  were drawn, so that each sample provides a simulated estimation of the sample mean difference  $\hat{\Delta}_n$ , varying  $n$  from  $n = 30$  up to  $n = 960$ , doubling each time the sample size. The confidence level is initially fixed at  $1 - \alpha = 0.99$ , successively a range of values for  $\alpha$  is considered.

For each group of  $B = 10,000$  samples, the following synthetic values were evaluated:

- the relative frequencies  $\hat{p}_{\hat{\Delta}_n}$  and  $\hat{p}_{\bar{X}_n}$ ;
- the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$  derived from the pseudo-generated distribution of the sample mean difference  $\hat{\Delta}^{sim}$ , to be compared with their asymptotic values, respectively given by:
 
$$\hat{\Delta}_{\alpha/2} = \Delta - z_{1-\alpha/2} \sqrt{Var(\hat{\Delta}_n)} \quad \text{and} \quad \hat{\Delta}_{1-\alpha/2} = \Delta + z_{1-\alpha/2} \sqrt{Var(\hat{\Delta}_n)};$$
- the simulated quantiles  $\bar{X}_{\alpha/2}^{sim}$  and  $\bar{X}_{1-\alpha/2}^{sim}$  to be compared with their asymptotic values, respectively given by:
 
$$\bar{X}_{\alpha/2} = \mu - z_{1-\alpha/2} \sqrt{Var(\bar{X}_n)} \quad \text{and} \quad \bar{X}_{1-\alpha/2} = \mu + z_{1-\alpha/2} \sqrt{Var(\bar{X}_n)};$$
- the mean value  $M_l(\hat{\Delta}^{sim})$ , to be compared with  $\Delta$ ;
- the median  $Me(\hat{\Delta}^{sim})$  of the 10,000 estimations;
- the standard deviation  $\sigma(\hat{\Delta}^{sim})$  of the 10,000 estimations;
- the standardized third moment of  $\hat{\Delta}^{sim}$  as an asymmetry index;
- the standardized fourth moment of  $\hat{\Delta}^{sim}$  as a kurtosis index.

All simulation results, for each chosen continuous population, and for each given sample size  $n$ , are summarized in a Table. In each row of the Table, with reference to a specific value of the risk  $\alpha$ , some synthetic values for the group of  $B$  samples are provided. For each group of  $B$  samples a graphical

representation of the simulated distribution  $\hat{\Delta}^{sim}$  is also reported.

As said before, the aim of simulations is to investigate the minimum sampling sizes for which  $\hat{p}_{\hat{\Delta}_n}$  is a good estimation for  $p_{\Delta_n}$ , for applications.

The approximation will be considered good if, for  $\alpha = 0.05$ , the simulated risk is in the range 0.04 – 0.06. More generally, a percent absolute deviation less than 20% between the nominal risk and the simulated one will be considered a good approximation (Vesserau 1957; Zenga, 1974).

#### 4.1 Normal distribution

Let  $X$  be the normal c.r.v. with probability density function:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \text{for } \sigma > 0, \mu \in \mathbb{P}.$$

The values of the parameters can be fixed by  $\mu = 0$ ,  $\sigma^2 = 25$ . With this choice<sup>2</sup>, the value of Gini's mean difference is  $\Delta = 5.64190$ .

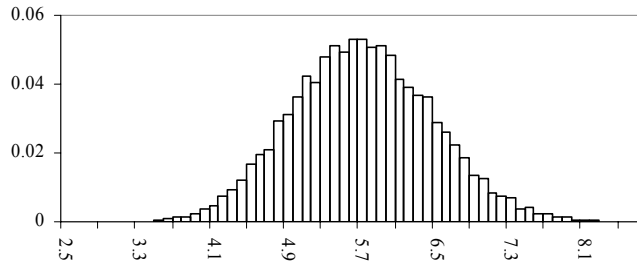
**Table 1.** Simulation data for the evaluation of the tails of the distribution of the s..m.d.  $\hat{\Delta}_n$ , sampling from  $N(\mu=0.0; \sigma=5.0)$  ( $n=30$ ,  $B=10,000$  samples)

$\hat{\Delta}_{\alpha/2}^{sim}$	$\hat{\Delta}_{1-\alpha/2}^{sim}$	$\hat{\Delta}_{\alpha/2}$	$\hat{\Delta}_{1-\alpha/2}$	$\hat{p}_{\hat{\Delta}_n}$	$1-\alpha$	Acc. range	$\bar{X}_{\alpha/2}^{sim}$	$\bar{X}_{1-\alpha/2}^{sim}$	$\bar{X}_{\alpha/2}$	$\bar{X}_{1-\alpha/2}$	$\hat{p}_{\bar{X}_n}$
3.847	7.697	3.705	7.579	<b>0.990</b>	0.99	0.988-0.992	-2.279	2.365	-2.351	2.351	<b>0.991</b>
4.049	7.432	3.956	7.328	<b>0.974</b>	0.975	0.970-0.980	-2.027	2.078	-2.046	2.046	<b>0.974</b>
4.229	7.198	4.168	7.116	<b>0.950</b>	0.95	0.940-0.960	-1.792	1.830	-1.789	1.789	<b>0.947</b>
4.357	7.026	4.303	6.981	<b>0.925</b>	0.925	0.910-0.940	-1.617	1.665	-1.625	1.625	<b>0.920</b>
4.435	6.926	4.405	6.879	<b>0.900</b>	0.90	0.880-0.920	-1.500	1.549	-1.502	1.502	<b>0.896</b>
4.579	6.753	4.559	6.725	<b>0.847</b>	0.85	0.820-0.880	-1.331	1.350	-1.314	1.314	<b>0.843</b>

Table 1 shows the simulated data obtained by generating  $B = 10,000$  samples of size  $n = 30$ . In each row a different value of the risk  $\alpha$  is considered. The values of the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$ , in the first and second column, are very close to their asymptotic values  $\hat{\Delta}_{\alpha/2}$  and  $\hat{\Delta}_{1-\alpha/2}$ , respectively given in the third and fourth column. All simulated probabilities  $\hat{p}_{\hat{\Delta}_n}$ , evaluated by (8) with reference to the interval

<sup>2</sup> Actually, as the analytic expression of the third and fourth standardized moment in the normal distribution -as well as in the rectangular and in the exponential- does not depend on the parameters, any choice for their values would yield the same results on  $\hat{p}_{\hat{\Delta}_n}$  and on  $\hat{p}_{\bar{X}_n}$ .

$\Delta \pm z_{1-\alpha/2} \sqrt{Var(\hat{\Delta}_n)}$ , are in the range of acceptable values indicated in the central column (in this and in the following tables, acceptable values are highlighted). This means that  $n = 30$  already assures a good behaviour of the tails of the distribution of  $\hat{\Delta}_n$ . Their performance is almost as good as that of the tails of  $\bar{X}_n$ , whose distribution is known to be normal. The probabilities  $\hat{p}_{\bar{X}_n}$  were indeed expected to be in the acceptable range, as the last column of Table 1 shows. All simulation data for  $\hat{\Delta}_n$ , with  $n = 30$ , are shown in Figure 1.



**Figure 1.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by  $B=10,000$  samples,  $n=30$ , drawn from  $N(\mu = 0.0; \sigma = 5.0)$ )

Table 2 shows some moments and indices of the simulated distribution of  $\hat{\Delta}_n$  and offers a comparison with the same quantities obtained on its asymptotic distribution: the mean value  $M_1(\hat{\Delta}^{sim})$ , to be compared with  $\Delta$  and with the median  $Me(\hat{\Delta}^{sim})$  of the 10,000 estimations; the standard deviation  $\sigma(\hat{\Delta}^{sim})$ , the standardized third moment  $\alpha_3(\hat{\Delta}^{sim})$  as an asymmetry index and the standardized fourth moment  $\beta_2(\hat{\Delta}^{sim})$  as a kurtosis index.

**Table 2.** Some indices of the asymptotic distribution  $N(\mu = \Delta; \sigma = \sqrt{Var(\hat{\Delta}_n)})$  of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B = 10,000$  samples from  $N(\mu = 0.0; \sigma = 5.0)$ ,  $n = 30$

	Simulated values	Asymptotic values
$M_1(\hat{\Delta}^{sim})$	5.65428	$E(\hat{\Delta}_n) = 5.64190$
$Me(\hat{\Delta}^{sim})$	5.63574	5.64190
$\alpha_3(\hat{\Delta}^{sim})$	0.14260	0.0
$\beta_2(\hat{\Delta}^{sim})$	2.98157	3.0
$\sigma(\hat{\Delta}^{sim})$	0.75651	$\sqrt{Var(\hat{\Delta}_n)} = 0.75204$

Actually, the approximation of  $\hat{\Delta}_n$  to the normal is good enough for all considered sample sizes. Indeed, the main simulation results for higher samples sizes  $n = 60, 120, 240, 480$  and  $960$  are shown in Table 3.

**Table 3.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$ , sampling from  $N(\mu=0.0; \sigma=5.0)$  (by  $B=10,000$  samples)

		$n = 60$		$n = 120$		$n = 240$		$n = 480$		$n = 960$	
$1-\alpha$	Acc. range	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$
0.99	0.988-0.992	<b>0.991</b>	<b>0.991</b>	<b>0.992</b>	<b>0.988</b>	<b>0.990</b>	<b>0.991</b>	<b>0.991</b>	<b>0.990</b>	<b>0.991</b>	<b>0.988</b>
0.975	0.970-0.980	<b>0.976</b>	<b>0.976</b>	<b>0.975</b>	<b>0.975</b>	<b>0.977</b>	<b>0.975</b>	<b>0.975</b>	<b>0.977</b>	<b>0.976</b>	<b>0.974</b>
0.95	0.940-0.960	<b>0.951</b>	<b>0.951</b>	<b>0.949</b>	<b>0.948</b>	<b>0.954</b>	<b>0.950</b>	<b>0.952</b>	<b>0.951</b>	<b>0.949</b>	<b>0.950</b>
0.925	0.910-0.940	<b>0.929</b>	<b>0.927</b>	<b>0.922</b>	<b>0.925</b>	<b>0.927</b>	<b>0.927</b>	<b>0.928</b>	<b>0.926</b>	<b>0.924</b>	<b>0.925</b>
0.90	0.880-0.920	<b>0.904</b>	<b>0.902</b>	<b>0.899</b>	<b>0.900</b>	<b>0.902</b>	<b>0.903</b>	<b>0.904</b>	<b>0.902</b>	<b>0.901</b>	<b>0.900</b>
0.85	0.820-0.880	<b>0.856</b>	<b>0.854</b>	<b>0.852</b>	<b>0.849</b>	<b>0.851</b>	<b>0.852</b>	<b>0.852</b>	<b>0.850</b>	<b>0.847</b>	<b>0.848</b>

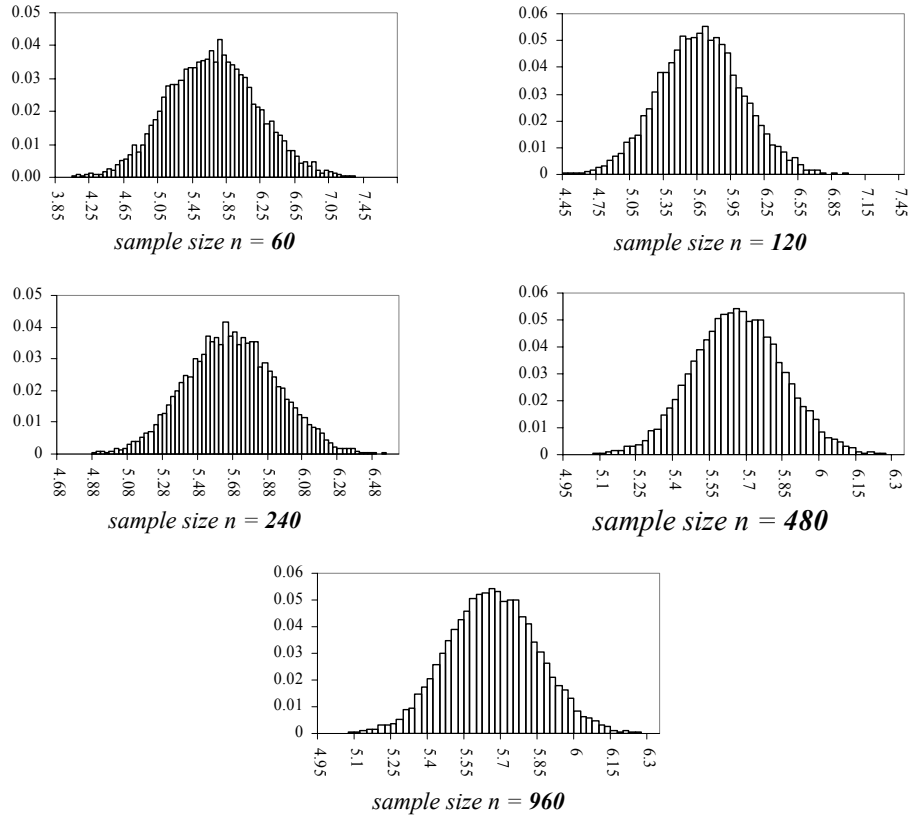
All the simulated data have shown a very good approximation of the probability associated with the tails of the distribution of the sample mean difference, even from the minimum sample size. The simulated distributions can also be synthesized by the indices presented in Table 4, in which the median of the simulated values  $Me(\hat{\Delta}^{sim})$  is increasing as  $n$  increases, approaching their mean  $M_I(\hat{\Delta}^{sim})$ . The decreasing values of the third standardized simulated moment  $\alpha_3(\hat{\Delta}^{sim})$ , as  $n$  increases, denote a slight and decreasing asymmetry. A very low degree of platykurtosis is roughly assessed by the standardized fourth moment  $\beta_2(\hat{\Delta}^{sim})$ , for all values of the sample size  $n$ . The specific expression for the variance of  $\hat{\Delta}_n$ , sampling from the normal distribution, can be found in Nair (1936) (see also Kendall *et al.* (1994), Zenga *et al.* (2004)).

**Table 4.** Some indices of the asymptotic distribution of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B=10,000$  s.w.r. from  $N(\mu=0.0; \sigma=5.0)$ , and the theoretical values for  $\sqrt{Var(\hat{\Delta}_n)}$  ( $n=30, \dots, 960$ )

Simul. values	$n = 30$	$n = 60$	$n = 120$	$n = 240$	$n = 480$	$n = 960$	Asympt. values
$M_I(\hat{\Delta}^{sim})$	5.65428	5.63921	5.63672	5.64677	5.64245	5.64265	$E(\hat{\Delta}_n)=5.64190$
$Me(\hat{\Delta}^{sim})$	5.63574	5.63816	5.63384	5.64157	5.64157	5.64185	5.64190
$\alpha_3(\hat{\Delta}^{sim})$	0.14260	0.06974	0.08224	0.05474	0.01770	0.01654	0.0
$\beta_2(\hat{\Delta}^{sim})$	2.98157	2.95252	2.96829	2.94844	2.96384	2.94698	3.0
$\sigma(\hat{\Delta}^{sim})$	0.75651	0.52196	0.36934	0.25969	0.18342	0.13075	
$\sqrt{Var(\hat{\Delta}_n)}$	0.75204	0.52623	0.37018	0.26108	0.18438	0.13029	



The distribution of  $\hat{\Delta}^{sim}$  for all chosen values of the sample size  $n > 30$  is shown in Figure 2.



**Figure 2.** Simulated distribution of the s.m.d.  $\hat{\Delta}$  (by 10,000 samples, drawn from the  $Normal(\mu=0; \sigma=5)$ )

While  $n$  increases, the range of values assumed by  $\hat{\Delta}^{sim}$  reduces and concentrates around  $\Delta = 5.641896$ . The simulated distribution is progressively more symmetric and bell-shaped as  $n$  increases, and a slightly unwieldy behaviour on the tails is gradually smoothing. In any case, it is not far from symmetry even for the lowest value of the sample size ( $n = 30$ ).

#### 4.2 Rectangular distribution

Let  $X$  be the c.r.v. with probability density function:

$$f(x) = \frac{1}{b-a} \quad \text{for } a \leq x \leq b, (b > a).$$

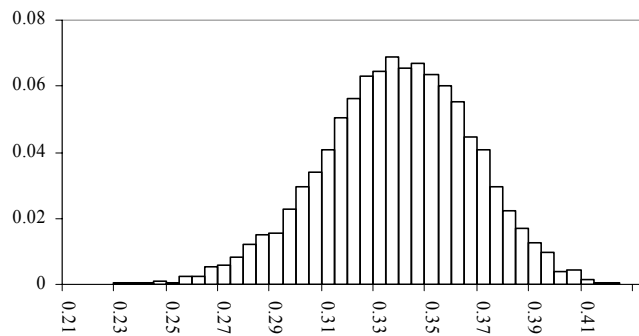
Let us fix the parameters  $a$  and  $b$  of this distribution by  $a = 0$  and  $b = 1$ , hence the main characteristics of the distribution are given by  $\mu = 0.5$ ,  $\sigma^2 = 0.083$  and  $\Delta = 0.333$ . The simulated data obtained for this distribution, for  $n = 30$ , are shown in Table 5.

**Table 5.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$ , sampling from the Rectangular ( $a=0, b=1$ ) ( $n=30, B=10,000$  samples)

$\hat{\Delta}_{\alpha/2}^{sim}$	$\hat{\Delta}_{1-\alpha/2}^{sim}$	$\hat{\Delta}_{\alpha/2}$	$\hat{\Delta}_{1-\alpha/2}$	$\hat{p}_{\hat{\Delta}_n}$	$1-\alpha$	Acc. range	$\bar{X}_{\alpha/2}^{sim}$	$\bar{X}_{1-\alpha/2}^{sim}$	$\bar{X}_{\alpha/2}$	$\bar{X}_{1-\alpha/2}$	$\hat{p}_{\bar{X}_n}$
0.254	0.403	0.259	0.408	<b>0.991</b>	0.99	0.988-0.992	0.371	0.634	0.364	0.636	<b>0.992</b>
0.265	0.394	0.268	0.398	<b>0.975</b>	0.975	0.970-0.980	0.386	0.618	0.382	0.618	<b>0.978</b>
0.274	0.388	0.276	0.390	<b>0.949</b>	0.95	0.940-0.960	0.402	0.602	0.397	0.603	<b>0.957</b>
0.279	0.384	0.282	0.385	<b>0.923</b>	0.925	0.910-0.940	0.409	0.593	0.406	0.594	<b>0.931</b>
0.284	0.380	0.286	0.381	<b>0.898</b>	0.90	0.880-0.920	0.416	0.587	0.413	0.587	<b>0.905</b>
0.291	0.375	0.292	0.375	<b>0.851</b>	0.85	0.820-0.880	0.425	0.576	0.424	0.576	<b>0.854</b>

The values of the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$  are slightly lower than their asymptotic values  $\hat{\Delta}_{\alpha/2}$  and  $\hat{\Delta}_{1-\alpha/2}$  (columns 1-4 of Table 5). The simulated probabilities  $\hat{p}_{\hat{\Delta}_n}$  are steadily in the range of acceptable values indicated in the central column (all values are hence highlighted). This means that  $n = 30$  already assures a good behaviour of the tails of the distribution of  $\hat{\Delta}_n$ . Their performance is similar to that of the tails of  $\bar{X}_n$  measured by  $\hat{p}_{\bar{X}_n}$ , referring to the sample mean, shown and highlighted (acceptable values) in last column of Table 5.

The simulated distribution of  $\hat{\Delta}_n$ , for  $n = 30$ , is represented in Figure 3.



**Figure 3.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by  $B=10,000$  s.w.r.,  $n=30$ ), drawn from the Rectangular distribution ( $a=0, b=1$ )

Table 6 shows all the results for different values of the sample size  $n$ :

**Table 6.** Simulation data for the evaluation of the tails of the distr. of the s.m.d.  $\hat{\Delta}_n$  sampling from the Rectangular distribution ( $a = 0, b = 1$ ) (by  $B=10,000$  samples)

$1-\alpha$	Acc. range	$n = 60$		$n = 120$		$n = 240$		$n = 480$		$n = 960$	
		$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$
0.99	0.988-0.992	<b>0.991</b>	<b>0.991</b>	<b>0.992</b>	<b>0.990</b>	<b>0.991</b>	<b>0.990</b>	<b>0.990</b>	<b>0.990</b>	<b>0.989</b>	<b>0.990</b>
0.975	0.970-0.980	<b>0.976</b>	<b>0.977</b>	<b>0.976</b>	<b>0.975</b>	<b>0.976</b>	<b>0.973</b>	<b>0.976</b>	<b>0.975</b>	<b>0.974</b>	<b>0.974</b>
0.95	0.940-0.960	<b>0.952</b>	<b>0.955</b>	<b>0.953</b>	<b>0.951</b>	<b>0.949</b>	<b>0.948</b>	<b>0.954</b>	<b>0.950</b>	<b>0.950</b>	<b>0.951</b>
0.925	0.910-0.940	<b>0.927</b>	<b>0.930</b>	<b>0.929</b>	<b>0.924</b>	<b>0.924</b>	<b>0.922</b>	<b>0.928</b>	<b>0.925</b>	<b>0.922</b>	<b>0.927</b>
0.90	0.880-0.920	<b>0.903</b>	<b>0.903</b>	<b>0.905</b>	<b>0.902</b>	<b>0.897</b>	<b>0.898</b>	<b>0.901</b>	<b>0.897</b>	<b>0.898</b>	<b>0.902</b>
0.85	0.820-0.880	<b>0.851</b>	<b>0.853</b>	<b>0.856</b>	<b>0.851</b>	<b>0.849</b>	<b>0.849</b>	<b>0.849</b>	<b>0.850</b>	<b>0.848</b>	<b>0.850</b>

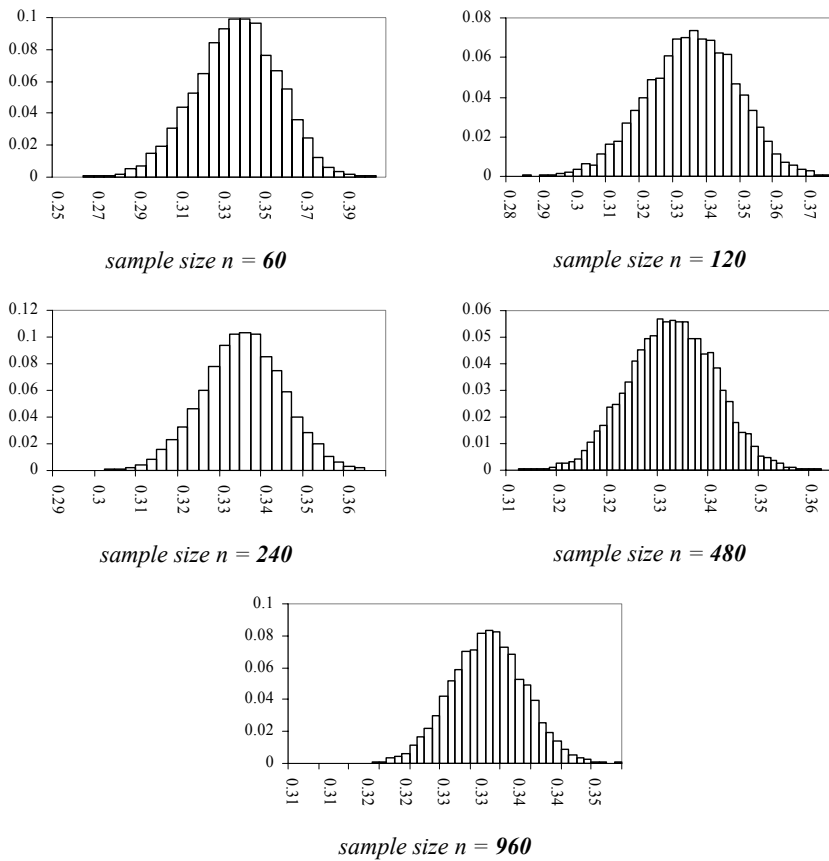
All the simulated probabilities perform well (highlighted values), for each of the chosen values of the risk  $\alpha$ , and for each value of the sample size  $n$ . The sampling distribution of  $\hat{\Delta}_n$  might be reasonably approximated by a normal curve already for  $n = 30$ , when the sample is drawn from the Rectangular distribution.

A synthetic overview of some indices of the simulated distribution is provided by Table 7, jointly with their theoretical or asymptotic values for meaningful comparisons. The sign of the third standardized simulated moment  $\alpha_3(\hat{\Delta}^{sim})$  assesses some slight negative asymmetry, and it becomes even slighter as the sample size increases. In any case, the degree of asymmetry is so weak that it does not affect the coverage probabilities of Table 6. The kurtosis index  $\beta_2(\hat{\Delta}^{sim})$  roughly indicates normal peakedness and tails. The specific expression for the standard error of  $\hat{\Delta}_n$  in the case of sampling from the Rectangular distribution can be found in Nair (1936) (it is also recalled in Kendall *et al.* (1994) or Zenga *et al.* (2004)).

**Table 7.** Values of some indices of the asymptotic distribution of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B=10,000$  s.w.r. from the Rectangular distribution ( $a=0, b=1$ ) and the theoretical values of  $\sqrt{Var(\hat{\Delta}_n)}$

Simul. values	$n = 30$	$n = 60$	$n = 120$	$n = 240$	$n = 480$	$n = 960$	Asympt. values
$M_1(\hat{\Delta}^{sim})$	0.33411	0.33334	0.33334	0.33324	0.33340	0.33330	$E(\hat{\Delta}_n)=0.33333$
$Me(\hat{\Delta}^{sim})$	0.33500	0.33400	0.33362	0.33331	0.33342	0.33333	0.33333
$\alpha_3(\hat{\Delta}^{sim})$	-0.23364	-0.16642	-0.13204	-0.06169	-0.03328	-0.03578	0.0
$\beta_2(\hat{\Delta}^{sim})$	3.05468	2.99126	2.94896	3.00188	2.92557	3.07993	3.0
$\sigma(\hat{\Delta}^{sim})$	0.02906	0.01989	0.01367	0.00967	0.00684	0.00486	
$\sqrt{Var(\hat{\Delta}_n)}$	0.02903	0.01988	0.01382	0.00970	0.00686	0.00480	

Figure 4 gives a graphical representations of the simulated distributions, for sample sizes  $n = 60, 120, 240, 480$  and  $960$ .



**Figure 4.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples, drawn from the Rectangular ( $a=0, b=1$ )).

### 4.3 Exponential distribution

Let  $X$  be the c.r.v. with probability density function:

$$f(x) = \theta e^{-\theta x} \quad \text{for } x \geq 0, \text{ where } \theta > 0.$$

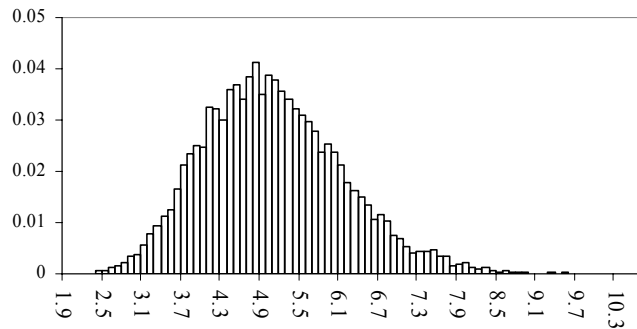
For this distribution, choosing  $\theta = 0.2$ , it holds that  $\mu = 5$ ,  $\sigma^2 = 25$  and  $\Delta = 5$ . Table 8 reports the simulated data obtained for this distribution.

**Table 8.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$  sampling from the Exponential with  $\theta = 0.2$  ( $n=30$ ,  $B=10,000$  samples)

$\hat{\Delta}_{\alpha/2}^{sim}$	$\hat{\Delta}_{1-\alpha/2}^{sim}$	$\hat{\Delta}_{\alpha/2}$	$\hat{\Delta}_{1-\alpha/2}$	$\hat{p}_{\hat{\Delta}_n}$	$1-\alpha$	Acc. range	$\bar{X}_{\alpha/2}^{sim}$	$\bar{X}_{1-\alpha/2}^{sim}$	$\bar{X}_{\alpha/2}$	$\bar{X}_{1-\alpha/2}$	$\hat{p}_{\bar{X}_n}$
2.720	8.159	2.262	7.739	<b>0.988</b>	0.99	0.988-0.992	2.975	7.705	2.649	7.351	<b>0.988</b>
2.965	7.720	2.617	7.383	<b>0.972</b>	0.975	0.970-0.980	3.179	7.327	2.954	7.046	<b>0.974</b>
3.180	7.381	2.916	7.084	<b>0.952</b>	0.95	0.940-0.960	3.392	6.967	3.211	6.789	<b>0.948</b>
3.305	7.083	3.107	6.893	<b>0.931</b>	0.925	0.910-0.940	3.518	6.777	3.375	6.625	<b>0.923</b>
3.418	6.883	3.251	6.749	<b>0.907</b>	0.90	0.880-0.920	3.605	6.648	3.498	6.502	<b>0.899</b>
3.588	6.632	3.470	6.530	<b>0.859</b>	0.85	0.820-0.880	3.748	6.416	3.686	6.314	<b>0.846</b>

The values of the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$ , in the first and second column, are not far from their asymptotic values  $\hat{\Delta}_{\alpha/2}$  and  $\hat{\Delta}_{1-\alpha/2}$ , respectively given in the third and fourth column. The simulated probabilities  $\hat{p}_{\hat{\Delta}_n}$ , evaluated by (8) with reference to the interval  $\Delta \pm z_{1-\alpha/2} \sqrt{Var(\hat{\Delta}_n)}$ , are all in the range of acceptable values indicated in the central column (good values are highlighted). This means that  $n = 30$  already assures a good behaviour of the tails of the distribution of  $\hat{\Delta}_n$ . Their performance is almost as good as that of the tails of  $\bar{X}_n$ : the behaviour of  $\hat{p}_{\bar{X}_n}$ , referring to the sample mean, is in the acceptable range, for  $\alpha > 0.01$ .

A simulated distribution of  $\hat{\Delta}_n$  for  $n = 30$  is graphically represented in Figure 5.



**Figure 5.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by  $B=10,000$  samples,  $n=30$ , drawn from the Exponential distribution with  $\theta = 0.2$ )

As before, simulation data for different values of the sample size  $n$  were obtained. They are shown in detail in Table 9.

**Table 9.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$ , sampling from the Exponential distr. with  $\theta = 0.2$  (by  $B=10,000$  samples)

$1-\alpha$	Acc. range	$n = 60$		$n = 120$		$n = 240$		$n = 480$		$n = 960$	
		$\hat{P}_{\hat{\Delta}_n}$	$\hat{P}_{\bar{X}_n}$	$\hat{P}_{\hat{\Delta}_n}$	$\hat{P}_{\bar{X}_n}$	$\hat{P}_{\hat{\Delta}_n}$	$\hat{P}_{\bar{X}_n}$	$\hat{P}_{\hat{\Delta}_n}$	$\hat{P}_{\bar{X}_n}$	$\hat{P}_{\hat{\Delta}_n}$	$\hat{P}_{\bar{X}_n}$
0.99	0.988-0.992	<b>0.988</b>	<b>0.989</b>	0.987	<b>0.988</b>	<b>0.990</b>	<b>0.990</b>	<b>0.990</b>	<b>0.989</b>	<b>0.990</b>	<b>0.990</b>
0.975	0.970-0.980	<b>0.975</b>	<b>0.978</b>	<b>0.973</b>	<b>0.974</b>	<b>0.973</b>	<b>0.976</b>	<b>0.976</b>	<b>0.973</b>	<b>0.976</b>	<b>0.977</b>
0.95	0.940-0.960	<b>0.952</b>	<b>0.955</b>	<b>0.951</b>	<b>0.949</b>	<b>0.946</b>	<b>0.951</b>	<b>0.951</b>	<b>0.948</b>	<b>0.953</b>	<b>0.951</b>
0.925	0.910-0.940	<b>0.928</b>	<b>0.928</b>	<b>0.926</b>	<b>0.921</b>	<b>0.923</b>	<b>0.925</b>	<b>0.926</b>	<b>0.925</b>	<b>0.927</b>	<b>0.926</b>
0.90	0.880-0.920	<b>0.905</b>	<b>0.903</b>	<b>0.902</b>	<b>0.896</b>	<b>0.899</b>	<b>0.899</b>	<b>0.899</b>	<b>0.901</b>	<b>0.901</b>	<b>0.901</b>
0.85	0.820-0.880	<b>0.856</b>	<b>0.854</b>	<b>0.851</b>	<b>0.845</b>	<b>0.852</b>	<b>0.8482</b>	<b>0.853</b>	<b>0.847</b>	<b>0.851</b>	<b>0.852</b>

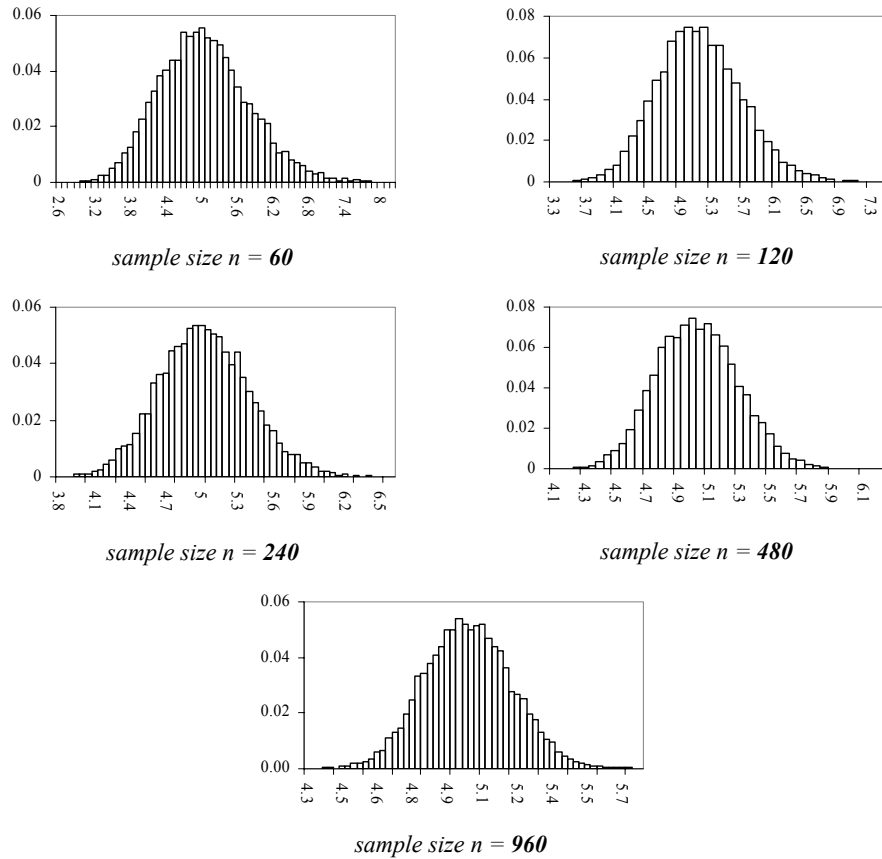
All simulated probabilities perform well, for each of the chosen values of the risk  $\alpha$ , whenever  $n \geq 30$ .

Table 10 gives, for each value of the sample size  $n$ , a synthetic view of the some characteristics of the simulated values  $\hat{\Delta}^{sim}$ : the mean, the median, the standard deviation, the standardized third and fourth moments. The standard error of  $\hat{\Delta}_n$  can be evaluated from Nair (1936) (see also Kendall *et al.* (1994) or Zenga *et al.* (2004)).

**Table 10.** Values of some indices of the asymptotic distribution of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B = 10,000$  s.w.r. from the Exponential distribution with  $\theta = 0.2$ , and the theoretical values for  $\sqrt{Var(\hat{\Delta}_n)}$

Simul. values	$n = 30$	$n = 60$	$n = 120$	$n = 240$	$n = 480$	$n = 960$	Asympt. values
$M_l(\hat{\Delta}^{sim})$	5.01117	4.99431	4.99754	4.99839	5.00434	5.00127	$E(\hat{\Delta}_n) = 5.0$
$Me(\hat{\Delta}^{sim})$	4.92976	4.95351	4.97316	4.98802	4.99710	4.99805	5.0
$\alpha_3(\hat{\Delta}^{sim})$	0.46358	0.36723	0.26683	0.16338	0.13911	0.09340	0.0
$\beta_2(\hat{\Delta}^{sim})$	3.30478	3.17737	3.15296	3.03846	2.91580	2.98820	3.0
$\sigma(\hat{\Delta}^{sim})$	1.06718	0.75048	0.53430	0.37523	0.26399	0.18635	
$\sqrt{Var(\hat{\Delta}_n)}$	1.06314	0.74851	0.52815	0.37307	0.26366	0.18639	

The observed difference between the mean and the median of  $\hat{\Delta}^{sim}$  denotes the asymmetry of the distribution of  $\hat{\Delta}_n$ . However, this asymmetry becomes slighter as the sample size increases, as the decreasing values of the third moment  $\alpha_3(\hat{\Delta}^{sim})$  show. The shape indices of the distribution of  $\hat{\Delta}^{sim}$  reveal some minor platykurtic behaviour for lower values of  $n$ , while  $\hat{\Delta}^{sim}$  is slightly mesokurtic whenever  $n > 240$ . These remarks are confirmed by the graphical representations of the simulated distributions, for different values of the sample size  $n$ , shown in Figure 6.



**Figure 6.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples, drawn from the Exponential,  $\theta = 0.2$ )

#### 4.4 Pareto distribution (first case: with parameters $x_0=1, \theta=3$ )

Let  $X$  be the c.r.v. with probability density function:

$$f(x) = \theta x_0^\theta x^{-(\theta+1)} \quad \text{for } x \geq x_0, \text{ where } x_0 > 0 \text{ and } \theta > 0.$$

As the third and fourth standardized moment depends on  $\theta$ , we have to consider the entire range of values  $\theta > 2$  (in order to assume the existence of  $Var(X)$ ). Actually, even if a broad set of different simulations were done incrementing  $\theta$ , only two cases are presented in this work because they depict clearly what happens on the distribution of  $\hat{\Delta}_n$ .

For this first case, let the value of the parameters be given by  $x_0 = 1$  and

$\theta = 3$ , so that  $X$  has finite moments of order  $r$  only for  $r < 3$ . With this choice of parameters, the main characteristics of  $X$  are:  $\mu = 1.5$ ,  $\sigma^2 = 0.75$  and  $\Delta = 0.6$ .

The aim, here, is to observe the behaviour of  $\hat{\Delta}_n$  when the sample is drawn from an asymmetric continuous variable that does not possess all finite moments. In other words, our purpose here –and in the following section– is to observe how the heaviness of the tails, expressed by  $x^{-(\theta+1)}$ , influences the convergence of the sample mean difference.

Table 11 summarizes the simulation results for a sample size  $n = 30$ .

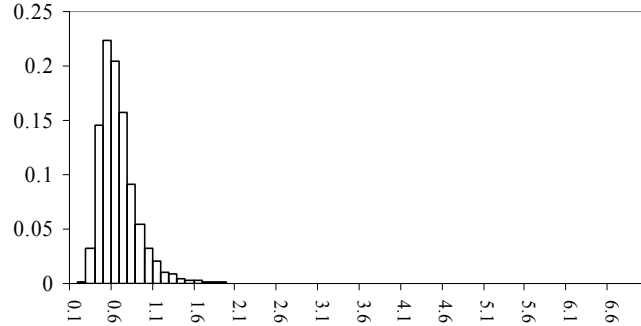
**Table 11.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$ , sampling from the Pareto ( $x_0 = 1$ ,  $\theta = 3$ ) ( $n=30$ ,  $B=10,000$  samples)

$\hat{\Delta}_{\alpha/2}^{sim}$	$\hat{\Delta}_{1-\alpha/2}^{sim}$	$\hat{\Delta}_{\alpha/2}$	$\hat{\Delta}_{1-\alpha/2}$	$\hat{p}_{\hat{\Delta}_n}$	$1-\alpha$	Acc. range	$\bar{X}_{\alpha/2}^{sim}$	$\bar{X}_{1-\alpha/2}^{sim}$	$\bar{X}_{\alpha/2}$	$\bar{X}_{1-\alpha/2}$	$\hat{p}_{\bar{X}_n}$
0.241	1.723	-0.078	1.278	0.981	0.99	0.988-0.992	1.244	2.138	1.093	1.907	0.983
0.263	1.415	0.010	1.190	<b>0.973</b>	0.975	0.970-0.980	1.266	1.949	1.146	1.854	<b>0.975</b>
0.288	1.216	0.084	1.116	0.966	0.95	0.940-0.960	1.285	1.855	1.190	1.810	0.965
0.305	1.093	0.131	1.069	0.959	0.925	0.910-0.940	1.296	1.800	1.218	1.782	0.956
0.319	1.026	0.167	1.033	0.951	0.90	0.880-0.920	1.307	1.764	1.240	1.760	0.944
0.340	0.942	0.221	0.979	0.935	0.85	0.820-0.880	1.325	1.716	1.272	1.728	0.915

The values of the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$  are considerably higher than their asymptotic values  $\hat{\Delta}_{\alpha/2}$  and  $\hat{\Delta}_{1-\alpha/2}$  (respectively in columns 1-2 and 3-4 of Table 11). The simulated probabilities  $\hat{p}_{\hat{\Delta}_n}$  are outside the range of acceptable values indicated in the central column (except for the case of  $\alpha = 0.025$ , for which the good value is highlighted). More specifically,  $\hat{p}_{\hat{\Delta}_n}$  is greater than the right threshold of acceptable ranges for  $\alpha > 0.025$ , and it is lower than the left threshold for  $\alpha = 0.01$ . This suggests a high asymmetry in the distribution of  $\hat{\Delta}_n$ , so that the tails of the distribution of  $\hat{\Delta}_n$  do not approximately behave as the tails of the normal distribution. The results observed for  $\hat{p}_{\hat{\Delta}_n}$  are confirmed also by  $\hat{p}_{\bar{X}_n}$ , indicating an analogous behaviour for the tails of  $\bar{X}_n$ .

The simulated distribution of  $\hat{\Delta}_n$  for  $n = 30$  is represented in Figure 7.





**Figure 7.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples,  $n=30$ , drawn from the Pareto distribution with  $x_0 = 1$ ,  $\theta = 3$ )

The range of values for  $\hat{\Delta}^{sim}$  is really wide: from a minimum value of 0.18221 to a maximum of 6.92917. The long right tail is confirmed also by the other sample sizes, shown in Figure 8 below, even if this behaviour is expected to decrease as the sample size  $n$  increases.

As before, simulation data for different values of the sample size  $n$  were obtained. They are shown in detail in Table 12.

**Table 12.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$  sampling from the Pareto ( $x_0 = 1$ ,  $\theta = 3$ ) (by  $B = 10,000$  samples)

$1-\alpha$	Acc. range	$n = 60$		$n = 120$		$n = 240$		$n = 480$		$N = 960$	
		$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\bar{X}_n}$
0.99	0.988-0.992	0.982	0.984	0.982	0.983	0.984	0.986	0.987	<b>0.988</b>	0.988	<b>0.988</b>
0.975	0.970-0.980	<b>0.974</b>	<b>0.976</b>	<b>0.973</b>	<b>0.975</b>	<b>0.976</b>	<b>0.977</b>	<b>0.976</b>	<b>0.977</b>	<b>0.977</b>	<b>0.977</b>
0.95	0.940-0.960	0.967	0.964	0.963	<b>0.956</b>	0.964	0.962	<b>0.959</b>	<b>0.958</b>	<b>0.958</b>	<b>0.955</b>
0.925	0.910-0.940	0.957	0.952	0.950	0.941	0.948	0.945	0.942	<b>0.936</b>	<b>0.937</b>	<b>0.934</b>
0.90	0.880-0.920	0.949	0.939	0.935	0.925	0.933	0.922	0.924	<b>0.914</b>	<b>0.917</b>	<b>0.911</b>
0.85	0.820-0.880	0.925	0.900	0.902	0.885	0.895	<b>0.880</b>	0.886	<b>0.871</b>	<b>0.873</b>	<b>0.865</b>

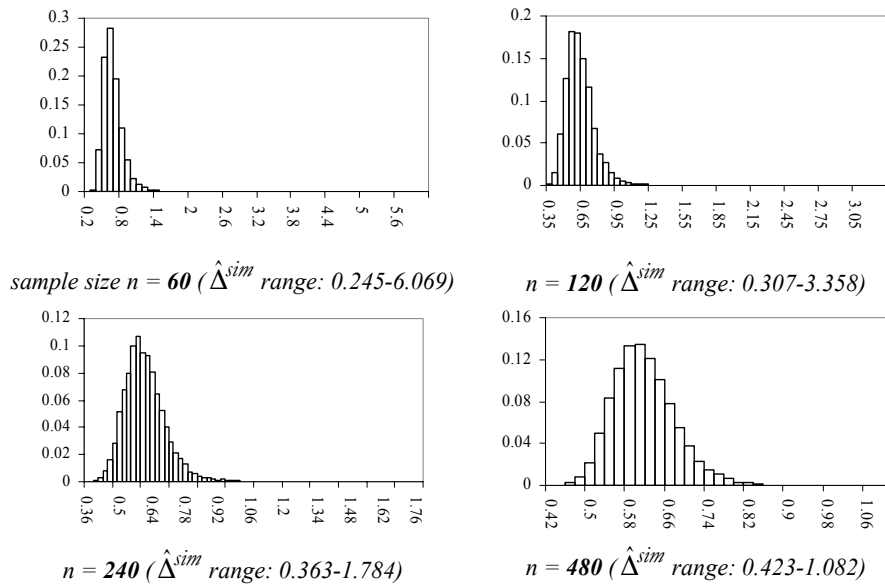
The simulated probabilities are often out from the accepted ranges determined by the values of  $\alpha$ , but the situation is somewhat better for  $\hat{p}_{\bar{X}_n}$  than for  $\hat{p}_{\hat{\Delta}_n}$ . Both cases show some slight improvements as the sample size  $n$  increases.

A summary of some moments and indices of the simulated distribution is reported in Table 13, providing also their theoretical or asymptotic values for meaningful comparisons:

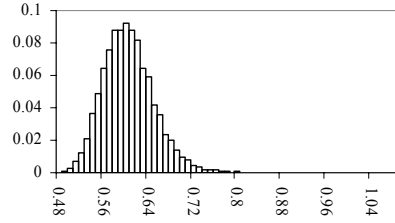
**Table 13.** Values of some indices of the asymptotic distribution of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B=10,000$  s.w.r. from the Pareto distribution ( $x_0=1, \theta=3$ ) and the theoretical values of  $\sqrt{\text{Var}(\hat{\Delta}_n)}$

Simul.	$n = 30$	$n = 60$	$n = 120$	$n = 240$	$n = 480$	$n = 960$	Asympt.values
$M_1(\hat{\Delta}^{sim})$	0.59570	0.59816	0.60031	0.59929	0.60006	0.59925	$E(\hat{\Delta}_n) = 0.6$
$Me(\hat{\Delta}^{sim})$	0.54532	0.56664	0.58115	0.58846	0.59355	0.59575	0.6
$\sigma(\hat{\Delta}^{sim})$	0.26512	0.18952	0.13269	0.09111	0.06406	0.04591	
$\sqrt{\text{Var}(\hat{\Delta}_n)}$	0.26340	0.18609	0.13153	0.09298	0.06574	0.04649	

Because of the choice of the parameters for the underlying Pareto distribution, the third and the fourth simulated moments of  $\hat{\Delta}_n$  do not appear in Table 13. Figure 8 shows that the simulated distributions  $\hat{\Delta}^{sim}$ , actually, are still simultaneously peaked and have a heavy right tail, also for higher sample sizes  $n = 60, 120, 240, 480$  and  $960$ . These characteristics weaken for increasing values of the sample size. It is interesting to observe also how the range of values of  $\hat{\Delta}^{sim}$  shrinks as  $n$  increases.



**Figure 8.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples, drawn from the Pareto ( $x_0=1, \theta=3$ ))



$n = 960$  ( $\hat{\Delta}^{sim}$  range: 0.471-1.081)

**Figure 8.** (continuation) Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples, drawn from the Pareto ( $x_0=1, \theta=3$ ))

The last distribution seems roughly to approach a good approximation to the normal curve, as the acceptable results about  $\hat{p}_{\hat{\Delta}_n}$  and  $\hat{p}_{\bar{X}_n}$  denote, in the last columns of Table 12.

**4.5 Pareto distribution (second case: with parameters  $x_0=2, \theta=4$ )**

For this c.r.v., fixing now the values of the parameters by  $x_0 = 2$  and  $\theta = 4$ , the main characteristics of the distribution are  $\mu = 2.66667$ ,  $\sigma^2 = 0.88889$  and  $\Delta = 0.76191$ . The third moment exists (but the fourth still does not exist), so we expect a better behaviour than in the previous case. The simulated data obtained for this distribution, choosing a sample size  $n = 30$ , are presented in Table 14.

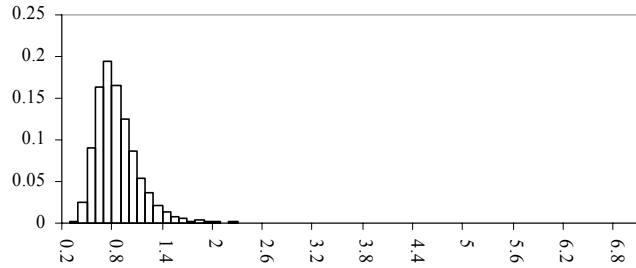
**Table 14.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$ , sampling from the Pareto with  $x_0 = 2, \theta = 4$  ( $n=30, B = 10,000$  samples)

$\hat{\Delta}_{\alpha/2}^{sim}$	$\hat{\Delta}_{1-\alpha/2}^{sim}$	$\hat{\Delta}_{\alpha/2}$	$\hat{\Delta}_{1-\alpha/2}$	$\hat{p}_{\hat{\Delta}_n}$	$1-\alpha$	Acc. range	$\bar{X}_{\alpha/2}^{sim}$	$\bar{X}_{1-\alpha/2}^{sim}$	$\bar{X}_{\alpha/2}$	$\bar{X}_{1-\alpha/2}$	$\hat{p}_{\bar{X}_n}$
0.335	1.822	0.087	1.437	0.979	0.99	0.988-0.992	2.344	3.284	2.223	3.110	0.983
0.369	1.567	0.175	1.349	<b>0.971</b>	0.975	0.970-0.980	2.372	3.162	2.281	3.052	<b>0.973</b>
0.398	1.386	0.249	1.275	<b>0.959</b>	0.95	0.940-0.960	2.399	3.058	2.329	3.004	<b>0.959</b>
0.419	1.292	0.296	1.228	0.948	0.925	0.910-0.940	2.415	3.007	2.360	2.973	0.942
0.436	1.230	0.331	1.193	0.936	0.90	0.880-0.920	2.429	2.972	2.384	2.950	0.924
0.463	1.144	0.385	1.139	0.905	0.85	0.820-0.880	2.449	2.919	2.419	2.914	0.882

Also in this case the values of the simulated quantiles  $\hat{\Delta}_{\alpha/2}^{sim}$  and  $\hat{\Delta}_{1-\alpha/2}^{sim}$  are higher than their asymptotic values  $\hat{\Delta}_{\alpha/2}$  and  $\hat{\Delta}_{1-\alpha/2}$  (see columns 1-4 of Table 14). The simulated probabilities  $\hat{p}_{\hat{\Delta}_n}$  are outside the range of

acceptable values indicated in the central column (except for the case of  $\alpha = 0.025$  and  $\alpha = 0.05$ , for which the good values are highlighted). More specifically  $\hat{p}_{\hat{\Delta}_n}$  is greater than the right threshold of acceptable ranges for  $\alpha > 0.05$ , and it is lower than the left threshold for  $\alpha = 0.01$ . This denotes a high asymmetry in the distribution of  $\hat{\Delta}_n$ , so that the tails of the distribution of  $\hat{\Delta}_n$  remarkably differ from the tails of a normal distribution. The results observed for  $\hat{p}_{\hat{\Delta}_n}$  are confirmed also by  $\hat{p}_{\hat{\bar{X}}_n}$ , indicating an analogous behaviour for the tails of  $\bar{X}_n$ .

The simulated distribution of  $\hat{\Delta}_n$  for  $n = 30$  is graphically represented in Figure 9.



**Figure 9.** Simulated distribution of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples,  $n=30$ , drawn from the Pareto distribution with  $x_0 = 1$ ,  $\theta = 3$ )

The values for  $\hat{\Delta}_n^{sim}$  vary in a wide range: from a minimum value of 0.20564 to a maximum of 7.01243. The heavy right tail is confirmed by the asymmetry index, shown in Table 16, evaluated for each chosen sample size  $n$ . Clearly, we expect that this behaviour will weaken as  $n$  increases.

As before, simulation data for other values of the sample size  $n$  are shown in detail in Table 15.

**Table 15.** Simulation data for the evaluation of the tails of the distribution of the s.m.d.  $\hat{\Delta}_n$  sampling from the Pareto distr. with  $x_0 = 2$ ,  $\theta = 4$  (by  $B = 10,000$  samples)

$1-\alpha$	Acc. range	$n = 60$		$n = 120$		$n = 240$		$n = 480$		$n = 960$	
		$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\hat{\bar{X}}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\hat{\bar{X}}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\hat{\bar{X}}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\hat{\bar{X}}_n}$	$\hat{p}_{\hat{\Delta}_n}$	$\hat{p}_{\hat{\bar{X}}_n}$
0.99	0.988-0.992	0.981	0.983	0.984	0.985	<b>0.988</b>	<b>0.990</b>	0.986	<b>0.988</b>	<b>0.988</b>	<b>0.989</b>
0.975	0.970-0.980	<b>0.971</b>	<b>0.972</b>	<b>0.972</b>	<b>0.974</b>	<b>0.977</b>	<b>0.977</b>	<b>0.974</b>	<b>0.977</b>	<b>0.973</b>	<b>0.974</b>
0.95	0.940-0.960	<b>0.957</b>	<b>0.953</b>	<b>0.958</b>	<b>0.953</b>	<b>0.958</b>	<b>0.958</b>	<b>0.953</b>	<b>0.955</b>	<b>0.949</b>	<b>0.948</b>
0.925	0.910-0.940	0.941	<b>0.935</b>	0.941	<b>0.934</b>	<b>0.940</b>	<b>0.934</b>	<b>0.931</b>	<b>0.929</b>	<b>0.926</b>	<b>0.924</b>
0.90	0.880-0.920	0.927	<b>0.915</b>	0.921	<b>0.915</b>	0.921	<b>0.912</b>	<b>0.910</b>	<b>0.903</b>	<b>0.903</b>	<b>0.901</b>

0.85 0.820-0.880 0.890 0.870 0.881 0.871 0.874 **0.865 0.862 0.854 0.853 0.851**

The simulated probabilities are better than those evaluated for sampling from the first case of Pareto distribution; they are out of the accepted ranges roughly only in 75% of all cases determined by a different value of  $\alpha$  and  $n$ . The situation is somewhat better for  $\hat{p}_{\bar{X}_n}$  than for  $\hat{p}_{\Delta_n}$ . Both improve as the sample size  $n$  increases.

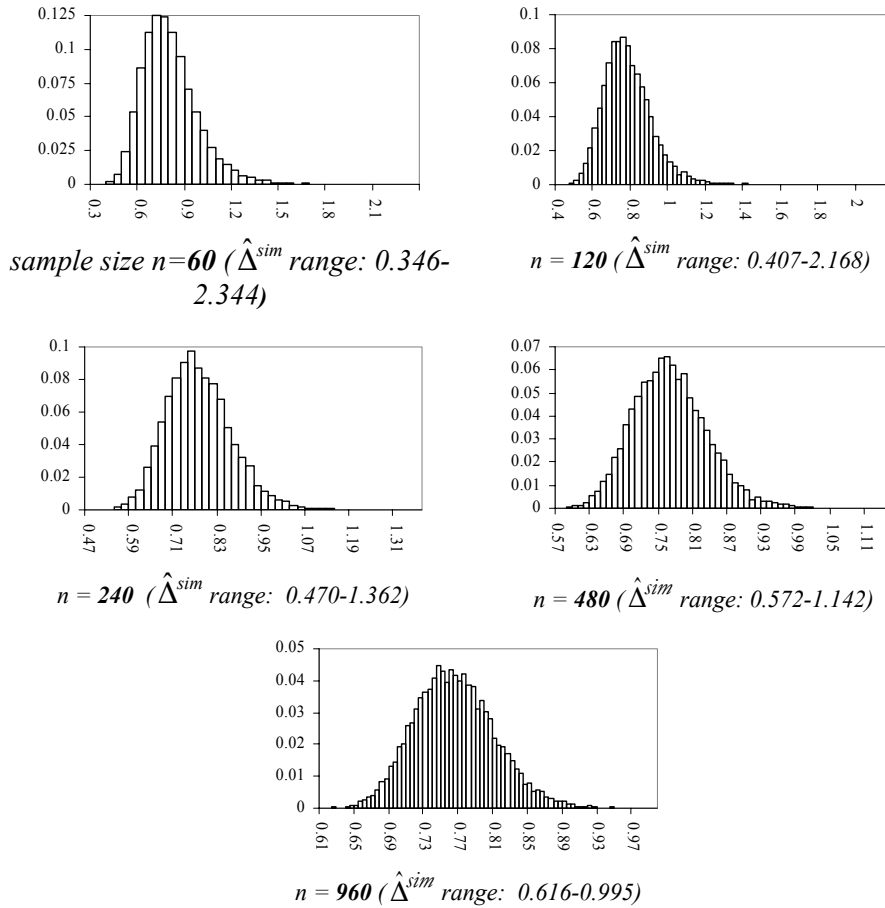
Trying to summarize the observed results on  $\hat{\Delta}_n$  for all values of  $\alpha$  and  $n$  and to compare the influence of heavier tails in the underlying distribution, for the Pareto with right tail given by  $x^{-4}$ ,  $\hat{p}_{\Delta_n}$  has only 11 acceptable values (cfr. Tables 11-12) while for the Pareto with right tail given by  $x^{-5}$ ,  $\hat{p}_{\Delta_n}$  achieve 21 good values (see Tables 14-15). Analogously, the same behaviour is observed on the tails of  $\bar{X}_n$ , as  $\hat{p}_{\bar{X}_n}$  has only 18 good evaluations in the first Pareto case and 28 acceptable values reported in this section.

A synthetic overview of the main characteristics of the simulated distribution is given in Table 16, compared to their theoretical or asymptotic values.

**Table 16.** Values of some indices of the asymptotic distribution of the s.m.d.  $\hat{\Delta}_n$ , compared with the simulated values obtained by  $B = 10,000$  s.w.r. from the Pareto ( $x_0 = 2, \theta = 4$ ) and the theoretical values for  $\sqrt{\text{Var}(\hat{\Delta}_n)}$

Simul. values	$n = 30$	$n = 60$	$n = 120$	$n = 240$	$n = 480$	$n = 960$	Asympt. values.
$M_I(\hat{\Delta}^{sim})$	0.76174	0.76242	0.76229	0.76183	0.76147	0.76123	$E(\hat{\Delta}_n)=0.76191$
$Me(\hat{\Delta}^{sim})$	0.71468	0.73524	0.74758	0.75369	0.75702	0.75912	0.76191
$\alpha_3(\hat{\Delta}^{sim})$	2.76495	1.25334	0.97897	0.61670	0.49237	0.31965	0.0
$\sigma(\hat{\Delta}^{sim})$	0.26833	0.18370	0.12912	0.09006	0.06521	0.04647	
$\sqrt{\text{Var}(\hat{\Delta}_n)}$	0.26193	0.18494	0.13068	0.09237	0.06531	0.04617	

The sign of the third simulated moments of  $\hat{\Delta}_n$  denotes positive asymmetry, even if they decrease as the sample size augments. The simulated distributions of  $\hat{\Delta}_n$  are simultaneously peaked and have a long right tail, as one can see in Figure 10, for sample sizes  $n = 60, 120, 240, 480$  and  $960$ , in which also the range of values of  $\hat{\Delta}^{sim}$  is indicated.



**Figure 10.** Simulated distributions of the s.m.d.  $\hat{\Delta}_n$  (by 10,000 samples, drawn from the Pareto ( $x_0 = 2$ ,  $\theta = 4$ ))

The last distribution seems roughly to approach a good approximation with the normal curve, as the acceptable results about  $\hat{p}_{\bar{x}_n}$  and  $\hat{p}_{\hat{\Delta}_n}$  indicate, in the last columns of Table 15.

Other simulations were done with different values of the parameter  $\theta$ : the high asymmetry and the long right tail in the distribution of the s.m.d.  $\hat{\Delta}_n$  gradually decreases as the value of  $\theta$  increases.

## 5. Concluding remarks

A simulation study has been developed to explore the minimum sample sizes required for a good approximation of the sample mean difference distribution to the normal curve. In general, how large  $n$  has to be for a good approximation depends on the population distribution from which the samples are drawn. This work shows that the shape of the underlying distribution can be very critical to construct confidence intervals.

In particular, when samples are drawn from symmetrical continuous variables, some very slight asymmetry in the distribution of the sample mean difference does not affect the desirable behaviour of its tails and the normal approximation is good, already for sample sizes  $n = 30$ , for all considered values of the risk  $\alpha$ , from 0.01 to 0.15.

Sampling from the Exponential distribution, i.e. considering an asymmetric distribution, still leads to low minimum required sample sizes:  $n = 30$  assures a good approximation to the tails of the Normal distribution for the sample mean difference for all chosen values of the nominal confidence.

For these three continuous distributions, a substantial agreement is observed between the behaviour of the sample mean difference and that of the sample mean, with reference to the goodness of the tails approximation by their asymptotic distribution.

Conversely, in sampling from the Pareto distribution that does not possess all finite moments, the sample mean difference presents a different behaviour from that of the sample mean: the minimum sample sizes assuring a good approximation for the sample mean are strongly lower than those needed for the sample mean difference. Remarkably, the values of the risk  $\alpha$  seriously affect the results about the sample mean difference. Actually, good probability coverage is attained, for sample sizes  $n > 480$ , only for the Pareto distribution that possesses the third moment.

Naturally we do not have enough elements to generalize these results, but surely this topic deserves further investigation, by the analysis of more simulations coming from a wider set of continuous distributions. A deeper understanding about the effect that asymmetry and heavy tails (in the underlying distribution) carry on the distribution of the sample mean difference is still needed.

Further analysis can be carried out considering that, in a non parametric context, the variance of  $\hat{\Delta}_n$  is not known and must hence be estimated. In this case the unbiased estimator  $\hat{V}ar(\hat{\Delta}_n)$  can be used. As one can show, the studentized  $\hat{\Delta}_n$  still converges to the normal, and it will be the aim of future work.

### Acknowledgements

The authors would like to thank two anonymous referees for their valuable comments and suggestions. This work is partially supported by F.A.R. 2005 - Università degli Studi di Milano-Bicocca.

### References

- Ali M.M. (1969). Distribution of linear combination of order statistics from rectangular population. *Bull.Inst. Statist. Res.*, tr. **3**, 1-21.
- Cowell F.A. (1989). Sampling variance and decomposable inequality measures. *Journal of Econometrics*, **42**, 27-41.
- Crocetta C., Loperfido N. (2005). The exact sampling distribution of L - statistics. *Metron*, **63**, 2, 213-224.
- Gon J. (2004). La distribuzione campionaria della differenza media di Gini. *Tesi di laurea*, Facoltà di Economia, Università degli Studi di Milano-Bicocca, Milano.
- Halmos P.R. (1946). The theory of unbiased estimation. *Annals of Mathematical Statistics*, **17**, 34-43.
- Hoefding W. (1948). A Class of Statistics with Asymptotically Normal Distribution. *The Annals of Mathematical Statistics*, **19**, No.3, 293-325.
- Kendall M.G., Stuart A., Ord J.K. (1994). *Distribution theory*. Halsted press, New York.
- Koroljuk V.S., Borovskich Yu.V. (1994) *Theory of U-Statistics*. Kluwer Academic Publishers, Dordrecht.
- Lomnicki Z.A. (1952). The standard error of Gini's mean difference. *The Annals of Mathematical Statistics*, **23**, No.4, 635-637.
- Michetti B., Dall'aglio G. (1957). La differenza semplice media. *Statistica*, **17**, 159-254.
- Nair U.S. (1936). The standard error of Gini's mean difference. *Biometrika*, **28**, 428-436.
- Vesserau A. (1957). Sur les conditions de l'application du criterium  $X^2$  de Poisson. Acte de la XXX<sup>e</sup> Session de l'I.I.S. (Stockolm) in *Bulletin de l'Institut International de Statistique*, **36** (3<sup>e</sup> livraison).
- Zenga M. (1974). L'approssimazione della distribuzione della statistica  $X^2$  di K.Pearson con la distribuzione di Chi-Quadrato. *Rivista Internazionale di Scienze sociali*. Anno LXXXII, 545-556.
- Zenga M., Poliscchio M., Greselin, F. (2004). The variance of Gini's mean difference and its estimators. *Statistica*, **3**, 2004, 455-475.