## Near-consistent robust estimations of moments for unimodal distributions

Tuban Leea,1

10

11

12

13

14

15

18

19

21

22

23

26

27

28

29

31

33

35

<sup>a</sup>Institute of Biomathematics, Macau SAR 999078, China

This manuscript was compiled on March 21, 2023

Descriptive statistics for parametric models currently heavily rely on the accuracy of distributional assumptions. Here, based on the invariant structures of unimodal distributions, a series of sophisticated yet efficient estimators, robust to both gross errors and departures from parametric assumptions, are proposed for estimating mean and central moments with insignificant asymptotic biases for common unimodal distributions. This article also illuminates the understanding of the common nature of probability distributions and the measures of them.

orderliness | invariant | unimodal | adaptive estimation | U-statistics

he asymptotic inconsistencies between sample mean  $(\bar{x})$ and nonparametric robust location estimators in asymmetric distributions on the real line have been noticed for more than two centuries (1), yet remain unsolved. Strictly speaking, it is unsolvable as by trimming, some information about the original distribution is removed, making it impossible to estimate the values of the removed parts without distributional assumptions. Newcomb (1886, 1912) provided the first modern approach to robust parametric estimation by developing a class of estimators that gives "less weight to the more discordant observations" (2, 3). In 1964, Huber (4) used the minimax procedure to obtain M-estimator for the contaminated normal distribution, which has played a pre-eminent role in the later development of robust statistics. However, as previously demonstrated, under growing asymmetric departures from normality, the bias of the Huber M-estimator increases rapidly. This is a common issue in parameter estimations. For example, He and Fung (1999) constructed (5) a robust M-estimator for the two-parameter Weibull distribution, from which all moments can be calculated. Nonetheless, it is inadequate for the gamma, Perato, lognormal, and the generalized Gaussian distributions (SI Dataset S1). Another old and interesting approach is arithmetically computing the parameters using one or more L-statistics as inputs, such as percentile estimators. Examples of percentile estimators for the Weibull distribution, the reader is referred to Menon (1963) (6), Dubey (1967) (7), Hassanein (1971) (8), Marks (2005) (9), and Boudt, Caliskan, and Croux (2011) (10)'s works. At the outset of the study of percentile estimators, it was known that they arithmetically utilizes the invariant structures of probability distributions (6, 11, 12). Maybe such estimators can be named as I-statistics. Formally, an estimator is classified as an *I*-statistic if it asymptotically satisfies  $I(LE_1, \dots, LE_l) = (\theta_1, \dots, \theta_q)$  for the distribution it is consistent, where LEs are calculated with the use of L-statistics, I is defined using arithmetic operations and constants, but it may also incorporate other functions, and  $\theta$ s are the population parameters it estimates. A subclass of I-statistics, arithmetic I-statistics, is defined as LEs are L-statistics, I is solely defined using arithmetic operations and constants.

Since some percentile estimators use the logarithmic function to transform all random variables before computing the L-statistics, a percentile estimator might not always be an arithmetic I-statistic (7). In this article, two subclasses of *I*-statistics are introduced, arithmetic *I*-statistics and quantile I-statistics. Examples of quantile I-statistics will be discussed later. Based on L-statistics, I-statistics are naturally robust. Compared to probability density functions (pdfs) and cumulative distribution functions (cdfs), the quantile functions of many parametric distributions are more elegant. Since the expectation of an L-statistic can be expressed as an integral of the quantile function, I-statistics are often analytically obtainable. However, the performance of the aforementioned examples is often worse than that of the robust M-statistics when the distributional assumption is violated (SI Dataset S1). Even when distributions such as the Weibull and gamma belong to the same larger family, the generalized gamma distribution, a misassumption can still result in substantial biases, rendering the approach ill-suited.

41

42

43

47

48

49

50

51

52

53

54

55

56

57

58

In previous research on semiparametric robust mean estimation, the binomial mean  $(BM_{\epsilon})$  is still inconsistent for any skewed distribution, despite having much smaller asymptotic biases than other weighted averages. All robust location estimators commonly used are symmetric due to the universality of the symmetric distributions. One can construct an asymmetric weighted average that is consistent for a semiparametric class of skewed distributions. This approach has been investigated previously, but its lack of symmetry makes it suitable only for certain applications (13). Shifting from semiparametrics to parametrics, an ideal robust location estimator would have a non-sample-dependent breakdown point (defined in Subsection ??) and be consistent for any symmetric distribution and a skewed distribution with finite second moments. This is called an invariant mean. Based on the mean-symmetric weighted

## **Significance Statement**

Bias, variance, and contamination are the three main errors in statistics. Consistent robust estimation is unattainable without parametric assumptions. Here, based on a paradigm shift inspired by mean-median-mode inequality, Bickel-Lehmann spread, and adaptive estimation, invariant moments are proposed as a means of achieving near-consistent and robust estimations of moments, even in scenarios where moderate violations of distributional assumptions occur, while the variances are sometimes smaller than those of the sample moments.

T.L. designed research, performed research, analyzed data, and wrote the paper. The author declares no competing interest.

<sup>&</sup>lt;sup>1</sup> To whom correspondence should be addressed. E-mail: tl@biomathematics.org

average-median inequality, the recombined mean is defined as

$$rm_{d,\epsilon,n} := \lim_{c \to \infty} \left( \frac{(\text{SWA}_{\epsilon,n} + c)^{d+1}}{\left( median + c \right)^d} - c \right),$$

where d is the key factor for bias correction,  $SWA_{\epsilon,n}$  is  $BM_{\epsilon,n}$  in the first three Subsections, but other symmetric weighted averages can also be used in practice as long as the inequalities hold. The following theorem shows the significance of this arithmetic I-statistic.

Theorem .1. If the second moments are finite,  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$  is a consistent mean estimator for the exponential and any symmetric distributions and the Pareto distribution with quantile function  $Q(p)=x_m(1-p)^{-\frac{1}{\alpha}},\,x_m>0,\,$  when  $\alpha\to\infty$ .

*Proof.* Finding d and  $\epsilon$  that make  $rm_{d,\epsilon}$  a consistent mean estimator is equivalent to finding the solution of  $E\left[rm_{d,\epsilon}\right] = E\left[X\right]$ . Rearranging the definition,  $rm_{d,\epsilon} = \lim_{c\to\infty} \left(\frac{(\mathrm{BM}_{\epsilon}+c)^{d+1}}{(median+c)^d} - c\right) = (d+1)\,\mathrm{BM}_{\epsilon} - d\mathrm{median} = \mu$ . So,  $d = \frac{\mu - \mathrm{BM}_{\epsilon}}{\mathrm{BM}_{\epsilon} - median}$ . The quantile function of the exponential distribution is  $Q(p) = \ln\left(\frac{1}{1-p}\right)\lambda$ .  $E\left[x\right] = \lambda$ .  $E\left[median\right] = \frac{1}{2}$  $Q\left(\frac{1}{2}\right) = \ln 2\lambda$ . For the exponential distribution, the expectation of  $\mathrm{BM}_{\frac{1}{8}}$  is  $E\left[\mathrm{BM}_{\frac{1}{8}}\right]=\lambda\left(1+\ln\left(\frac{46656}{8575\sqrt{35}}\right)\right)$ . Obviously, the scale parameter  $\lambda$  can be canceled out,  $d\approx0.375$ . The proof of the second assertion follows directly from the coincidence property. For any symmetric distribution with a finite second moment,  $E[BM_{\epsilon}] = E[median] = E[X]$ . Then  $E\left[rm_{d,\epsilon}\right] = \lim_{c\to\infty} \left(\frac{(E[X]+c)^{d+1}}{(E[X]+c)^d} - c\right) = E\left[X\right]$ . The proof for the Pareto distribution is more general. The mean of the Pareto distribution is given by  $\frac{\alpha x_m}{\alpha - 1}$ . The d value with two unknown percentiles  $p_1$  and  $p_2$  for the Pareto distribution is  $d_{Perato} = \frac{\frac{\alpha x_m}{\alpha - 1} - x_m (1 - p_1)^{-\frac{1}{\alpha}}}{x_m (1 - p_1)^{-\frac{1}{\alpha}} - x_m (1 - p_2)^{-\frac{1}{\alpha}}}.$  Since any weighted average  $\frac{1}{x_m} \frac{1}{x_m} \frac{1}{$ age can be expressed as an integral of the quantile function,  $\lim_{\alpha\to\infty}\frac{\frac{\alpha}{\alpha-1}-(1-p_1)^{-1/\alpha}}{(1-p_1)^{-1/\alpha}-(1-p_2)^{-1/\alpha}}=-\frac{\ln(1-p_1)+1}{\ln(1-p_1)-\ln(1-p_2)},$  the d value for the Pareto distribution approaches that of the exponential distribution as  $\alpha \to \infty$ , regardless of the type of weighted average used. This completes the demonstration.  $\Box$ 

Theorem .1 implies that for the Weibull, gamma, Pareto, lognormal and generalized Gaussian distribution,  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$  is consistent for at least one particular case. The biases of  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$  for distributions with skewness between those of the exponential and symmetric distributions are tiny (SI Dataset S1).  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$  exhibits excellent performance for all these common unimodal distributions (SI Dataset S1).

Besides introducing the concept of invariant mean, the purpose of this paper is to demonstrate that, in light of previous works, the estimation of central moments can be transformed into a location estimation problem by using U-statistics, the central moment kernel distributions possess desirable properties, and a series of sophisticated yet efficient robust estimators can be constructed whose biases are typically smaller than the variances (as seen in Table  $\ref{Table 1}$  for n=5400) for unimodal distributions.

## **Background and Main Results**

**A. Invariant mean.** It has long been known that a theoretical model can be adjusted to fit the first two moments of the observed data. A continuous distribution belonging to a location–scale family takes the form  $F(x) = F_0\left(\frac{x-\mu}{\lambda}\right)$ , where  $F_0$  is a "standard" distribution. Therefore,  $F(x) = Q^{-1}(x) \to x = Q(p) = \lambda Q_0(p) + \mu$ . Thus, any weighted average can be expressed as  $\lambda \mathrm{WA}_0(\epsilon) + \mu$ , where  $\mathrm{WA}_0(\epsilon)$  is an integral of  $Q_0(p)$  according to the definition of the weighted average. The substitution of the weighted average. The invariant consequently, the roles of  $\mathrm{BM}_\epsilon$  and median in  $rm_{d,\epsilon}$  can be replaced by any weighted averages, although only symmetric weighted averages are considered in defining the invariant mean.

The performance in heavy-tailed distributions can be further improved by constructing the quantile mean as

$$qm_{d,\epsilon,n} \coloneqq \hat{Q}_n\left(\left(\hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right) - \frac{1}{2}\right)d + \hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right)\right),$$

provided that  $\hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right) \geq \frac{1}{2}$ , where  $\hat{F}_n\left(x\right)$  is the empirical cumulative distribution function of the sample,  $\hat{Q}_n$  is the sample quantile function. The most popular method for computing the sample quantile function was proposed by Hyndman and Fan in 1996 (14). To minimize the finite sample bias, here,  $\hat{F}_n\left(x\right) \coloneqq \frac{1}{n}\left(\frac{x-X_{sp}}{X_{sp+1}-X_{sp}}+sp\right)$ , where  $sp=\sum_{i=1}^n 1_{X_i\leq x}, 1_A$  is the indicator of event A. The solution of  $\hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right)<\frac{1}{2}$  is reversing the percentile by  $1-\hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right)$ , the obtained percentile is also reversed. Without loss of generality, in the following discussion, only the case where  $\hat{F}_n\left(\mathrm{SWA}_{\epsilon,n}\right)\geq\frac{1}{2}$  is considered. Moreover, in extreme heavy-tailed distributions, the calculated percentile can exceed the breakdown point of  $\mathrm{SWA}_{\epsilon}$ , so the percentile will be modified to  $1-\epsilon$  if this occurs. The quantile mean uses the location-scale invariant in a different way as shown in the following proof.

**Theorem A.1.**  $qm_{d\approx 0.321,\epsilon=\frac{1}{8}}$  is a consistent mean estimator for the exponential, Pareto  $(\alpha\to\infty)$  and any symmetric distributions provided that the second moments are finite.

Proof. Similarly, rearranging the definition,  $d=\frac{F(\mu)-F(\mathrm{BM}_\epsilon)}{F(\mathrm{BM}_\epsilon)-\frac{1}{2}}$ . The cdf of the exponential distribution is  $F(x)=1-e^{-\lambda^{-1}x}$ ,  $\lambda\geq 0,\ x\geq 0$ , the expectation of  $\mathrm{BM}_\epsilon$  can be expressed as  $\lambda\mathrm{BM}_0(\epsilon)$ , so  $F(\mathrm{BM}_\epsilon)$  is free of  $\lambda$ . When  $\epsilon=\frac{1}{8},\ d=\frac{-e^{-1}+e^{-\left(1+\ln\left(\frac{46656}{8575\sqrt{35}}\right)\right)}}{\frac{1}{2}-e^{-\left(1+\ln\left(\frac{46656}{8575\sqrt{35}}\right)\right)}}\approx 0.321$ . The proof of the symmetric case is similar. Since for any symmetric distribution with a finite second moment,  $F(E[\mathrm{BM}_\epsilon])=F(\mu)=\frac{1}{2}$ .

with a finite second moment,  $F(E[\mathrm{BM}_{\epsilon}]) = F(\mu) = \frac{1}{2}$ . Then, the expectation of the quantile mean is  $qm_{d,\epsilon} = F^{-1}\left(\left(F(\mu) - \frac{1}{2}\right)d + F(\mu)\right) = F^{-1}\left(0 + F(\mu)\right) = \mu$ .

For the assertion related to the Pareto distribution, the cdf of it is  $1 - \left(\frac{x_m}{x}\right)^{\alpha}$ . So, the d value with two unknown percentile  $p_1$  and  $p_2$  is

$$d_{Pareto} = \frac{1 - \left(\frac{x_m}{\alpha x_m}\right)^{\alpha} - \left(1 - \left(\frac{x_m}{x_m(1-p_1)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right)}{\left(1 - \left(\frac{x_m}{x_m(1-p_1)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right) - \left(1 - \left(\frac{x_m}{x_m(1-p_2)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right)} = 152$$

 $\frac{1-\left(\frac{\alpha-1}{\alpha}\right)^{\alpha}-p_1}{p_1-p_2}$ . When  $\alpha \to \infty$ ,  $\left(\frac{\alpha-1}{\alpha}\right)^{\alpha}=\frac{1}{e}$ . The d value for the exponential distribution is identical, since  $d_{exp}=$ 

2 |

$$\frac{\left(1-e^{-1}\right)-\left(1-e^{-\ln\left(\frac{1}{1-p_1}\right)}\right)}{\left(1-e^{-\ln\left(\frac{1}{1-p_1}\right)}\right)-\left(1-e^{-\ln\left(\frac{1}{1-p_2}\right)}\right)}=\frac{1-\frac{1}{e}-p_1}{p_1-p_2}. \text{ All results}$$
156 are now proven.

157

158

159

160

161

162

163

164

165

166

167

169

170

172

173

174

175

177

178

180

181

182

183

184

185

186

187

188

189

191

192

193

195

196

197

198

199

200

201

204

205

206

207

210

The definitions of location and scale parameters are such that they must satisfy  $F(x; \lambda, \mu) = F(\frac{x-\mu}{\lambda}; 1, 0)$ . By recalling  $x = \lambda Q_0(p) + \mu$ , it follows that the percentile of any weighted average is free of  $\lambda$  and  $\mu$ , which guarantees the validity of the quantile mean. The quantile mean is a quantile I-statistic. Specifically, an estimator is classified as a quantile I-statistic if LEs are percentiles of a distribution obtained by plugging L-statistics into a cumulative distribution function and I is defined with arithmetic operations, constants and quantile functions.  $qm_{d\approx 0.321,\epsilon=\frac{1}{2}}$  works better in the fat-tail scenarios (SI Dataset S1). Theorem .1 and A.1 show that  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$ and  $qm_{d\approx 0.321,\epsilon=\frac{1}{2}}$  are both consistent mean estimators for any symmetric distribution and a skewed distribution with finite second moments. It's obvious that the breakdown points of  $rm_{d\approx 0.375,\epsilon=\frac{1}{8}}$  and  $qm_{d\approx 0.321,\epsilon=\frac{1}{9}}$  are both  $\frac{1}{8}$ . Therefore they are all invariant means.

To study the impact of the choice of SWAs in rm and qm, it is constructive to recall that a symmetric weighted average is a linear combination of symmetric quantile averages. While using a less-biased symmetric weighted average can generally enhance performance (SI Dataset S1), there is a greater risk of violation in the semiparametric framework. However, the mean-SWA-median inequality is robust to slight fluctuations of the SQA function of the underlying distribution. Suppose the SQA function is generally decreasing in [0, u], but increasing in  $[u, \frac{1}{2}]$ , since  $1-2\epsilon$  of the symmetric quantile averages will be included in the computation of SWA<sub> $\epsilon$ </sub>, as long as  $\frac{1}{2}-u \ll 1-2\epsilon$ , and other portions of the SQA function satisfy the inequality constraints that define the  $\nu$ th orderliness on which the SWA $_{\epsilon}$ is based, the mean-SWA<sub> $\epsilon$ </sub>-median inequality will still hold. This is due to the violation being bounded (15) and therefore cannot be extreme for unimodal distributions. For instance, the SQA function is non-monotonic when the shape parameter of the Weibull distribution  $\alpha>\frac{1}{1-\ln(2)}\approx 3.259$  as shown in the previous article, the violation of the third orderliness starts near this parameter as well, yet the mean-BM  $_{\frac{1}{a}}\text{-median}$ inequality is still valid when  $\alpha \leq 3.322$ . Another key factor in determining the risk of violation is the skewness of the distribution. Previously, it was demonstrated that in a family of distributions differing by a skewness-increasing transformation in van Zwet's sense, the violation of orderliness, if it happens, often only occurs when the distribution is nearly symmetrical (16). The over-corrections in rm and qm are dependent on the  $SWA_{\epsilon}$ -median difference, which can be a reasonable measure of skewness (17, 18), implying that the over-correction is often tiny with a moderate d. This qualitative analysis provides another perspective, in addition to the bias bounds (15), that rm and qm based on the mean-SWA<sub> $\epsilon$ </sub>-median inequality are generally safe.

**B. Robust estimations of the central moments.** In 1979, Bickel and Lehmann, in their final paper of the landmark series *Descriptive Statistics for Nonparametric Models* (19), generalized a class of estimators called "measures of spread," which "does not require the assumption of symmetry." From that, a popular

efficient scale estimator, the Rousseeuw-Croux scale estimator (20), was derived in 1993, but the importance of tackling the symmetry assumption has been greatly underestimated. While they had already considered one version of the trimmed standard deviation in the third paper of that series (21), in the final section of that paper (19), they explored another two possible versions, which were modified here for comparison,

$$\left[n\left(\frac{1}{2} - \epsilon\right)\right]^{-\frac{1}{2}} \left[\sum_{i=\frac{n}{2}}^{n(1-\epsilon)} \left[X_i - X_{n-i+1}\right]^2\right]^{\frac{1}{2}}, \qquad [1] \quad \text{21}$$

211

212

213

214

217

221

222

223

224

225

226

227

229

230

231

232

233

234

236

237

238

239

240

241

242

244

245

246

247

248

249

250

251

252

253

254

255

256

and 219

$$\left[ \binom{n}{2} \left( 1 - \epsilon - \gamma \epsilon \right) \right]^{-\frac{1}{2}} \left[ \sum_{i=\binom{n}{2}\gamma\epsilon}^{\binom{n}{2}(1-\epsilon)} \left( X - X' \right)_i^2 \right]^{\frac{1}{2}}, \quad [2] \quad {}_{220}$$

where  $(X - X')_1 \leq \ldots \leq (X - X')_{\binom{n}{2}}$  are the order statistics of the "pseudo-sample",  $X_i - X_j$ , i < j. The paper ended with, "We do not know a fortiori which of the measures [1] or [2] is preferable and leave these interesting questions open."

Observe that the kernel of the unbiased estimation of the second central moment by using U-statistic is  $\psi_2(x_1, x_2) = \frac{1}{2}(x_1 - x_2)^2$ . If adding the  $\frac{1}{2}$  term in [2], as  $\epsilon \to 0$ , the result is equivalent to the standard deviation estimated by using U-statistic (also noted by Janssen, Serfling, and Veraverbeke in 1987) (22). In fact, they also showed that, when  $\epsilon$  is 0, [2] is  $\sqrt{2}$  times the standard deviation.

To address their open question, the nomenclature used in this paper is introduced as follows:

Nomenclature. Given a robust estimator  $\hat{\theta}$  with an adjustable breakdown point which can be infinitesimal. The name of  $\hat{\theta}$  is composed of two parts: the first part denotes the type of estimator, and the second part is the name of the population parameter  $\theta$  that the estimator is consistent with as  $\epsilon \to 0$ . The abbreviation of the estimator is formed by combining the initial letter(s) of the first part with the common abbreviation of the consistent estimator that measures the population parameter. If the estimator is symmetric, the asymptotic breakdown point,  $\epsilon$ , is indicated in the subscript of the abbreviation of the estimator, except the median. For asymmetric estimators based on quantile average, the corresponding  $\gamma$  is also indicated after  $\epsilon$ . Note that  $\epsilon$  is the right breakdown point (defined in Subsection ??), while the left breakdown point should be further calculated.

In the previous article on semiparametric robust mean estimation, it was shown that the bias of a robust estimator with an adjustable breakdown point is often monotonic with respect to the breakdown point in a semiparametric distribution. Naturally, the estimator's name should correspond to the population parameter with which it is consistent as  $\epsilon \to 0$ . The trimmed standard deviation following this nomenclature

is 
$$\operatorname{Tsd}_{\epsilon=1-\sqrt{1-\epsilon_0},\gamma,n}:=\left[\operatorname{TM}_{\epsilon_0,\gamma}\left((\psi_2\left(X_{N_1},X_{N_2}\right))_{N=1}^{\binom{n}{2}}\right)\right]^{-\frac{1}{2}},$$
 where  $\operatorname{TM}_{\epsilon_0,\gamma}(Y)$  denotes the  $\epsilon_0,\gamma$ -trimmed mean with the sequence  $(\psi_2\left(X_{N_1},X_{N_2}\right))_{N=1}^{\binom{n}{2}}$  as an input, the proof of the breakdown point is given in Subsection ??. Removing the square root yields the trimmed variance  $(\operatorname{T}\!\mathit{var}_{\epsilon,\gamma,n})$ . It is now very clear that this definition, essentially the same as

[2], should be preferable. Not only because it is essentially a trimmed U-statistic for the standard deviation but also because the  $\gamma$ -orderliness of the second central moment kernel distribution is ensured by the next exciting theorem.

**Theorem B.1.** The second central moment kernel distribution generated from any unimodal distribution is  $\gamma$ -ordered.

Proof. The monotonic increasing of the pairwise difference distribution was first implied in its unimodality proof done by Hodges and Lehmann in 1954 (23). Whereas they used induction to get the result in Theorem ??, Dharmadhikari and Jogdeo in 1982 (24) provided a modern proof of the unimodality using Khintchine's representation (25). Assuming absolute continuity, Purkayastha (26) introduced a much simpler proof in 1998. Transforming the pairwise difference distribution by squaring and multiplying by  $\frac{1}{2}$  does not change the monotonicity, making the pdf become monotonically decreasing with mode at zero. In the previous article, it was proven that a right skewed distribution with a monotonic decreasing pdf is always  $\gamma$ -ordered, which gives the desired result.

Previously, it was shown that any symmetric distribution with a finite second moment is  $\nu$ th ordered, indicating that orderliness does not require unimodality, e.g., a symmetric bimodal distribution is also ordered. An analysis of the Weibull distribution showed that unimodality does not guarantee orderliness. Theorem B.1 reveals another profound relationship between unimodality and orderliness, which is sufficient for trimming inequality.

In 1928, Fisher constructed k-statistics as unbiased estimators of cumulants (27). Halmos (1946) proved that the functional  $\theta$  admits an unbiased estimator if and only if it is a regular statistical functional of degree k and showed a relation of symmetry, unbiasness and minimum variance (28). In 1948, Hoeffding generalized U-statistics (29) which enable the derivation of a minimum-variance unbiased estimator from each unbiased estimator of an estimable parameter. Heffernan (1997) (30) obtained an unbiased estimator of the kth central moment by using U-statistics and demonstrated that it is the minimum variance unbiased estimator for distributions with finite moments (31, 32). In 1976, Saleh generalized the Hodges-Lehmenn (H-L) estimator (33) to the trimmed H-L mean (which he named "Wilcoxon one-sample statistic") (34). In 1984, Serfling pointed out the speciality of Hodges-Lehmann estimator, which is neither a simple L-statistic nor a U-statistic, and considered the generalized L-statistics and Ustatistic structure (35). Also in 1984, Janssen and Serfling and Veraverbeke (36) showed that the Bickel-Lehmann spread also belongs to the same class. It gradually became clear that the Hodges-Lehmenn estimator, trimmed H-L mean and trimmed standard deviation are all trimmed U-statistics (37–39).

Extending the trimmed U-statistic to weighted U-statistic, i.e., replacing the trimmed mean with weighted average. The weighted kth central moment  $(k \leq n)$  is defined as,

$$Wkm_{\epsilon=1-(1-\epsilon_0)^{\frac{1}{k}},\gamma,n} := WA_{\epsilon_0,\gamma,n} \left( \left( \psi_k \left( X_{N_1}, \cdots, X_{N_k} \right) \right)_{N=1}^{\binom{n}{k}} \right),$$

where  $X_{N_1}, \dots, X_{N_k}$  are the n choose k elements from X,  $\psi_k(x_1, \dots, x_k) = \sum_{j=0}^{k-2} (-1)^j \left(\frac{1}{k-j}\right) \sum_{j=0}^{k-1} \left(x_{i_1}^{k-j} \dots x_{i_{(j+1)}}\right) + 3 \left(-1\right)^{k-1} (k-1) x_1 \dots x_k$ , the second summation is over  $x_1, \dots, x_{j+1} = 1$  to k with  $x_1 < \dots < x_{j+1}$  (30). Despite the

complexity, the structure of the kth central moment kernel distributions can be elucidated by decomposing.

**Theorem B.2.** For each pair  $(Q(p_i), Q(p_j))$  of the original distribution such that  $Q(p_i) < Q(p_j)$ , let  $x_1 = Q(p_i)$  and  $x_k = Q(p_j)$ ,  $\Delta = Q(p_i) - Q(p_j)$ , the kth central moment kernel distribution, k > 2, can be seen as a mixture distribution and each of the components has the support  $\left(-\left(\frac{k}{3+(-1)^k}\right)^{-1}(-\Delta)^k, \frac{1}{k}(-\Delta)^k\right)$ .

Proof. Without loss of generality, generating the distribution of the k-tuple  $(Q(p_{i_1}),\ldots,Q(p_{i_k}))$  under continuity,  $k>2,\ i_1<\ldots< i_k,\ p_{i_1}<\ldots< p_{i_k},$  the corresponding probability density is  $f_{X,\ldots,X}(Q(p_{i_1}),\ldots,Q(p_{i_k}))=k!f(Q(p_{i_1}))\ldots f(Q(p_{i_k})).$  Transforming the distribution of the k-tuple by the function  $\psi_k\left(x_1,\ldots,x_k\right)$ , denoting  $\bar{\Delta}=\psi_k\left(Q(p_{i_1}),\ldots,Q(p_{i_k})\right).$  The probability  $f_{\Xi_k}(\bar{\Delta})=\sum_{\bar{\Delta}=\psi_k\left(Q(p_{i_1}),\ldots,Q(p_{i_k})\right)}f_{X,\ldots,X}(Q(p_{i_1}),\ldots,Q(p_{i_k}))$  is the summation of the probabilities of all k-tuples such that  $\bar{\Delta}$  is equal to  $\psi_k\left(Q(p_{i_1}),\ldots,Q(p_{i_k})\right).$  The following  $\Xi_k$  is equivalent.

 $\Xi_k$ : Every pair with a difference equal to  $\Delta = Q(p_{i_1}) - Q(p_{i_k})$  can generate a pseudodistribution (but the integral is not equal to 1, so "pseudo") such that  $x_2, \ldots, x_{k-1}$  exhaust all combinations under the inequality constraints, i.e.,  $Q(p_{i_1}) = x_1 < x_2 < \ldots < x_{k-1} < x_k = Q(p_{i_k})$ . The combination of all the pseudodistributions with the same  $\Delta$  is  $\xi_{\Delta}$ . The combination of  $\xi_{\Delta}$ , i.e., from  $\Delta = 0$  to Q(0) - Q(1), is  $\Xi_k$ .

The support of  $\xi_{\Delta}$  is the extrema of  $\psi_k$  subject to the inequality constraints. Using the Lagrange multiplier, one can easily determine the only critical point at  $x_1=\ldots=x_k=0$ , where  $\psi_k=0$ . Other candidates are within the boundaries, i.e.,  $\psi_k\left(x_1=x_1,x_2=x_k,\cdots,x_k=x_k\right)$ , ...,  $\psi_k\left(x_1=x_1,\cdots,x_i=x_1,x_{i+1}=x_k,\cdots,x_k=x_k\right)$ , ...,  $\psi_k\left(x_1=x_1,\cdots,x_i=x_1,x_{i+1}=x_k,\cdots,x_{k-1}=x_1,x_k=x_k\right)$ , where  $\psi_k=0$  is the proof of the proof

4 |

370  $(k-i)\left(\frac{(-1)^k}{i-k} + (-1)^k\right) = (-1)^{k+1} + (k-i)(-1)^k$ 371 and  $\sum_{g=k-i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-k+i-1} = \frac{1}{2} \int_0^1 \sum_{g=k-i+1}^{k-1} (-1)^{g+1} i \binom{i-1}{g-k+i-1} t^{k-g} dt = \frac{1}{2} \int_0^1 \left(i(-1)^{k-i} (t-1)^{i-1} - i(-1)^{k+1}\right) dt = (-1)^k (i-1).$ 374 If  $j < \frac{k+1-i}{2}$ , i > k-1, if i = k,  $\psi_k = 0$ , if  $\frac{k+1-i}{2} \le j \le \frac{k-1}{2}$ ,  $\frac{k+1}{2} \le i \le k-1$ , the summed coefficient of  $x_1^i x_k^{k-i}$  is  $(-1)^{k-1} (k-1) + \sum_{g=k-i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-k+i-1} + \sum_{g=i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} (k-i) \binom{k-i-1}{g-i-1}$ , the same as above. If 377 i+j < k, since  $\binom{i}{k-j} = 0$ , the related terms can be ignored, so,  $\begin{aligned} i+j &< k, \text{ since } \binom{k-j}{j} = 0, \text{ the related terms can be ignored, so,} \\ \text{using the binomial theorem and beta function, the summed coefficient of } x_1^{k-j} x_k^j \text{ is } \sum_{g=j+1}^{i+j} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-j-1} \binom{k-i}{j} = i \binom{k-i}{j} \int_0^1 \sum_{g=j+1}^{i+j} (-1)^{g+1} \binom{i-1}{g-j-1} t^{k-g} dt & = \binom{k-i}{j} i \int_0^1 \left( (-1)^j t^{k-j-1} \left( \frac{t}{t-1} \right)^{1-i} \right) dt & = \binom{k-i}{j} i \frac{(-1)^{j+i+1} \Gamma(i) \Gamma(k-j-i+1)}{\Gamma(k-j+1)} & = \frac{(-1)^{j+i+1} i! (k-j-i)! (k-i)!}{(k-j)! j! (k-j-i)!} & = \binom{k-i}{j!} i! \frac{i! (k-i)!}{k!} \frac{k!}{(k-j)! j!} & = \binom{k}{i!} i! \binom{k-j-i}{j!} i! & = \binom{k-i}{j!} i! \binom{k-j-i}{j!} i! & = \binom{k-i}{j!} i! & = \binom{k-i}{j!}$ The coefficient of  $x_1^i x_k^{k-i}$  in  $\binom{k}{i}^{-1} (-1)^{1+i} (x_1 - x_k)^k$  is  $\binom{k}{i}^{-1} (-1)^{1+i} \binom{k}{i} (-1)^{k-i} = (-1)^{k+1}$ , same as the 385 386 summed coefficient if i+j=k. If i+j< k, the coefficient of  $x_1^{k-j}x_k^j$  is  $\binom{k}{i}^{-1}(-1)^{1+i}\binom{k}{j}(-1)^j$ , same as the corresponding summed coefficient. There-387 389 fore,  $\psi_k (x_1 = x_1, \dots, x_i = x_1, x_{i+1} = x_k, \dots, x_k = x_k) = \binom{k}{i}^{-1} (-1)^{1+i} (x_1 - x_k)^k$ , the maximum and minimum of  $\psi_k$ 390 391 follow directly from the properties of the binomial coeffi-392 cient.

 $\xi_{\Delta}$  is closely related to  $f_{\Xi}(\Delta)$ , which is the pairwise difference distribution, since the probability density of  $\xi_{\Delta}$  can be exence distribution, since the pressed as  $f_{\Xi_k}(\bar{\Delta}|\Delta)$  and  $\sum_{\bar{\Delta}=-\left(\frac{3+(-1)^k}{2}\right)^{-1}(-\Delta)^k}^{\frac{1}{k}(-\Delta)^k}f_{\Xi_k}(\bar{\Delta}|\Delta)=$ 

 $f_{\Xi}(\Delta) = \int_{0}^{\infty} 2f(t) f(t-\Delta) dt$ . The support of the original distribution is assumed to be  $[0, \infty)$  for simplicity. Recall that  $f_{\Xi}(\Delta)$  is monotonic increasing with a mode at the origin if the original distribution is unimodal. Thus, in general, ignoring the shape of  $\xi_{\Delta}$ ,  $\Xi_k$  is monotonic left and right around zero. In fact, the median of  $\Xi_k$  also exhibits a strong tendency to be close to zero, as it can be cast as a weighted mean of the medians of  $\xi_{\Delta}$ . When  $\Delta$  is small, all values of  $\xi_{\Delta}$  are close to zero, resulting in the median of  $\xi_{\Delta}$  being close to zero as well. When  $\Delta$  is large, the median of  $\xi_{\Delta}$  depends on its skewness, but the corresponding weight is much smaller, so even if  $\xi_{\Delta}$  is highly skewed, the median of  $\Xi_k$  will only be slightly shifted from zero. Denote the median of  $\Xi_k$  as  $m_{\Xi_k}$ , for the five parametric distributions here,  $|m_{\Xi_k}|$ s are all  $\leq 0.1\sigma$  for  $\Xi_3$  and  $\Xi_4$  (SI Dataset S1). Assuming  $m_{\Xi_k} = 0$ , for the even ordinal central moment kernel distribution, the average probability density on the left side of zero is greater than that on the right side, since  $\frac{\frac{1}{2}}{{\binom{k}{2}}^{-1}(Q(0)-Q(1))^k} > \frac{\frac{1}{2}}{\frac{1}{k}(Q(0)-Q(1))^k}$ . This means that, on average, the inequality  $f(Q(\epsilon)) \geq f(Q(1-\epsilon))$  holds. For the odd ordinal distribution, the discussion is more challenging since it is generally symmetric. Just consider  $\Xi_3$ , let  $x_1 = Q(p_i)$ and  $x_3 = Q(p_j)$ , changing the value of  $x_2$  from  $Q(p_i)$  to  $Q(p_j)$  will monotonically change the value of  $\psi_3(x_1, x_2, x_3)$ , since  $\frac{\partial \psi_3(x_1, x_2, x_3)}{\partial x_2} = -\frac{x_1^2}{2} - x_1 x_2 + 2x_1 x_3 + x_2^2 - x_2 x_3 - \frac{x_3^2}{2},$   $-\frac{3}{4} (x_1 - x_3)^2 \leq \frac{\partial \psi_3(x_1, x_2, x_3)}{\partial x_2} \leq -\frac{1}{2} (x_1 - x_3)^2 \leq 0. \text{ If the original distribution is right-skewed,}$   $\xi_{\Delta}$  will be left-skewed, so, for  $\Xi_3$ , the average probability density of the right side of zero will be greater than that of the left side, which means, on average, the inequality  $f(Q(\epsilon)) \leq f(Q(1-\epsilon))$  holds (the same result can be inferred from the definition of central moments, where the positivity of the odd order central moment is directly related to the left-skewness of the corresponding kernel distribution). In all, the monotonicity of the pairwise difference distribution guides the general shape of the kth central moment kernel distribution, k > 2, forcing it to be unimodal-like with the mode and median close to zero, then, the inequality  $f(Q(\epsilon)) \le f(Q(1-\epsilon))$  or  $f(Q(\epsilon)) \ge f(Q(1-\epsilon))$ holds in general. If a distribution is ordered and all of its central moment kernel distributions are also ordered, it is called completely ordered. Although strict complete orderliness is difficult to prove, even if the inequality may be violated in a small range, as discussed in Subsection A, the mean-SWAmedian inequality remains valid, in most cases, for the central moment kernel distribution.

423

424

425

426

427

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

447

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

Another crucial property of the central moment kernel distribution, location invariant, is introduced in the next theorem. The proof is provided in the SI Text.

**Theorem B.3.** 
$$\psi_k (x_1 = \lambda x_1 + \mu, \dots, x_k = \lambda x_k + \mu) = \lambda^k \psi_k (x_1, \dots, x_k).$$

Consider two continuous distributions belonging to the same location-scale family, their corresponding kth central moment kernel distributions only differ in scaling. So d is invariant, as shown in Subsection A. The recombined kth central moment, based on rm, is defined by,

$$rkm_{d,\epsilon=1-(1-\epsilon_0)^{\frac{1}{k}},n} := (d+1)\operatorname{SW}km_{\epsilon,n} - dmkm_n,$$

where  $SWkm_{\epsilon,n}$  is using the binomial kth central moment  $(Bkm_{\epsilon_0,n})$  here,  $mkm_n$  is the median kth central moment. Since  $SWkm_{\epsilon,n}$  is an L-statistic, the resulting  $rkm_{d,\epsilon,n}$  is an arithmetic I-statistic. Similarly, the quantile will not change after scaling. The quantile kth central moment is thus defined

$$qkm_{d,\epsilon,n} := \hat{Q}_n\left(\left(pSWkm_{\epsilon,n} - \frac{1}{2}\right)d + pSWkm_{\epsilon,n}\right),$$

where  $pSWkm_{\epsilon,n} = \hat{F}_{\psi,n} (SWkm_{\epsilon,n}), \hat{F}_{\psi,n}$  is the empirical cumulative distribution function of the corresponding central moment kernel distribution.  $qkm_{d,\epsilon,n}$  is a quantile *I*-statistic.

Finally, for standardized moments, quantile skewness and

Finally, for standardized moments, quantile skewness are defined to be  $qskew_{d,\epsilon,n} := \frac{qtm_{d,\epsilon,n}}{qsd_{d,\epsilon,n}^3}$ and  $qkurt_{d,\epsilon,n} := \frac{qfm_{d,\epsilon,n}}{qsd_{d,\epsilon,n}^4}$ . Quantile standard deviation  $(qsd_{d,\epsilon,n})$ , recombined standard deviation  $(rsd_{d,\epsilon,n})$ , quantile third central moment  $(qtm_{d,\epsilon,n})$ , quantile fourth central moment  $(qfm_{d,\epsilon,n})$ , recombined third central moment  $(rtm_{d,\epsilon,n})$ , recombined fourth central moment  $(rfm_{d,\epsilon,n})$ , recombined skewness  $(rskew_{d,\epsilon,n})$ , and recombined kurtosis  $(rkurt_{d,\epsilon,n})$ are all defined similarly as above and not repeated here. The transformation to a location problem can also empower related statistical tests. From the better performance of the quantile mean in heavy-tailed distributions, quantile central moments are generally better than recombined central moments regarding asymptotic bias.

To avoid confusion, it should be noted that the robust location estimations of the kernel distributions discussed in this paper differ from the approach taken by Joly and Lugosi

393

394

396

397

402

403

404

405

406

407

409

410

411

412

413

417

418

419

420

(2016) (40), which is computing the median of all U-statistics from different disjoint blocks. Compared to bootstrap median U-statistics, this approach can produce two additional kinds of finite sample bias, one arises from the limited numbers of blocks, another is due to the size of the U-statistics (consider the mean of all U-statistics from different disjoint blocks, it is definitely not identical to the original U-statistic, except when the kernel is the Hodges-Lehmann kernel). Laforgue, Clemencon, and Bertail (2019)'s median of randomized U-statistics (41) is more sophisticated and can overcome the limitation of the number of blocks, but the second kind of bias remains unsolved.

466

467

468

469

472

473

474

475

476

478

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509 510

511

512

513

514

515

516

517

518

519

520

521

**C. Congruent distribution.** In the realm of nonparametric statistics, the precise values of robust estimators are of secondary importance. What is of primary importance is their relative differences or orders. Based on this principle, in the absence of contamination, as the parameters of the distribution vary, all reasonable nonparametric location estimates should asymptotically change in the same direction. Otherwise if the results obtained based on the trimmed mean are completely different from those based on the median, a contradiction arises. However, such contradictions are possible, as in the case of the Weibull distribution,  $E[m] = \lambda \sqrt[\alpha]{\ln(2)}$ ,  $E[\mu] = \lambda \Gamma \left(1 + \frac{1}{\alpha}\right)$ , then, when  $\alpha = 1$ ,  $E[m] = \lambda \ln(2) \approx 0.693\lambda$ ,  $E[\mu] = \lambda$ , but when  $\alpha = \frac{1}{2}$ ,  $E[m] = \lambda \ln^2(2) \approx 0.480\lambda$ ,  $E[\mu] = 2\lambda$ , the mean increases, but the median decreases. To study the conditions that avoid such scenarios by classifying distributions through the signs of derivatives, let the quantile average function of a parametric distribution be denoted as QA  $(\epsilon, \gamma, \alpha_1, \dots, \alpha_i, \dots, \alpha_k)$ , where  $\alpha_i$  represent the parameters of the distribution, then, a distribution is  $\gamma$ -congruent if and only if the sign of  $\frac{\partial QA}{\partial \alpha_i}$  remains the same for all  $0 \le \epsilon \le \frac{1}{1+\gamma}$ If this partial derivative is equal to zero or undefined, it can be considered both positive and negative, and thus does not impact the analysis. Asymptotically, any weighted average can be expressed as an integral of the quantile average function. Since the sign does not change after integration, the sign of  $\frac{\partial QA}{\partial \alpha_i}$  remains the same for all  $0 \le \epsilon \le \frac{1}{1+\gamma}$  implies that all  $\gamma$ -weighted averages change in the same direction as the parameters change, as long as they are not undefined. A distribution is completely  $\gamma$ -congruent if and only if it is  $\gamma$ -congruent and all its central moment kernel distributions are also  $\gamma$ -congruent. Setting  $\gamma = 1$  constitutes the definitions of congruence and complete congruence. Chebyshev's inequality implies that, for any probability distribution with finite moments, even if some weighted averages change in a direction different from that of the sample mean, the deviations are bounded. Furthermore, distributions with infinite moments can be  $\gamma$ -congruent, since the definition is based on the quantile average, not the sample

The following theorems show the conditions that a distribution is congruent or  $\gamma$ -congruent.

**Theorem C.1.** A symmetric distribution with a finite second moment is always congruent.

*Proof.* For any symmetric distribution with a finite second moment, all symmetric quantile averages coincide. The conclusion follows immediately.

**Theorem C.2.** A positive define location-scale distribution with a finite second moment is always  $\gamma$ -congruent.

*Proof.* As shown in discussions in Subsection A, for a location-scale distribution, any weighted average can be expressed as  $\lambda WA_0(\epsilon) + \mu$ , where  $WA_0(\epsilon)$  is an integral of  $Q_0(p)$  according to the definition of the weighted average. Therefore, the derivatives with respect to the parameters  $\lambda$  or  $\mu$  are always positive. By application of the definition, the desired outcome is obtained.

526

527

528

531

532

533

534

535

536

537

538

541

545

546

547

548

549

550

551

552

553

554

555

556

557

559

560

561

562

565

566

567

568

569

570

571

573

574

575

576

577

578

**Theorem C.3.** The second central moment kernal distribution derived from a continuous location-scale unimodal distribution with a finite second moment is always  $\gamma$ -congruent.

*Proof.* Theorem B.3 shows that the corresponding central moment kernel distribution is also a location-scale family distribution. Theorem B.1 shows that it is positively defined. Implementing Theorem C.2 yields the desired result.  $\Box$ 

For the Pareto distribution,  $\frac{\partial Q(p,\alpha)}{\partial \alpha} = \frac{x_m(1-p)^{-1/\alpha} \ln(1-p)}{\alpha^2}$ . Since  $\ln(1-p) < 0$  for all  $0 , <math>(1-p)^{-1/\alpha} > 0$  for all  $0 and <math>\alpha > 0$ , so  $\frac{\partial Q(p,\alpha)}{\partial \alpha} < 0$ , and therefore  $\frac{\partial QA(\epsilon,\gamma,\alpha)}{\partial \alpha} < 0$ , the Pareto distribution is  $\gamma$ -congruent. The derivative for the lognormal distribution is  $\frac{\partial SQA(\epsilon,\sigma)}{\partial \sigma} = \frac{-\text{erfc}^{-1}(2\epsilon)e^{\mu-\sqrt{2}\sigma\text{erfc}^{-1}(2\epsilon)}-\text{erfc}^{-1}(2-2\epsilon)e^{\mu-\sqrt{2}\sigma\text{erfc}^{-1}(2-2\epsilon)}}{\sqrt{2}}$ . Since the inverse complementary error function is positive when the input is smaller than 1, and negative when the input is larger than 1,  $\operatorname{erfc}^{-1}(2\epsilon) = -\operatorname{erfc}^{-1}(2-2\epsilon)$ ,  $e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2-2\epsilon)} > e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2\epsilon)}$ ,  $\frac{\partial\operatorname{SQA}(\epsilon,\sigma)}{\partial\sigma} > 0$ , the lognormal distribution is congruent. Theorem C.1 implies that the generalized Gaussian distribution is congruent. For the Weibull distribution, when  $\alpha$  changes from 1 to  $\frac{1}{2}$ , the average probability density on the left side of the median increases, since  $\frac{\frac{1}{2}}{\lambda \ln(2)} < \frac{\frac{1}{2}}{\lambda \ln^2(2)}$ , but the mean increases, indicating that the distribution is more heavy-tailed, the probability density of large values will also increase. The main reason for non-congruence of a right-skewed smooth partial bounded probability distribution lies in the simultaneous increase of probability densities on two opposite sides: one approaching the bound and the other approaching infinity. Note that the gamma distribution does not have this issue, it looks to be congruent.

Although some common parametric distributions are not congruent, Theorem C.2 establishes that  $\gamma$ -congruence always holds for a positive define location-scale family distribution and thus for the second central moment kernel distribution generated from a continuous location-scale unimodal distribution as shown in Theorem C.3. Theorem B.2 demonstrates that all their central moment kernel distributions are unimodal-like with mode and median close to zero, as long as they are unimodal distributions. Assuming finite moments and constant Q(0) - Q(1), increasing the mean of the kernel distribution will result in a more heavy-tailed distribution, i.e., the probability density of the values close to  $\frac{1}{k}(-\Delta)^k$  increases. While the total probability density on either side of zero remains unchanged as the median is generally close to zero and much less impacted by increasing the mean, the probability density of the values close to zero decreases. This transformation will increase nearly all symmetric weighted averages, in the general sense. Therefore, except for the median, which is assumed to be zero, nearly all symmetric weighted averages for all central moment kernel distributions derived from unimodal

6 | Lee

581

582

583

584

585

586

587

589

590

591

592

593

597

598

599

600

601

602

603

604

605

606

607

608

610

611

612

613

614

615

616

617

618

619

620

621

623

624

627

628

## D. A shape-scale distribution as the consistent distribution.

Up to this point, in this article, the consistent robust estimation has been limited to a location-scale distribution, with the location parameter often being omitted for simplicity. To construct probability distributions can be made to fit the observed skewness and kurtosis arbitrarily well, in 1894, Pearson (42) introduced a family of continuous probability distributions that are now often characterized by the square of the skewness and the kurtosis. If the skewness and the kurtosis are interrelated by a shape parameter, a distribution specified by a shape parameter (denoted as  $\alpha$ ) and a scale parameter (denoted as  $\lambda$ ) is often referred to as a shape-scale distribution. Weibull, gamma, Pareto, lognormal, and generalized Gaussian distributions (when  $\mu$  is a constant) are all shape-scale unimodal distributions. Moreover, if  $\alpha$  or skewness or kurtosis is a constant, the shape-scale distribution is reduced to a locationscale distribution. The above discussion shows that, due to the invariant property, if a location-scale distribution is chosen as the consistent distribution, the type of invarient moments and their related weighted moments are given, there should exist a unique k-tuple  $(d_{im}, \ldots, d_{ikm})$  calibrated by the distribution and the corresponding kernel distributions generated from this distribution. For a right skewed shape-scale distribution, let  $D(|skewness|, kurtosis, k, etype, dtype, n) = d_{ikm}$  denote these relations, where the first input is the absolute value of the skewness, the second input is the kurtosis, the third is the order of the central moment (if k = 1, the mean), the fourth is the type of estimator, the fifth is the type of consistent distribution, and the sixth input is the sample size. For simplicity, the last three inputs will be omitted in the following discussion. Hold in awareness that due to the invariant property of scale, specifying d values for a shape-scale distribution only requires either skewness or kurtosis, while the other may be also omitted. Since many common shape-scale distributions are always right skewed (if not, only the right skewed or left skewed part is used for calibration, while the other part is omitted), the absolute value of the skewness should be identical to the skewness for them and it can also handle the left skew scenario well.

For recombined moments, the object of using a shape-scale distribution as the consistent distribution is to find solutions for the system of equations

$$\begin{cases} rm\left(\text{SWA}, median, D(|rskew|, rkurt, 1)\right) = \mu \\ rvar\left(\text{SW}var, mvar, D(|rskew|, rkurt, 2)\right) = \mu_2 \\ rtm\left(\text{SW}tm, mtm, D(|rskew|, rkurt, 3)\right) = \mu_3 \\ rfm\left(\text{SW}fm, mfm, D(|rskew|, rkurt, 4) = \mu_4 \\ rskew = \frac{\mu_3}{\frac{3}{2}} \\ \frac{\mu_2^2}{rkurt} = \frac{\mu_4}{\mu_2^2} \end{cases}$$

where  $\mu_2$ ,  $\mu_3$  and  $\mu_4$  are the population second, third and fourth central moments. |rskew| and rkurt should be the invariant points of the functions  $\varsigma(|rskew|) = \left| \frac{rtm(SWtm,mtm,D(|rskew|,3))}{r} \right|$  and

tions  $\zeta(|rskew|) = \left| \frac{rtm(SWtm,mtm,D(|rskew|,3))}{rvar(SWvar,mvar,D(|rskew|,2))^{\frac{3}{2}}} \right|$  and  $\varkappa(rkurt) = \frac{rfm(SWfm,mfm,D(rkurt,4))}{rvar(SWvar,mvar,D(rkurt,2))^2}$ . Clearly, this is an overdetermined nonlinear system of equations, given that the skewness and kurtosis are interrelated for a shape-scale

distribution. Since an overdetermined system constructed with random coefficients is almost always inconsistent, it is natural to optimize them separately using the fixed-point iteration (see Algorithm 1, only rkurt is provided, others are the same).

633

634

635

636

637

638

639

640

641

644

645

646

647

648

649

650

651

653

654

655

656

657

658

659

660

661

662

664

665

666

667

668

670

Algorithm 1 rkurt for a shape-scale distribution

Input: D; SWvar; SWfm; mvar; mfm; maxit;  $\delta$  Output:  $rkurt_{i-1}$ 

i = 0

2:  $rkurt_i \leftarrow \varkappa(kurtosis_{max}) \triangleright \text{Using the maximum kurtosis}$  available in D as an initial guess.

repeat

number.

4: i = i + 1  $rkurt_{i-1} \leftarrow rkurt_i$ 6:  $rkurt_i \leftarrow \varkappa(rkurt_{i-1})$ 

until i > maxit or  $|rkurt_i - rkurt_{i-1}| < \delta \implies maxit$  is the maximum number of iterations,  $\delta$  is a small positive

The following theorem shows the validaty of Algorithm 1.

**Theorem D.1.** Assuming the mkms are all equal to zero, |rskew| and rkurt, defined as the largest attracting fix points of the functions  $\varsigma(|rskew|)$  and  $\varkappa(rkurt)$ , are consistent estimators of  $\tilde{\mu}_3$  and  $\tilde{\mu}_4$  for a shape-scale distribution whose central moment kernel distributions are all congruent, as long as they are within the domain of D, where  $\tilde{\mu}_3$  and  $\tilde{\mu}_4$  are the population skewness and kurtosis.

*Proof.* Without loss of generality, only rkurt is considered here, while the logic for |rskew| is the same. Also, according to the property of invariance, the second central moments of the sample and consistent distribution are all assumed to be 1. From the definition of D,  $\frac{\varkappa(rkurt_D)}{rkurt_D} =$ 

$$\frac{\frac{\mu_{4_D} - \text{SW} f m_D}{\text{SW} f m_D - m f m_D} (\text{SW} f m - m f m) + \text{SW} f m}{\text{SW} f m_D - m f m_D} (\text{SW} f m - m f m) + \text{SW} f m} \frac{\mu_{2_D} - \text{SW} var_D}{\text{SW} var_D - m var_D} (\text{SW} var - m var) + \text{SW} var} \right)^2, \text{ where the sub-}$$

script  $\hat{D}$  indicates that the estimates are from the central moment kernel distributions generated from the consistent distribution used to calibrate the d values, while other estimates are from the sample.

Then, assuming the mkms are all equal to zero,  $\frac{\varkappa(rkurt_D)}{rkurt_D} =$ 

$$\frac{\frac{\mu_{4_D} - \text{SW}fm_D}{\text{SW}fm_D}(\text{SW}fm) + \text{SW}fm}{\text{SW}trt_D} \left(\frac{\text{SW}var}{\text{SW}var_D}\right)^2} = \frac{\left(\frac{\mu_{4_D} - \text{SW}fm_D}{\text{SW}fm_D} + 1\right)(\text{SW}fm)}{\text{SW}fm_D}}{\frac{\text{SW}fm}{\text{SW}var_D}} = \frac{\frac{\text{SW}fm}{\text{SW}var_D}}{\text{SW}fm_D}}{\frac{\text{SW}fm}{\text{SW}var_D}} = \frac{\text{SW}kurt}{\text{SW}kurt_D}.$$

From the definitions, SWkurt is also a measure of kurtosis, so an increase in  $rkurt_D$  will also result in an increase in SWkurt<sub>D</sub>,  $\lim_{rkurt_D\to\infty}\frac{\varkappa(rkurt_D)}{rkurt_D}<1$ . As a result, if there is at least one fix point, let the largest one be  $fix_{max}$ , then it is attracting since  $|\frac{\partial(\varkappa(rkurt_D))}{\partial(rkurt_D)}|<1$  for all  $rkurt_D \in [fix_{max}, kurtosis_{max}]$ .

Asymptotically, consider any  $SWkurt_D > SWkurt$ ,  $\frac{\varkappa(rkurt_D)}{rkurt_D} < 1$ , the same logic applies, a consistent estimator must be the last attracting fix point,  $fix_{max}$  is the consistent estimator.

**Data Availability.** Data for Table ?? are given in SI Dataset S1. All codes have been deposited in GitHub.

- **ACKNOWLEDGMENTS.** I gratefully acknowledge the construc-671 672 tive comments made by the editor which substantially improved the clarity and quality of this paper. 673
  - CF Gauss, Theoria combinationis observationum erroribus minimis obnoxiae. (Henricus Dieterich), (1823),

674 675

676

677

678

679

680

681

682

683

684

687 688

689 690

695 696

697

698

706 707

708

709

710

719

721 722

727

728

731

732

738 739

- S Newcomb, A generalized theory of the combination of observations so as to obtain the best result. Am. journal Math. 8, 343-366 (1886).
- S Newcomb, Researches on the motion of the moon. part ii, the mean motion of the moon and other astronomical elements derived from observations of eclipses and occultations extending from the period of the babylonians until ad 1908. United States. Naut. Alm. Off. Astron. paper; v. 99. 1 (1912).
- PJ Huber, Robust estimation of a location parameter. Ann. Math. Stat. 35, 73-101 (1964).
- X He, WK Fung, Method of medians for lifetime data with weibull models. Stat. medicine 18, 1993-2009 (1999).
- M Menon, Estimation of the shape and scale parameters of the weibull distribution. Techno-685 metrics 5, 175-182 (1963). 686
  - SD Dubey, Some percentile estimators for weibull parameters. *Technometrics* **9**, 119–129 (1967).
  - KM Hassanein. Percentile estimators for the parameters of the weibull distribution. Biometrika 58, 673-676 (1971).
- NB Marks, Estimation of weibull parameters from common percentiles. J. applied Stat. 32. 691 692 17-24 (2005).
- K Boudt, D Caliskan, C Croux, Robust explicit estimators of weibull parameters. Metrika 73. 693 694 187-209 (2011).
  - SD Dubey, Contributions to statistical theory of life testing and reliability. (Michigan State University of Agriculture and Applied Science. Department of statistics), (1960)
  - 12. LJ Bain, CE Antle, Estimation of parameters in the weibdl distribution. Technometrics 9, 621-627 (1967).
- 699 RV Hogg, Adaptive robust procedures: A partial review and some suggestions for future 700 applications and theory. J. Am. Stat. Assoc. 69, 909-923 (1974).
- 701 RJ Hyndman, Y Fan, Sample quantiles in statistical packages. The Am. Stat. 50, 361-365 702 (1996)
- C Bernard, R Kazzi, S Vanduffel, Range value-at-risk bounds for unimodal distributions under 703 15. partial information. Insur. Math. Econ. 94, 9-24 (2020). 704
- 705 WR van Zwet, Convex Transformations of Random Variables: Nebst Stellingen. (1964).
  - AL Bowley, Elements of statistics. (King) No. 8, (1926).
  - 18. RA Groeneveld, G Meeden, Measuring skewness and kurtosis. J. Royal Stat. Soc. Ser. D (The Stat. 33, 391-399 (1984).
  - PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models iv. spread in Selected Works of EL Lehmann. (Springer), pp. 519-526 (2012).
- 711 20. PJ Rousseeuw, C Croux, Alternatives to the median absolute deviation. J. Am. Stat. association 88, 1273-1283 (1993). 712
- 713 PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models. iii. dispersion in Selected works of EL Lehmann. (Springer), pp. 499-518 (2012). 714
- 715 P Janssen, R Serfling, N Veraverbeke, Asymptotic normality of u-statistics based on trimmed 716 samples. J. statistical planning inference 16, 63-74 (1987).
- 717 J Hodges, E Lehmann, Matching in paired comparisons. The Annals Math. Stat. 25, 787-791 718
- S Dharmadhikari, K Jogdeo, Unimodal laws and related in A Festschrift For Erich L. Lehmann. (CRC Press), p. 131 (1982). 720
  - AY Khintchine, On unimodal distributions. Izv. Nauchno-Isled. Inst. Mat. Mech. 2, 1-7 (1938). S Purkayastha, Simple proofs of two results on convolutions of unimodal distributions. Stat. & probability letters 39, 97-100 (1998).
- 723 724 RA Fisher, Moments and product moments of sampling distributions, Proc. Lond. Math. Soc. 2, 199-238 (1930). 725
- 726 28. PR Halmos, The theory of unbiased estimation. The Annals Math. Stat. 17, 34-43 (1946).
  - W Hoeffding, A class of statistics with asymptotically normal distribution. The Annals Math. Stat. 19, 293-325 (1948).
- PM Heffernan, Unbiased estimation of central moments by using u-statistics. J. Royal Stat. 729 Soc. Ser. B (Statistical Methodol. 59, 861-863 (1997). 730
  - D Fraser, Completeness of order statistics. Can. J. Math. 6, 42-45 (1954).
    - 32. AJ Lee, U-statistics: Theory and Practice. (Routledge), (2019).
- J Hodges Jr, E Lehmann, Estimates of location based on rank tests. The Annals Math. Stat. 733 **34**, 598-611 (1963). 734
- 34. A Ehsanes Saleh, Hodges-lehmann estimate of the location parameter in censored samples. 735 Annals Inst. Stat. Math. 28, 235-247 (1976). 736
- RJ Serfling, Generalized I-, m-, and r-statistics. The Annals Stat. 12, 76-86 (1984). 737
  - P Janssen, R Serfling, N Veraverbeke, Asymptotic normality for a general class of statistical functions and applications to measures of spread. The Annals Stat. 12, 1369-1379 (1984).
- MG Akritas. Empirical processes associated with v-statistics and a class of estimators under 740 random censoring. The Annals Stat. 14, 619-637 (1986). 741
- I Gijbels, P Janssen, N Veraverbeke, Weak and strong representations for trimmed u-statistics. 742 Probab, theory related fields 77, 179-194 (1988). 743
- J Choudhury, R Serfling, Generalized order statistics, bahadur representations, and sequential 744 nonparametric fixed-width confidence intervals. J. Stat. Plan. Inference 19, 269-282 (1988). 745
- 40. E Joly, G Lugosi, Robust estimation of u-statistics. Stoch. Process. their Appl. 126, 3760-3773 746 747 (2016).
- 748 P Laforque, S Clémencon, P Bertail, On medians of (randomized) pairwise means in Interna-749 tional Conference on Machine Learning. (PMLR), pp. 1272-1281 (2019).
- 750 K Pearson, Contributions to the mathematical theory of evolution. Philos. Transactions Royal Soc. London, A 185, 71-110 (1894). 751

8 | Lee