

# Near-consistent robust estimations of moments for unimodal distributions

Tuban Lee<sup>a,1</sup>

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

This manuscript was compiled on March 2, 2023

**Descriptive statistics for parametric models currently rely heavily 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 continuous 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

The 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) but unsolved. Strictly speaking, it is unsolvable because 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 this problem 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 contaminated normal distribution, which has played a pre-eminent role in the later development of robust statistics. As previously demonstrated, under growing asymmetric departures from normality, the bias of the Huber M-estimator increases rapidly. This is a common issue in parametric estimations. For example, He and Fung (1999) constructed (5) a robust M-estimator for the two-parameter Weibull distribution. All moments can be calculated from its estimated parameters. As expected, it is inadequate for the gamma, Perato, lognormal, and especially the generalized Gaussian distributions, because the logarithmic function does not produce a result for negative inputs (SI Dataset S1). Instead of minimizing the residuals, another old and interesting approach is arithmetically computing the parameters using one or more  $L$ -statistics as input values, e.g., the percentile estimators. Examples 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 clear that this class of estimators arithmetically utilizes the invariant structures of probability distributions (6, 11, 12). Maybe it can be named as  $I$ -statistics. Formally, an estimator is classified as an  $I$ -statistic if asymptotically it satisfies  $I(WA_1, \dots, WA_l) = (\theta_1, \dots, \theta_q)$  for the distribution it is consistent with, where WAs are weighted averages,  $\theta$ s are the population parameters it estimates. If the function  $I$  is solely defined through addition and/or subtraction, it is also an  $L$ -statistic. In the previous article, it is shown that quantile average is fundamental for all weighted

averages. Based on the quantile function,  $I$ -statistic is naturally robust. For many parametric distributions, the quantile functions are much more elegant than the pdfs and cdfs. So  $I$ -statistics are often analytically obtainable. However, the performance of the above 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.

In previous work on semiparametric robust mean estimation, although greatly shrinking the asymptotic biases, binomial mean ( $BM_\epsilon$ ) is still inconsistent for any skewed distribution if  $\epsilon > 0$  (if  $\epsilon \rightarrow 0$ , since the alternating sum of binomial coefficients is zero,  $BM \rightarrow \mu$ ). All robust location estimators commonly used are symmetric due to the universality of the symmetric distributions. One can construct an asymmetric trimmed mean that is consistent for a semiparametric class of skewed distributions. This approach was investigated previously, but it is not symmetric and therefore only suitable for some special applications (13). From semiparametric to parametric, an ideal robust location estimator would have a non-sample-dependent breakdown point (defined in Subsection F) and be consistent with any symmetric distribution and a skewed distribution with finite second moments. This is called an invariant mean. Based on the mean-symmetric weighted average-median inequality, the recombined mean is defined as

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

where  $d$  is for bias correction,  $SWA_{\epsilon,n}$  is  $BM_{\epsilon,n}$  in the first three Subsections, while other symmetric weighted averages can also be used in practice as long as the inequalities hold. The next theorem shows the significance of this composite estimator.

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

**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  $\alpha \rightarrow \infty$ .

*Proof.* Finding  $d, \epsilon$  values that make  $rm_{d, \epsilon}$  a consistent mean estimator is equivalent to finding the solution of  $E[rm_{d, \epsilon}] = E[X]$ . Rearranging the definition,  $rm_{d, \epsilon} = \lim_{c \rightarrow \infty} \left( \frac{(BM_{\epsilon} + c)^{d+1}}{(median + c)^d} - c \right) = (d+1)BM_{\epsilon} - dmedian = \mu$ . So,  $d = \frac{\mu - BM_{\epsilon}}{BM_{\epsilon} - median}$ . The pdf of the exponential distribution is  $f(x) = \lambda^{-1}e^{-\lambda^{-1}x}$ ,  $\lambda \geq 0$ ,  $x \geq 0$ , the cdf is  $F(x) = 1 - e^{-\lambda^{-1}x}$ ,  $x \geq 0$ . The quantile function is  $Q(p) = \ln\left(\frac{1}{1-p}\right)\lambda$ .  $E[x] = \lambda$ .  $E[median] = Q\left(\frac{1}{2}\right) = \ln 2\lambda$ . For the exponential distribution, the expectation of  $BM_{\frac{1}{8}}$  is  $E\left[BM_{\frac{1}{8}}\right] = \lambda\left(1 + \ln\left(\frac{46656}{8575\sqrt{35}}\right)\right)$ . Obviously, the scale parameter  $\lambda$  can be canceled out,  $d \approx 0.375$ . 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  $E[rm_{d, \epsilon}] = \lim_{c \rightarrow \infty} \left( \frac{(E[X] + c)^{d+1}}{(E[X] + c)^d} - c \right) = E[X]$ . 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  $d_{Pareto} = \frac{\frac{\alpha x_m}{\alpha - 1} - x_m(1-p_1)^{-\frac{1}{\alpha}}}{x_m(1-p_1)^{-\frac{1}{\alpha}} - x_m(1-p_2)^{-\frac{1}{\alpha}}}$ . Since any weighted average can be expressed as an integral of the quantile function,  $\lim_{\alpha \rightarrow \infty} \frac{\frac{\alpha}{\alpha-1} - (1-p_1)^{-1/\alpha}}{(1-p_1)^{-1/\alpha} - (1-p_2)^{-1/\alpha}} = -\frac{\ln(1-p_1)+1}{\ln(1-p_1)-\ln(1-p_2)}$ , the  $d$  value for the Pareto distribution approaches that of the exponential distribution as  $\alpha \rightarrow \infty$ , regardless of the type of weighted average used. This completes the demonstration.  $\square$

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 of these two-parameter distributions. 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).  $rm_{d \approx 0.375, \epsilon = \frac{1}{8}}$  has 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 the estimation of central moments can be transformed into a location estimation problem by using  $U$ -statistics, the central moment kernel distributions have nice properties, and, in light of previous works, a series of sophisticated yet efficient robust estimators can be constructed whose biases are typically smaller than the variances ( $n = 5400$ , Table ??) 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 has the form  $F(x) = F_0\left(\frac{x-\mu}{\lambda}\right)$ , where  $F_0$  is a "standard" distribution. Then,  $F(x) = Q^{-1}(x) \rightarrow x = Q(p) = \lambda Q_0(p) + \mu$ . So, 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. The simultaneous

cancellation of  $\mu$  and  $\lambda$  in  $\frac{(\lambda\mu_0 + \mu) - (\lambda BM_0(\epsilon) + \mu)}{(\lambda BM_0(\epsilon) + \mu) - (\lambda median_0 + \mu)}$  ensures that  $d$  is a constant. Consequently, the roles of  $BM_{\epsilon}$  and median in  $rm_{d, \epsilon}$  can be replaced by any weighted averages, although for the definition of invariant mean, only symmetric weighted averages are considered here.

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

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

provided that  $\hat{F}_n(SWA_{\epsilon, n}) \geq \frac{1}{2}$ , where  $\hat{F}_n(x)$  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,  $\hat{F}_n(x) := \frac{1}{n} \left( \frac{x - Q_n\left(\frac{sp}{n}\right)}{Q_n\left(\frac{1}{n}(sp+1)\right) - Q_n\left(\frac{sp}{n}\right)} + sp \right)$ , where  $sp = \sum_{i=1}^n \mathbf{1}_{X_i \leq x}$ ,  $\mathbf{1}_A$  is the indicator of event  $A$ . The solution of  $\hat{F}_n(SWA_{\epsilon, n}) < \frac{1}{2}$  is reversing the percentile by  $1 - \hat{F}_n(SWA_{\epsilon, n})$ , the obtained percentile is also reversed. Without loss of generality, in the following discussion, only the  $\hat{F}_n(SWA_{\epsilon, n}) \geq \frac{1}{2}$  case will be considered. Moreover, in extreme heavy-tailed distributions, the calculated percentile can exceed the breakdown point of  $SWA_{\epsilon}$ , so the percentile will be modified to  $1 - \epsilon$  if this happens. 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 ( $\alpha \rightarrow \infty$ ) and any symmetric distributions provided that the second moments are finite.

*Proof.* Similarly, rearranging the definition,  $d = \frac{F(\mu) - F(BM_{\epsilon})}{F(BM_{\epsilon}) - \frac{1}{2}}$ .

Recall the cdf is  $F(x) = 1 - e^{-\lambda^{-1}x}$ ,  $x \geq 0$ , the expectation of  $BM_{\epsilon}$  can be expressed as  $\lambda BM_0(\epsilon)$ , so  $F(BM_{\epsilon})$  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,  $F(E[BM_{\epsilon}]) = F(\mu) = \frac{1}{2}$ . 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

$$d_{Pareto} = \frac{1 - \left(\frac{x_m}{\frac{\alpha x_m}{\alpha - 1}}\right)^{\alpha} - \left(1 - \left(\frac{x_m}{x_m(1-p_1)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right)}{\left(1 - \left(\frac{x_m}{x_m(1-p_1)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right) - \left(1 - \left(\frac{x_m}{x_m(1-p_2)^{-\frac{1}{\alpha}}}\right)^{\alpha}\right)} = \frac{1 - \left(\frac{\alpha - 1}{\alpha}\right)^{\alpha} - p_1}{p_1 - p_2}. \text{ When } \alpha \rightarrow \infty, \left(\frac{\alpha - 1}{\alpha}\right)^{\alpha} = \frac{1}{e}. \text{ The } d \text{ value for the exponential distribution is identical, since } d_{exp} = \frac{(1 - e^{-1}) - \left(1 - e^{-\ln\left(\frac{1}{1-p_1}\right)}\right)}{\left(1 - e^{-\ln\left(\frac{1}{1-p_1}\right)}\right) - \left(1 - e^{-\ln\left(\frac{1}{1-p_2}\right)}\right)} = \frac{1 - \frac{1}{e} - p_1}{p_1 - p_2}. \text{ All results are now proven. } \square$$

The definitions of location and scale parameters are such that they must satisfy  $F(x; \lambda, \mu) = F\left(\frac{x-\mu}{\lambda}; 1, 0\right)$ . Recall that

152  $x = \lambda Q_0(p) + \mu$ , so the percentile of any weighted average  
 153 is free of  $\lambda$  and  $\mu$ , guaranteeing the validity of the quantile  
 154 mean.  $qm_{d \approx 0.321, \epsilon = \frac{1}{8}}$  works better in the fat-tail scenarios (SI  
 155 Dataset S1). Theorem .1 and A.1 show that  $rm_{d \approx 0.375, \epsilon = \frac{1}{8}}$   
 156 and  $qm_{d \approx 0.321, \epsilon = \frac{1}{8}}$  are both consistent mean estimators for  
 157 any symmetric distribution and a skewed distribution with  
 158 finite second moments. It's obvious that the breakdown points  
 159 of  $rm_{d \approx 0.375, \epsilon = \frac{1}{8}}$  and  $qm_{d \approx 0.321, \epsilon = \frac{1}{8}}$  are both  $\frac{1}{8}$ . Therefore  
 160 they are all invariant means.

161 To study the impact of the choice of SWAs in  $rm$  and  $qm$ ,  
 162 it is constructive to consider a symmetric weighted average  
 163 as a mixture of symmetric quantile averages. Although using  
 164 a less-biased symmetric weighted average can generally  
 165 improve performance (SI Dataset S1), there is a higher risk  
 166 of violation in the semiparametric framework. However, even  
 167 the SQA function is not strictly monotonic, suppose it is  
 168 generally decreasing in  $[0, u]$ , but increasing in  $[u, \frac{1}{2}]$ , since  
 169  $1 - 2\epsilon$  of the symmetric quantile averages will be included  
 170 in the computation of  $SWA_\epsilon$ , as long as  $|u - \frac{1}{2}| \ll 1 - 2\epsilon$ ,  
 171 and other parts of the SQA function satisfy the inequality  
 172 constraints which define the  $\nu$ th orderliness, the mean- $SWA_\epsilon$ -  
 173 median inequality will still be valid (as an example, the SQA  
 174 function is non-monotonic when the shape parameter of the  
 175 Weibull distribution  $\alpha > \frac{1}{1 - \ln(2)} \approx 3.259$  as shown in the  
 176 previous article, yet the mean- $BM_{\frac{1}{8}}$ -median inequality is still  
 177 valid when  $\alpha \leq 3.322$ ). Another key factor determining the  
 178 risk of violation is the skewness of the distribution. In the  
 179 previous article, it is shown that in a family of distributions  
 180 that differ by a skewness-increasing transformation in van  
 181 Zwet's sense, the violation of orderliness, if it happens, often  
 182 only occurs when the distribution is near-symmetric (15).  
 183 The over-corrections in  $rm$  and  $qm$  are dependent on the  
 184  $SWA_\epsilon$ -median difference, which is correlated to the skewness  
 185 (16, 17), so the over-correction is often tiny with a moderate  
 186  $d$ . This qualitative analysis provides another perspective, in  
 187 addition to the bias bounds (18), that  $rm$  and  $qm$  based on  
 188 the mean- $SWA_\epsilon$ -median inequality are generally safe.

189 **B. Robust estimations of the central moments.** In 1976, Bickel  
 190 and Lehmann, in their third paper of the landmark series *De-*  
 191 *scriptive Statistics for Nonparametric Models* (19), generalized  
 192 a class of estimators called "measures of disperse," which is now  
 193 often named as Bickel-Lehmann dispersion. As an example,  
 194 they proposed a first version of the trimmed standard deviation,  
 195  $\hat{\tau}^2(F; \epsilon) \equiv \tau^2(F; \epsilon)$ , for independent and identically  
 196 distributed random variables  $X_i$  with a distribution  $F$ , where  
 197  $\tau^2(F; \epsilon) = \frac{1}{1-2\epsilon} \int_{Q(\epsilon)}^{Q(1-\epsilon)} y dG(y)$ ,  $Q$  is the quantile function  
 198 of  $G$ ,  $G$  is the distribution of  $Y = X^2$ . Obviously, when  
 199  $\epsilon = 0$ , the result is equivalent to the second raw moment.  
 200 In 1979, in the same series (20), they explored another class  
 201 of estimators called "measures of spread," which "does not  
 202 require the assumption of symmetry." From that, a popular  
 203 efficient scale estimator, the Rousseeuw-Croux scale estimator  
 204 (21), was derived in 1993, but the importance of tackling the  
 205 symmetry assumption has been greatly underestimated. In  
 206 the final section of the paper, they considered another two  
 207 possible versions of the trimmed standard deviations, which  
 208 were modified here for comparison,

$$\left[ n \left( \frac{1}{2} - \epsilon \right) \right]^{-\frac{1}{2}} \left[ \sum_{k=\frac{n}{2}}^{n(1-\epsilon)} [X_k - X_{n-k+1}]^2 \right]^{\frac{1}{2}}, \quad [1] \quad 209$$

and 210

$$\left[ \binom{n}{2} (1 - \epsilon - \gamma\epsilon) \right]^{-\frac{1}{2}} \left[ \sum_{k=\binom{n}{2}\epsilon}^{\binom{n}{2}(1-\gamma\epsilon)} (X - X')_k^2 \right]^{\frac{1}{2}}, \quad [2] \quad 211$$

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

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

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

225 *Nomenclature.* Given a robust estimator  $\hat{\theta}$ . The first part of  
 226 the name of the robust statistic defined in this paper is a prefix  
 227 that indicates the type of estimator, and the second part is  
 228 the name of the population parameter  $\theta$  that the estimator is  
 229 consistent with as  $\epsilon \rightarrow 0$ . The abbreviation of the estimator is  
 230 the initial letter(s) of the first part plus the common abbrevia-  
 231 tion of the consistent estimator that measures the population  
 232 parameter. If the estimator is symmetric and not a  $U$ -statistic,  
 233 the breakdown point,  $\epsilon$ , is indicated in the subscript of the  
 234 abbreviation of the estimator. If the estimator is asymmetric,  
 235 the corresponding  $\gamma$  is also indicated after  $\epsilon$ . If the estimator  
 236 is a weighted  $U$ -statistic, the breakdown point of the location  
 237 estimator is indicated (except the median).

238 In the previous semiparametric robust mean article, it is  
 239 shown that the bias of a reasonable robust estimator should  
 240 be monotonic with respect to the breakdown point in a semi-  
 241 parametric distribution and naturally, its name should align  
 242 with the consistent estimator. Naturally, the trimmed stan-  
 243 dard deviation following this nomenclature is  $Tsd_{\epsilon, \gamma, n} :=$   
 244  $\left[ TM_{\epsilon, \gamma} \left( (\psi_2(X_{N_1}, X_{N_2}))_{N=1}^{\binom{n}{2}} \right) \right]^{-\frac{1}{2}}$ , where  $TM_{\epsilon, \gamma}(Y)$  denotes  
 245 the  $\epsilon, \gamma$ -trimmed mean with the sequence  $(\psi_2(X_{N_1}, X_{N_2}))_{N=1}^{\binom{n}{2}}$   
 246 as an input. If the square root is removed, it is named as the  
 247 trimmed variance ( $Tvar_{\epsilon, \gamma, n}$ ). It is now very clear that this  
 248 definition, essentially the same as [2], should be preferable.  
 249 Not only because it is essentially a trimmed  $U$ -statistic for  
 250 the standard deviation but also because the  $\gamma$ -orderliness of  
 251 the pseudo-sample distribution is ensured by the next exciting  
 252 theorem.

253 **Theorem B.1.** *The second central moment kernel distribution*  
 254 *generated from any continuous unimodal distribution is  $\gamma$ -*  
 255 *ordered, if  $\gamma \geq 1$ .*

256 *Proof.* Let  $Q(p)$ ,  $0 \leq p \leq 1$ , denote the quantile of the contin-  
 257 uous unimodal distribution  $f_X(x)$ . The corresponding proba-  
 258 bility density is  $f(Q(p))$ . Generating the distribution of the



pair  $(Q(p_i), Q(p_j))$ ,  $i < j$ ,  $p_i < p_j$ , the corresponding probability density is  $f_{X,X}(Q(p_i), Q(p_j)) = 2f(Q(p_i))f(Q(p_j))$ . Transforming the pair  $(Q(p_i), Q(p_j))$ ,  $i < j$ , by the function  $\Phi(x_1, x_2) = x_1 - x_2$ , the pairwise difference distribution has a mode that is arbitrary close to  $M - M = 0$ . 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, Dharmadhikari and Jogdeo in 1982 (24) gave 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  $\frac{1}{2}$  does not change the monotonicity, making the pdf become monotonically decreasing with mode at zero. In the previous semiparametric robust mean estimation article, it is proven that a right skewed distribution with a monotonic decreasing pdf is always  $\gamma$ -ordered, which gives the desired result.  $\square$

**Remark.** The assumption of continuity of distributions is important for monotonicity because, unlike in the continuous case, it is possible to get pairs with the same value for a discrete distribution. For example, let the probabilities of the singletons  $\{1\}$ ,  $\{2\}$ ,  $\{3\}$ ,  $\{4\}$  and  $\{5\}$  of a probability mass function of a discrete probability distribution be  $\frac{1}{11}$ ,  $\frac{4}{11}$ ,  $\frac{3}{11}$ ,  $\frac{2}{11}$ , and  $\frac{1}{11}$ , respectively. This is a unimodal distribution, but the corresponding  $\psi_2$  distribution is non-monotonic, whose singletons  $\{0\}$ ,  $\{0.5\}$ ,  $\{2\}$ ,  $\{4.5\}$  and  $\{8\}$  have probabilities  $\frac{21}{66}$ ,  $\frac{24}{66}$ ,  $\frac{2}{14}$ ,  $\frac{6}{66}$ , and  $\frac{1}{66}$ , respectively.

Previously, it was shown that any symmetric distribution with a finite second moment is  $\nu$ th ordered. That means the orderliness does not require unimodality, e.g., for a symmetric bimodal distribution, it is also ordered. Examples from the Weibull distribution show 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, unbiasedness 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  $k$ th 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-Lehmann estimator (33) to the trimmed H-L mean (he named "Wilcoxon one-sample statistic") (34). In 1984, Serfling pointed out the speciality of Hodges-Lehmann estimator, which is neither a simple  $L$ -statistics nor  $U$ -statistic, and considered the generalized  $L$ -statistics and  $U$ -statistic structure (35). Also in 1984, Janssen and Serfling and Veraverbeke (36) showed that the Bickel-Lehmann spread also belongs to the same class. It was gradually clear that the Hodges-Lehmann estimator, trimmed H-L mean and trimmed standard deviation are all trimmed  $U$ -statistics (37–39). Due to the combinatorial explosion, the bootstrap (40), introduced by Efron in 1979, is indispensable in large sample studies. In 1981, Bickel and Freedman (41) showed that the bootstrap

is asymptotically valid to approximate the original distribution in a wide range of situations, including  $U$ -statistics. The limit laws of bootstrapped trimmed  $U$ -statistics was proven by Helmers, Janssen, and Veraverbeke (1990) (42).

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

$$Wkm_{\epsilon, \gamma, n} := WA_{\epsilon, \gamma, n} \left( (\psi_k(X_{N_1}, \dots, X_{N_k}))_{N=1}^n \right),$$

where  $X_{N_1}, \dots, X_{N_k}$  are the  $n$  choose  $k$  elements from  $X$ ,  $\psi_k(x_1, \dots, x_k) = \sum_{j=0}^{k-2} (-1)^j \binom{1}{k-j} \sum (x_{i_1}^{k-j} \dots x_{i_{j+1}}) + (-1)^{k-1} (k-1) x_1 \dots x_k$ , the second summation is over  $i_1, \dots, i_{j+1} = 1$  to  $k$  with  $i_1 < \dots < i_{j+1}$  (30). Despite the complexity, the structure of the  $k$ th 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, let  $x_1 = Q(p_i)$  and  $x_k = Q(p_j)$ ,  $\Delta = Q(p_i) - Q(p_j)$ . The  $k$ th central moment kernel distribution,  $k > 2$ , can be seen as a mixture distribution and each of the components has the support  $(-\frac{k}{3+(-1)^k})^{-1}(-\Delta)^k, \frac{1}{k}(-\Delta)^k$ .

*Proof.* Generating the distribution of the  $k$ -tuple  $(Q(p_{i_1}), \dots, Q(p_{i_k}))$ ,  $k > 2$ ,  $i_1 < \dots < i_k$ ,  $p_{i_1} < \dots < p_{i_k}$ , the corresponding probability density is  $f_{X, \dots, X}(Q(p_{i_1}), \dots, Q(p_{i_k})) = k!f(Q(p_{i_1})) \dots f(Q(p_{i_k}))$ . Transforming the distribution of the  $k$ -tuple by the function  $\psi_k(x_1, \dots, x_k)$ , denoting  $\bar{\Delta} = \psi_k(Q(p_{i_1}), \dots, Q(p_{i_k}))$ . The probability  $f_{\Xi_k}(\bar{\Delta}) = \sum_{\bar{\Delta}=\psi_k(Q(p_{i_1}), \dots, Q(p_{i_k}))} f_{X, \dots, X}(Q(p_{i_1}), \dots, Q(p_{i_k}))$  is the summation of the probabilities of all  $k$ -tuples such that  $\bar{\Delta}$  is equal to  $\psi_k(Q(p_{i_1}), \dots, Q(p_{i_k}))$ . 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  $x_2, \dots, x_{k-1}$  exhaust all combinations under the inequality constraints, i.e.,  $Q(p_{i_1}) = x_1 < x_2 < \dots < x_{k-1} < x_k = Q(p_{i_k})$ . 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  $x_1 = \dots = x_k = 0$ , where  $\psi_k = 0$ . Other candidates are within the boundaries, i.e.,  $\psi_k(x_1 = x_1, x_2 = x_k, \dots, x_k = x_k)$ ,  $\dots$ ,  $\psi_k(x_1 = x_1, \dots, x_i = x_1, x_{i+1} = x_k, \dots, x_k = x_k)$ ,  $\dots$ ,  $\psi_k(x_1 = x_1, \dots, x_{k-1} = x_1, x_k = x_k)$ .  $\psi_k(x_1 = x_1, \dots, x_i = x_1, x_{i+1} = x_k, \dots, x_k = x_k)$  can be divided into  $k$  groups. If  $\frac{k+1-i}{2} \leq j \leq \frac{k-1}{2}$ , from  $j+1$ st to  $k-j$ th group, the  $g$ th group has  $i \binom{i-1}{g-j-1} \binom{k-i}{j}$  terms having the form  $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$ , from  $k-j+1$ th to  $i+j$ th group, the  $g$ th group has  $i \binom{i-1}{g-j-1} \binom{k-i-1}{j} + (k-i) \binom{k-i-1}{j-k+g-1} \binom{i}{k-j}$  terms having the form  $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$ . If  $j < \frac{k+1-i}{2}$ , from  $j+1$ st to  $i+j$ th group, the  $g$ th group has  $i \binom{i-1}{g-j-1} \binom{k-i}{j}$  terms having the form  $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$ . If  $j \geq \frac{k}{2}$ , from  $k-j+1$ st to  $j$ th group, the  $g$ th group has  $(k-i) \binom{k-i-1}{j-k+g-1} \binom{i}{k-j}$  terms having the form  $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$ , from  $j+1$ th to  $j+i$ th group,  $i+j < k$ , the  $g$ th group

has  $i \binom{i-1}{g-j-1} \binom{k-i}{k-g-1} \binom{i}{k-j}$  terms having the form  $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$ . The final  $k$ th group is the term  $(-1)^{k-1} (k-1) x_1^i x_k^{k-i}$ . So, if  $i+j = k$ ,  $j \geq \frac{k}{2}$ ,  $i \leq \frac{k}{2}$ , the summed coefficient of  $x_1^i x_k^{k-i}$  is  $(-1)^{k-1} (k-1) + \sum_{g=i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} (k-i) \binom{k-i-1}{g-i-1} + \sum_{g=k-i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-k+i-1} = (-1)^{k-1} (k-1) + (-1)^{k+1} + (k-i) (-1)^k + (-1)^k (i-1) = (-1)^{k+1}$ . The summation identities are  $\sum_{g=i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} (k-i) \binom{k-i-1}{g-i-1} = (k-i) \int_0^1 \sum_{g=i+1}^{k-1} (-1)^{g+1} \binom{k-i-1}{g-i-1} t^{k-g} dt = (k-i) \int_0^1 ((-1)^i (t-1)^{k-i-1} - (-1)^{k+1}) dt = (k-i) \left( \frac{(-1)^k}{i-k} + (-1)^k \right) = (-1)^{k+1} + (k-i) (-1)^k$ .  $\sum_{g=k-i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-k+i-1} = \int_0^1 \sum_{g=k-i+1}^{k-1} (-1)^{g+1} i \binom{i-1}{g-k+i-1} t^{k-g} dt = \int_0^1 (i (-1)^{k-i} (t-1)^{i-1} - i (-1)^{k+1}) dt = (-1)^k (i-1)$ . If  $j < \frac{k+1-i}{2}$ ,  $i > k-1$ , if  $i = k$ ,  $\psi_k = 0$ , if  $\frac{k+1-i}{2} \leq j \leq \frac{k-1}{2}$ ,  $\frac{k+1}{2} \leq i \leq k-1$ , 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  $i+j < k$ , since  $\binom{i}{k-j} = 0$ , the related terms can be ignored, so, using the binomial theorem and beta function, the summed coefficient of  $x_1^{k-j} x_k^j$  is  $\sum_{g=j+1}^{i+j} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-j-1} \binom{k-i}{j} = i \binom{k-i}{j} \int_0^1 \sum_{g=j+1}^{i+j} (-1)^{g+1} \binom{i-1}{g-j-1} t^{k-g} dt = \binom{k-i}{j} i \int_0^1 ((-1)^j t^{k-j-1} \left( \frac{t-1}{t} \right)^{i-1}) dt = \binom{k-i}{j} i \frac{(-1)^{j+i+1} \Gamma(i) \Gamma(k-j-i+1)}{\Gamma(k-j+1)} = \frac{(-1)^{j+i+1} i! (k-j-i)! (k-i)!}{(k-j)! j! (k-j-i)!} = (-1)^{j+i+1} \frac{i! (k-i)!}{k!} \frac{k!}{(k-j)! j!} = \binom{k-i}{i}^{-1} (-1)^{1+i} \binom{k}{j} (-1)^j$ . The coefficient of  $x_1^i x_k^{k-i}$  in  $\binom{k}{i}^{-1} (-1)^{1+i} (x_1 - x_k)^k$  is  $\binom{k-i}{i}^{-1} (-1)^{1+i} \binom{k}{i} (-1)^{k-i} = (-1)^{k+1}$ , same as the 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. Therefore,  $\psi_k(x_1 = x_1, \dots, x_i = x_1, x_{i+1} = x_k, \dots, x_k = x_k) = \binom{k}{i}^{-1} (-1)^{1+i} (x_1 - x_k)^k$ , the maximum and minimum of  $\psi_k$  follow directly from the properties of the binomial coefficient.  $\square$

$\xi_\Delta$  is closely related to  $f_\Xi(\Delta)$ , which is the pairwise difference distribution, since the probability density of  $\xi_\Delta$  is  $f_{\Xi_k}(\bar{\Delta}|\Delta) = \sum_{\bar{\Delta} = -(\frac{k}{2} - \Delta)^k}^{\frac{1}{k}(-\Delta)^k} f_{\Xi_k}(\bar{\Delta}|\Delta) = f_\Xi(\Delta)$ . 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  $\Xi_k$  is also 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  $\xi_\Delta$  close to zero. 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 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  $\frac{1}{\binom{k}{2}^{-1} (Q(0) - Q(1))^k} > \frac{1}{\frac{1}{k} (Q(0) - Q(1))^k}$ . This means that, on average, the inequality  $f(Q(\epsilon)) \geq f(Q(1 - \epsilon))$  holds. For the odd ordinal distribution, the discussion is harder 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  $\frac{\partial \psi_3(x_1, x_2, x_3)}{\partial x_2} = -\frac{x_1^2}{2} - x_1 x_2 + 2x_1 x_3 + x_2^2 - x_2 x_3 - \frac{x_3^2}{2}$ ,  $-\frac{3}{4} (x_1 - x_3)^2 \leq \frac{\partial \psi_3(x_1, x_2, x_3)}{\partial x_2} \leq -\frac{1}{2} (x_1 - x_3)^2 \leq 0$ . 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, the positive of 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  $k$ th central moment kernel distribution,  $k > 2$ , forcing it to be unimodal-like with mode and median close to zero, then, the inequality  $f(Q(\epsilon)) \leq f(Q(1 - \epsilon))$  or  $f(Q(\epsilon)) \geq f(Q(1 - \epsilon))$  holds in general. If a distribution is ordered and its all central moment kernel distributions are also ordered, it is called completely ordered. Although strict complete orderliness is hard to prove, the inequality may be violated in a small range, as discussed in Subsection A, the mean-SWA $_{\epsilon}$ -median inequality remains valid, in most cases, for the central moment kernel distribution.

Another key property of the central moment kernel distribution, location invariant, is introduced in the next theorem. The proof is given 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  $k$ th central moment kernel distributions only differ in scaling. So  $d$  is invariant, as shown in Subsection A. The recombined  $k$ th central moment, based on  $rm$ , is defined by,

$$rkm_{d,\epsilon,n} := (d+1) SWkm_{\epsilon,n} - d mkm_{\epsilon,n},$$

where  $SWkm_{\epsilon,n}$  is using the binomial  $k$ th central moment ( $Bkm_{\epsilon,n}$ ) here,  $mkm_{\epsilon,n}$  is the median  $k$ th central moment. Similarly, the quantile will not change after scaling. The quantile  $k$ th central moment is thus defined as

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

where  $pSWkm = \hat{F}_{\psi,n}(SWkm_{\epsilon,n})$ ,  $\hat{F}_{\psi,n}$  is the empirical cumulative distribution function of the corresponding central moment kernel distribution.

Finally, for standardized moments, quantile skewness and quantile kurtosis 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, the robust location estimations of the kernel distributions here are different from Joly and Lugosi (2016) (43)'s approach, which is computing the median of all  $U$ -statistics from different disjoint blocks based on the median of means technique, although asymptotically, as discussed in the previous article, it can be equivalent to the median  $U$ -statistic if the size of each block is equal to the degree of the kernel. Laforgue, Clemencon, and Bertail (2019) proposed the median of randomized  $U$ -statistics (43, 44), which is more sophisticated and closer to the median  $U$ -statistic if setting an additional constraint on the block size.

**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. In the absence of contamination, as the parameters of the distribution vary, a reasonable nonparametric location estimator will asymptotically change in the same direction as the other location estimators. Otherwise if the results based on trimmed mean are completely different from those based on median, a contradiction arises. A distribution satisfying this property for any symmetric weighted average is called a congruent distribution. If extending to any  $\epsilon, \gamma$ -weighted average, it is  $\gamma$ -congruent. A distribution is completely congruent if and only if it is congruent and its all central moment kernel distributions are also congruent. Complete  $\gamma$ -congruence is analogous. 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. Also, distributions with infinite moments can be congruent, since it is defined here that from infinity to infinity, the direction change can be interpreted as both increasing and decreasing. The following theorems show the conditions that a distribution is congruent or  $\gamma$ -congruent.

**Theorem C.1.** *Let the symmetric quantile average function of a parametric distribution be denoted as  $SQA(\epsilon, \alpha_1, \dots, \alpha_i, \dots, \alpha_k)$ , where  $\alpha_i$  represent the parameters of the distribution. This distribution is congruent if and only if the sign of  $\frac{\partial SQA(\epsilon, \alpha_i)}{\partial \alpha_i}$  remains the same (if equal to zero, it can be seen as both positive and negative and thus also impact the analysis) for all  $0 < \epsilon \leq \frac{1}{2}$ . Replacing  $SQA$  with  $QA_{\epsilon, \gamma}$  constitutes a necessary and sufficient condition for the distribution to be considered  $\gamma$ -congruent.*

*Proof.* Asymptotically, any symmetric weighted average can be expressed as an integral of the symmetric quantile average function. Since the sign won't change after integration, from definition, the sign of  $\frac{\partial SQA(\epsilon, \alpha_i)}{\partial \alpha_i}$  remains the same for all  $0 < \epsilon \leq \frac{1}{2}$  is equivalent to all symmetric weighted averages change in the same direction as the parameters change. The same logic applies to the  $\gamma$ -congruence case, as the constancy of the sign of  $\frac{\partial QA_{\epsilon, \gamma}(\alpha_i)}{\partial \alpha_i}$  for all  $0 < \epsilon \leq \frac{1}{2}$  is equivalent to the statement that all  $\gamma$ -weighted averages also change in the same direction. The proof is finished.  $\square$

**Theorem C.2.** *If a distribution is  $\gamma$ -congruent, it is congruent.*

*Proof.* Any symmetric weighted average is also a weighted average. This concludes the proof.  $\square$

**Theorem C.3.** *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.  $\square$

**Theorem C.4.** *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 Theorem C.1, the desired outcome is obtained.  $\square$

**Theorem C.5.** *The second central moment kernel 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.4 yields the desired result.  $\square$

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  $0 < p < 1$ ,  $(1-p)^{-1/\alpha} > 0$  for all  $0 < p < 1$  and  $\alpha > 0$ , 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  $\frac{\partial SQA(\epsilon, \sigma)}{\partial \sigma} = \frac{-\operatorname{erfc}^{-1}(2\epsilon)e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2\epsilon)} - \operatorname{erfc}^{-1}(2-2\epsilon)e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2-2\epsilon)}}{\sqrt{2}}$ . Since the inverse complementary error function is positive when the input is smaller than 1, and negative when the input is larger than 1,  $\operatorname{erfc}^{-1}(2\epsilon) = -\operatorname{erfc}^{-1}(2-2\epsilon)$ ,  $e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2-2\epsilon)} > e^{\mu-\sqrt{2}\sigma\operatorname{erfc}^{-1}(2\epsilon)}$ ,  $\frac{\partial SQA(\epsilon, \sigma)}{\partial \sigma} > 0$ , the lognormal distribution is congruent. Theorem C.3 implies that the generalized Gaussian distribution is congruent. For the Weibull distribution, just consider the median and mean,  $E[m] = \lambda \sqrt[\alpha]{\ln(2)}$ ,  $E[\mu] = \lambda \Gamma(1 + \frac{1}{\alpha})$ , 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. Therefore, it is not congruent. When  $\alpha$  changes from 1 to  $\frac{1}{2}$ , the average probability density on the left side of median increases, since  $\frac{\frac{1}{2}}{\lambda \ln(2)} < \frac{\frac{1}{2}}{\lambda \ln^2(2)}$ , but the mean increases, meaning that the distribution is more heavy-tailed, the probability density of large values will also increase. The reason for non-congruence lies in the simultaneous increase of probability densities on two opposite sides: one approaching zero and the other approaching infinity. Note that the gamma distribution does not have this issue, it looks to be congruent.

Although many common parametric distributions are not congruent, Theorem C.4 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.5. 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. This implies, align with Theorem B.3, that different kernel distributions mainly differ in scale and they are, in some senses, reduced to a location-scale family distribution. Assuming finite moments, if  $Q(0) - Q(1)$  remains constant, increasing the mean of the kernel distribution will result in a more heavy-tailed distribution, i.e., the probability density closer to  $\frac{1}{k}(-\Delta)^k$  will increase. While the total probability density on either side of zero will remain unchanged as the median is generally close to zero and much less impacted during the mean increasing, the probability density close to zero will decrease. This transformation will increase nearly all symmetric weighted averages, in the general sense, due to the heavy tail. As a result, nearly all symmetric weighted averages, except the median since it is assumed to be zero, for all central moment kernel distributions derived from unimodal distributions should change in the same direction when the parameters change.

## D. A two-parameter distribution as the consistent distribution.

Up to this point, the consistent robust estimation has been limited to a parametric location-scale distribution. The location parameter is often omitted for simplicity. A distribution specified by a shape parameter (denoted as  $\alpha$  here) and a scale parameter (denoted as  $\lambda$  here) is often referred to as a two-parameter distribution. Weibull, gamma, Pareto, lognormal, and generalized Gaussian distributions are all two-parameter unimodal distributions.  $\alpha$  can be converted to skewness or kurtosis, e.g., for the gamma distribution, the skewness is  $\frac{2}{\sqrt{\alpha}}$ , the kurtosis is  $\frac{6}{\alpha} + 3$ . If  $\alpha$  is a constant, the two-parameter distribution is reduced to a single-parameter distribution. The above discussion shows that, for a single-parameter distribution as the consistent distribution and a fixed  $\epsilon$ , there should be a  $k$ -tuple  $(d_{im}, \dots, d_{ikm})$  (using a distribution as the consistent distribution means the  $d$  values used are calibrated by the distribution and the corresponding kernel distributions generated from this distribution). For a two-parameter distribution, let  $D(kurtosis, |skewness|, k, etype, dtype, n) = d_{ikm}$  denote these relations, where the first input is the kurtosis, the second input is the absolute value of the skewness, 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, the sixth input is the sample size. For simplicity, the last three inputs will be omitted in the following discussion. Note that specifying  $d$  values for a two-parameter distribution requires only kurtosis or skewness.

Using a two-parameter distribution as the consistent distribution is a problem of robust estimation of parametric models. For recombined moments, the object is to find solutions for the system of equations

$$\begin{cases} rm(SWA, median, D(rkurt, |rskev|, 1)) = \mu \\ rvar(SWvar, mvar, D(rkurt, |rskev|, 2)) = \mu_2 \\ rtm(SWtm, mtm, D(rkurt, |rskev|, 3)) = \mu_3 \\ rfm(SWfm, mfm, D(rkurt, |rskev|, 4)) = \mu_4 \\ rskew = \frac{\mu_3}{\mu_2} \\ rkurt = \frac{\mu_4}{\mu_2} \end{cases}, \quad \text{where}$$

$\mu_2, \mu_3$  and  $\mu_4$  are the population second, third and fourth central moments.  $rkurt$  and  $|rskev|$  should be the invariant

points of the functions  $\kappa(rkurt) = \frac{rfm(SWfm, mfm, D(rkurt, 4))}{rvar(SWvar, mvar, D(rkurt, 2))^2}$

and  $\varsigma(|rskev|) = \frac{rtm(SWtm, mtm, D(|rskev|, 3))}{rvar(SWvar, mvar, D(|rskev|, 2))^{\frac{3}{2}}}$ . Clearly, this is an overdetermined nonlinear system of equations, because the skewness and kurtosis are interrelated for a two-parameter distribution. As 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).

### Algorithm 1 $rkurt$ for a two-parameter distribution

**Input:**  $D$ ;  $SWvar$ ;  $SWfm$ ;  $mvar$ ;  $mfm$ ;  $maxit$ ;  $\delta$

**Output:**  $rkurt_{i-1}$

```

1:  $i = 0$ 
2:  $rkurt_i \leftarrow \kappa(kurtosis_{max})$   $\triangleright$  Using the maximum kurtosis
   available in  $D$  as an initial guess.
3: repeat
4:    $i = i + 1$ 
    $rkurt_{i-1} \leftarrow rkurt_i$ 
5:    $rkurt_i \leftarrow \kappa(rkurt_{i-1})$ 
6:   until  $i > maxit$  or  $|rkurt_i - rkurt_{i-1}| < \delta$   $\triangleright maxit$  is
   the maximum number of iterations,  $\delta$  is a small positive
   number.
```

The following theorem shows the validity of Algorithm 1.

**Theorem D.1.**  $rkurt$  and  $|rskev|$ , defined as the largest attracting fix points of the functions  $\kappa(rkurt)$  and  $\varsigma(|rskev|)$ , are consistent estimators of  $\tilde{\mu}_4$  and  $\tilde{\mu}_3$  for a completely congruent two-parameter distribution with finite moments, as long as they are within the domain of  $D$ , where  $\tilde{\mu}_4$  and  $\tilde{\mu}_3$  are the population kurtosis and skewness.

*Proof.* Without loss of generality, only  $rkurt$  is considered here, while the logic for  $|rskev|$  is the same. From the definition of  $D$ ,  $\lim_{rkurt \rightarrow \infty} \frac{\kappa(rkurt)}{rkurt} =$

$$\lim_{rkurt \rightarrow \infty} \frac{\frac{rkurt \mu_{2cali}^2 - SWfm_{cali}}{SWfm_{cali} - mfm_{cali}} (SWfm - mfm) + SWfm}{rkurt \left( \frac{\mu_{2cali}^2 - SWvar_{cali}}{SWvar_{cali} - mvar_{cali}} (SWvar - mvar) + SWvar \right)^2}.$$

Since  $SWfm_{cali}$  and  $mfm_{cali}$  are from the same kernel distribution as  $\mu_{4cali} = rkurt \mu_{2cali}^2$ , so an increase in  $\mu_{4cali}$  will also result in an increase in  $SWfm_{cali}$  and hence  $SWfm_{cali} \gg SWfm$ . Furthermore, Theorem B.2 and qualitative discussion in Subsection B shows that  $mfm_{cali}$  is close to zero, the increases in  $SWfm_{cali}$  leads to an increase in  $(SWfm_{cali} - mfm_{cali})$ . According to the property of invariance, assuming  $rkurt = \mu_{4cali}$ ,  $\left( \frac{\mu_{2cali}^2 - SWvar_{cali}}{SWvar_{cali} - mvar_{cali}} (SWvar - mvar) + SWvar \right)^2 > 1$ , then

$\lim_{rkurt \rightarrow \infty} \frac{rkurt - SWfm_{cali}}{SWfm_{cali} - mfm_{cali}} (SWfm - mfm) + SWfm < 1$ . As a result, if there is at least one fix point, let the largest one be  $fix_{max}$ , then it is attracting since  $|\frac{\partial(\kappa(rkurt))}{\partial(rkurt)}| < 1$  for all  $rkurt \in [fix_{max}, kurtosis_{max}]$ . Asymptotically, consider any  $rkurt > \mu_4$ ,  $SWfm_{cali} > SWfm$ ,  $mfm_{cali} > mfm$ ,  $\frac{rkurt - SWfm_{cali}}{SWfm_{cali} - mfm_{cali}} (SWfm - mfm) + SWfm < rkurt$ , the same logic applies, a consistent estimator must be the last attracting fix point,  $fix_{max}$  is the consistent estimator.  $\square$

As a result of Theorem D.1, assuming continuity and congruence of the central moment kernel distributions, Algorithm

1 converges surely if a fix point exists within the domain of  $D$ . At this stage,  $D$  can only be approximated through a Monte Carlo study. Using linear interpolation can ensure continuity. A common encountered problem is that  $D$  has a domain depending on both the consistent distribution and the Monte Carlo study, so the iteration may halt at the boundary if the fix point is not within the domain. However, by setting a proper maximum number of iterations, the algorithm will return the boundary value which is optimal within  $D$ . For quantile moments, the logic is similar, if the percentiles do not exceed the breakdown point. If so, consistent estimation is impossible and the algorithm will stop due to the maximum number of iterations. The fix point iteration is, in principle, similar to the iterative reweighing in M-estimator, but an advantage of this algorithm is that the optimization is only related to the function of  $d$  value and is independent of the sample size (except for the quantile moments, which require re-computation of the quantile function, but this operation has a time complexity of  $O(1)$  for a sorted sample). Since  $|rskew|$  can specify  $d_{rm}$  after optimization, this enables the robust estimations of all four moments to reach a near-consistent level for all five unimodal distributions (Table ??, SI Dataset S1), just using the Weibull distribution as the consistent distribution.

**E. Variance.** As the fundamental theorem in statistics, the central limit theorem states that the standard deviation of the limiting form of the sampling distribution of the sample mean is  $\frac{\sigma}{\sqrt{n}}$ . The principle was later applied to the sampling distributions of the robust location estimators (2, 36, 45–52) and it was found that the efficiencies of the robust location estimators are sometimes very different from the arithmetic mean. Daniell (1920) stated (45) that the comparison of efficiencies of the various kinds of estimators is useless unless they all tend to coincide asymptotically. Bickel and Lehmann argued, also in the landmark series (51, 52), that meaningful comparisons can be made by studying the standardized variance, asymptotic variances, and sharp lower bounds of these estimators. Here, the scaled standard error (SSE) is proposed to estimate the variances of all estimators, including recombined/quantile moments, on a scale similar to that of the sample mean.

*Definition E.1 (Scaled standard error).* Let  $\mathcal{M}_{s_i s_j} \in \mathbb{R}^{i \times j}$  denote the sample-by-statistics matrix, i.e., the first column is the main statistics of interest,  $\widehat{\theta}_m$ , the second to the  $j$ th column are  $j - 1$  statistics required to scale,  $\widehat{\theta}_{r_1}$ ,  $\widehat{\theta}_{r_2}$ , ...,  $\widehat{\theta}_{r_{j-1}}$ . Then, the scaling factor  $\mathcal{S} = \left[ 1, \frac{\widehat{\theta}_{r_1}}{\widehat{\theta}_m}, \frac{\widehat{\theta}_{r_2}}{\widehat{\theta}_m}, \dots, \frac{\widehat{\theta}_{r_{j-1}}}{\widehat{\theta}_m} \right]^T$  is a  $j \times 1$  matrix, which  $\widehat{\theta}$  is the mean of the column. The normalized matrix is  $\mathcal{M}_{s_i s_j}^N = \mathcal{M}_{s_i s_j} \mathcal{S}$ . The SSEs are the unbiased standard deviations of the corresponding columns.

Setting the bootstrap moments as the main statistics of interest, the SSEs of all robust estimators proposed here are often between those of the median moments and those of the unbiased sample moments (SI Dataset S1). This is because similar monotonic relations between robustness and variance are also very common, e.g., Bickel and Lehmann (52) proved that a bound for the efficiency of  $\text{TM}_\epsilon$  to sample mean is  $(1 - 2\epsilon)^2$  and this monotonic bound is valid for any distribution. Lai, Robbins, and Yu (1983) proposed an estimator that adaptively chooses the mean or median in a symmetric distribution and showed that the result is typically as good

as the better of sample mean and median regarding variance (53). It may be interpreted as an attempt to use the variance version of mean-SWA-median inequality. While they used bootstrap standard error as the criterion, another approach can be dated back to Laplace (1812) (54) is using a linear combination of median and mean and the weight is deduced to achieve minimum variance; examples for symmetric distributions see Samuel-Cahn, Chan and He, and Damilano and Puig (55–57).

Scaled standard error enables the direct comparison of variances of different location estimators for asymmetric distributions. Here, two invariant means and related symmetric weighted averages ( $\text{BM}_{\frac{1}{8}}$ ,  $\text{SQM}_{\frac{1}{8}}$ ,  $\text{BM}_{\nu=2, \epsilon=\frac{1}{8}}$ ,  $\text{WM}_{\frac{1}{8}}$ ,  $\text{BWM}_{\frac{1}{8}}$ , and  $\text{TM}_{\frac{1}{8}}$  used here) can create twelve possible combinations. Each combination has an SSE for a single-parameter distribution, which can be inferred by a Monte Carlo study. Then, among twelve possible combinations, there is one that has the smallest SSE (if the percentiles of quantile moments exceed the breakdown point, this combination will be excluded). Theoretically, bootstrap is the optimal way to infer the variance-optimal choice without distributional assumptions, however, the computational cost is very high. Similar to Subsection D, let  $I(\text{kurtosis}, |\text{skewness}|, k, \text{dtype}, n) = \text{ikm}_{\text{WA}}$  denote these relations. Then since  $\lim_{rkurt \rightarrow \infty} \frac{I(rkurt, 4)}{I(rkurt, 2)^2 rkurt} < 1$ , the same fix point iteration algorithm can be used to choose the variance-optimum combination. The only problem is that unlike  $D$ ,  $I$  is defined to be discontinuous and also simulated by a Monte Carlo study, but linear interpolation can also be used to ensure the continuity. Using this approach, the result is often very close to the optimum choice (SI Datasets S1).

In 1958, Richtmyer proposed quasi-Monte Carlo simulation based on low-discrepancy sequences, which dramatically reduces the computational cost of large sample simulation (58). Quasi-Monte Carlo methods frequently employ Sobol sequences as the favored numerical sets (59). Do and Hall extended the principle to bootstrap in 1991 (60) and found that the quasi-random approach is competitive in terms of variance when compared with other bootstrap Monte Carlo procedures. By using quasi-sampling, the impact of the number of repetitions of the bootstrap, or bootstrap size, on variance is negligible (SI Dataset S1). An estimator based on the quasi-bootstrap approach can be seen as a very complex deterministic estimator which is not only computationally efficient, but also statistical efficient. The only drawback is that small bootstrap size can produce additional finite sample bias but this can be corrected by re-calibrating the  $d$  values. The default bootstrap size is setting as 18 thousand here as it balances computational cost and finite sample bias, except the asymptotic value calculation. In general, compared to the unbiased sample central moments (deduced by Cramér (61)), the variances of invariant central moments are much smaller (except the second central moment, Table ??).

**F. Robustness.** The measure of robustness to gross errors used here is the breakdown point proposed by Hampel (62) in 1968. However, the sample-dependent breakdown point has apparently not been defined previously.

*Definition F.1 (Sample-dependent breakdown point).* An estimator  $\hat{\theta}$  has a sample-dependent breakdown point if and only if its asymptotic breakdown point  $\epsilon(\hat{\theta}, R, \zeta, v)$  is zero and the empirical influence function of  $\hat{\theta}$  is bounded, where  $R$  is



the measure of badness,  $\zeta$  is the contaminating processes,  $v$  is the uncontaminated process. For a full formal definition of the asymptotic breakdown point, which is the breakdown point when  $n \rightarrow \infty$ , and the empirical influence function, the reader is referred to Genton and Lucas (2003) and Devlin, Gnanadesikan and Kettenring (1975)'s papers (63, 64).

Bear in mind that it differs from the "infinitesimal robustness" defined by Hampel, which is related to whether the asymptotic influence function is bounded (65–67). The proof of the consistency of MoM assumes that it is an estimator with a sample-dependent breakdown point since its breakdown point is  $\frac{b}{2n}$ , where  $b$  is the number of blocks, then  $\lim_{n \rightarrow \infty} \left(\frac{b}{2n}\right) = 0$ , if  $b$  is a constant and any changes in any one of the points of the sample cannot breakdown this estimator (68–70).

For the robust estimations of central moments or other weighted  $U$ -statistics based on a robust location estimator, the asymptotic breakdown points are suggested by the following theorem by extending the method in Donoho and Huber (1983)'s proof of the breakdown point of the Hodges-Lehmann estimator (71).

**Theorem F.1.** *Given  $n$  independent random variables  $(X_1, \dots, X_n)$  with the same distribution  $F$  and a  $U$ -statistic associated with a symmetric kernel of degree  $k$ . Then, assuming as  $n \rightarrow \infty$ ,  $k \ll n$ , the asymptotic breakdown point of the weighted  $U$ -statistic is  $1 - (1 - \epsilon)^{\frac{1}{k}}$ , where  $\epsilon$  is the breakdown point of the weighted average.*

*Proof.* According to the definition of  $\epsilon$ -contamination (71), suppose  $m$  contaminants are added to the sample. The fraction of bad values in the sample is  $\epsilon_U = \frac{m}{n+m}$ , while the original  $n$  data points are not impacted. In the distribution of the kernel,  $\binom{n}{k}$  of total  $\binom{n+m}{k}$  points are not corrupted. Then, the breakdown will not occur if the following inequality holds

$$\binom{n}{k} > \left(\frac{1}{\epsilon} - 1\right) \times \left(\binom{n+m}{k} - \binom{n}{k}\right),$$

Since  $\epsilon$  is the breakdown point of the weighted average,  $\frac{1}{2} \geq \epsilon \geq 0$ ,

$$\frac{1}{1 - \epsilon} > \frac{\binom{n+m}{k}}{\binom{n}{k}} = \frac{(n+m)(n+m-1)\dots(n+m-k+1)}{n(n-1)\dots(n-k+1)}.$$

For asymptotic breakdown point, assuming  $n \rightarrow \infty$ ,  $k \ll n$ ,  $\lim_{n \rightarrow \infty} \left(\frac{n+m-k+1}{n-k+1}\right) = \frac{n+m}{n} = x$ , then the above inequality does not hold when  $x \geq \left(\frac{1}{1-\epsilon}\right)^{\frac{1}{k}}$ . So, the breakdown point of the weighted  $U$ -statistic is  $\epsilon_U = \frac{m}{n+m} = 1 - \frac{n}{n+m} = 1 - \frac{1}{x} = 1 - (1 - \epsilon)^{\frac{1}{k}}$ .  $\square$

*Remark.* If  $k = 1$ ,  $1 - (1 - \epsilon)^{\frac{1}{k}} = \epsilon$ , so this formula also holds for the weighted average. If the weighted average is asymmetric, the breakdown point of it is the minimum of  $\epsilon$  and  $\gamma\epsilon$ . In addition, the numerical solutions for  $k = 2, 3, 4$ ,  $\epsilon = \frac{1}{8}$  are  $\approx 0.065, 0.044$ , and  $0.033$ , respectively. When  $\epsilon = \frac{1}{2}$ , the weighted  $U$ -statistic becomes  $U$ -quantile, which converges almost surely as proven by Choudhury and Serfling (39) in 1988.

Every statistic is based on certain assumptions. For instance, the sample mean assumes that the second moment of the underlying distribution is finite. If this assumption is violated, the variance of the sample mean becomes infinitely

large even the population mean is finite. Therefore, the sample mean not only has zero robustness to gross errors, but also has zero robustness to departures. If departures are unlimited, nearly any estimators can be broken, so posing a constraint on departures is necessary for comparison.

Bias bound (1) is the first approach to study the robustness to departures under regularity conditions, i.e., although all estimators can be biased under departures from the assumptions, but their standardized maximum biases can differ substantially (72, 73). In the previous semiparametric robust mean estimation article, it is shown that another way to qualitatively compare the estimators' robustness to departures from the symmetry assumption is constructing and comparing the corresponding semiparametric models. An estimator based on a smaller model is naturally more robust to asymmetric departures within that model. Although the comparison is limited to the smaller semiparametric model and is not universal, quite surprisingly, the results coincide with those obtained from the bias bound analysis. Bias bound is more universal since it is possible to deduce the bounds for distributions with finite moments without assuming unimodality (72, 73). However, the bias bounds are often hard to deduce for complex estimators. Also, sometimes there are discrepancies between maximum bias and average bias. For example, the maximum bias of  $rm_{d \approx 0.375, \epsilon = \frac{1}{8}}$  is higher than  $SQM_{\frac{1}{8}}$ , but it has much better average bias (SI Dataset S1). Since the estimators proposed here are all consistent under certain assumptions, measuring their biases is a convenient way of measuring the robustness to departures. Average asymptotic bias is thus defined as follows.

**Definition F.2** (Average asymptotic bias). For a single-parameter distribution, the average asymptotic bias (AAB) is just the standardized asymptotic bias  $\frac{|\hat{\theta} - \theta|}{\sigma}$ , where  $\hat{\theta}$  is the estimation of  $\theta$  and  $\sigma$  is the population standard deviation of the distribution. For a two-parameter distribution, the first step is setting the lower bound of the kurtosis range of interest  $\tilde{\mu}_4$ . Then, the average asymptotic bias is defined as

$$AAB_{\hat{\theta}} := \frac{1}{C} \sum_{\substack{\delta + \tilde{\mu}_4 \leq \tilde{\mu}_4 \leq C\delta + \tilde{\mu}_4 \\ \tilde{\mu}_4 \text{ is a multiple of } \delta}} E_{\hat{\theta}|\tilde{\mu}_4} \left[ \frac{|\hat{\theta} - \theta|}{\sigma} \right]$$

where  $\tilde{\mu}_4$  is the kurtosis specifying the two-parameter distribution,  $E_{\hat{\theta}|\tilde{\mu}_4}$  denotes the expected value given that the  $\tilde{\mu}_4$  is fixed.

Standardization is crucial for comparing the performances of estimators under different distributions. Currently, there are several options available, such as using the root mean square deviation from the mode (as in Gauss (1)), the mean absolute deviation, or standard deviation. The standard deviation is used here because of its central role in standard error estimation. The estimation of central moments based on the location estimations of the kernel distributions also enables the standardization of average biases (ABs, for finite sample scenarios) and average asymptotic biases (AABs) of robustified central moments. The only difference is that the population standard deviation is replaced by the asymptotic standard deviation of the kernel distribution ( $\sigma_{km}$ ).

In Table ??,  $\delta = 0.1$ ,  $C = 120$ . For the Weibull, gamma, lognormal and generalized Gaussian distributions, the kurtosis range is from 3 to 15 (there are two shape parameter solutions for the Weibull distribution, the lower one is used here). For

the Pareto distribution, the range is from 9 to 21. To provide a more practical and straightforward illustration, all results from five distributions are further weighted by the number of Google Scholar search results. Interestingly, the asymptotic biases of  $TM_{\frac{1}{8}}$  and  $WM_{\frac{1}{8}}$ , after averaging and weighting, are  $0.128\sigma$  and  $0.078\sigma$ , respectively, in line with the sharp bias bounds of  $TM_{2,14:15}$  and  $WM_{2,14:15}$  (a different subscript is used to indicate a sample size of 15, with the removal of the first and last order statistics.),  $0.173\sigma$  and  $0.126\sigma$ , for distributions with finite moments, without assuming unimodality (72, 73). This setting seems arbitrary, however, the orderliness of symmetric quantile averages ensures that the order among different SWAs remains generally the same as the parameters change, the range of kurtosis is in fact not very important for AAB. It is important if using the maximum biases within the range of kurtosis among all five unimodal distributions as a measure of robustness to departures, because different estimators reach their maximum biases at different parameters. Within the range of kurtosis setting, nearly all SWAs and  $SWkms$  proposed here reach or at least close to their maximum biases (SI Dataset S1). The pseudo-maximum bias is thus defined as the maximum value of the biases in the AAB computations for all five unimodal distributions. In most cases, the pseudo-maximum biases of invariant moments occur in lognormal or generalized Gaussian distributions (SI Dataset S1), since besides unimodality, the Weibull distribution differs entirely from them.

**Data Availability.** Data for Table ?? are given in SI Dataset S1. All codes have been deposited in [GitHub](#).

**ACKNOWLEDGMENTS.** I gratefully acknowledge the constructive comments made by the editor which substantially improved the clarity and quality of this paper.

1. CF Gauss, *Theoria combinationis observationum erroribus minimis obnoxiae*. (Henricus Dieterich), (1823).
2. S Newcomb, A generalized theory of the combination of observations so as to obtain the best result. *Am. journal Math.* **8**, 343–366 (1886).
3. 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. **9**, 1 (1912).
4. PJ Huber, Robust estimation of a location parameter. *Ann. Math. Stat.* **35**, 73–101 (1964).
5. X He, WK Fung, Method of medians for lifetime data with weibull models. *Stat. medicine* **18**, 1993–2009 (1999).
6. M Menon, Estimation of the shape and scale parameters of the weibull distribution. *Technometrics* **5**, 175–182 (1963).
7. SD Dubey, Some percentile estimators for weibull parameters. *Technometrics* **9**, 119–129 (1967).
8. KM Hassanein, Percentile estimators for the parameters of the weibull distribution. *Biometrika* **58**, 673–676 (1971).
9. NB Marks, Estimation of weibull parameters from common percentiles. *J. applied Stat.* **32**, 17–24 (2005).
10. K Boudt, D Caliskan, C Croux, Robust explicit estimators of weibull parameters. *Metrika* **73**, 187–209 (2011).
11. 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 weibull distribution. *Technometrics* **9**, 621–627 (1967).
13. RV Hogg, Adaptive robust procedures: A partial review and some suggestions for future applications and theory. *J. Am. Stat. Assoc.* **69**, 909–923 (1974).
14. RJ Hyndman, Y Fan, Sample quantiles in statistical packages. *The Am. Stat.* **50**, 361–365 (1996).
15. WR van Zwet, *Convex Transformations of Random Variables: Nebst Stellingen*. (1964).
16. AL Bowley, *Elements of statistics*. (King) No. 8, (1926).
17. RA Groeneveld, G Meeden, Measuring skewness and kurtosis. *J. Royal Stat. Soc. Ser. D (The Stat.)* **33**, 391–399 (1984).
18. C Bernard, R Kazzi, S Vanduffel, Range value-at-risk bounds for unimodal distributions under partial information. *Insur. Math. Econ.* **94**, 9–24 (2020).
19. PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models. iii. dispersion in *Selected works of EL Lehmann*. (Springer), pp. 499–518 (2012).

20. PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models iv. spread in *Selected Works of EL Lehmann*. (Springer), pp. 519–526 (2012).
21. PJ Rousseeuw, C Croux, Alternatives to the median absolute deviation. *J. Am. Stat. association* **88**, 1273–1283 (1993).
22. P Janssen, R Serfling, N Veraverbeke, Asymptotic normality of u-statistics based on trimmed samples. *J. statistical planning inference* **16**, 63–74 (1987).
23. J Hodges, E Lehmann, Matching in paired comparisons. *The Annals Math. Stat.* **25**, 787–791 (1954).
24. S Dharmadikari, K Jogdeo, Unimodal laws and related in *A Festschrift For Erich L. Lehmann*. (CRC Press), p. 131 (1982).
25. AY Khintchine, On unimodal distributions. *Izv. Nauchno-Issled. Inst. Mat. Mech.* **2**, 1–7 (1938).
26. S Purkayastha, Simple proofs of two results on convolutions of unimodal distributions. *Stat. & probability letters* **39**, 97–100 (1998).
27. RA Fisher, Moments and product moments of sampling distributions. *Proc. Lond. Math. Soc.* **2**, 199–238 (1930).
28. PR Halmos, The theory of unbiased estimation. *The Annals Math. Stat.* **17**, 34–43 (1946).
29. W Hoeffding, A class of statistics with asymptotically normal distribution. *The Annals Math. Stat.* **19**, 293–325 (1948).
30. PM Heffernan, Unbiased estimation of central moments by using u-statistics. *J. Royal Stat. Soc. Ser. B (Statistical Methodol.)* **59**, 861–863 (1997).
31. D Fraser, Completeness of order statistics. *Can. J. Math.* **6**, 42–45 (1954).
32. AJ Lee, *U-statistics: Theory and Practice*. (Routledge), (2019).
33. J Hodges Jr, E Lehmann, Estimates of location based on rank tests. *The Annals Math. Stat.* **34**, 598–611 (1963).
34. A Ehsanes Saleh, Hodges-lehmann estimate of the location parameter in censored samples. *Annals Inst. Stat. Math.* **28**, 235–247 (1976).
35. RJ Serfling, Generalized L-, m-, and r-statistics. *The Annals Stat.* **12**, 76–86 (1984).
36. 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).
37. MG Akritas, Empirical processes associated with v-statistics and a class of estimators under random censoring. *The Annals Stat.* pp. 619–637 (1986).
38. I Gijbels, P Janssen, N Veraverbeke, Weak and strong representations for trimmed u-statistics. *Probab. theory related fields* **77**, 179–194 (1988).
39. J Choudhury, R Serfling, Generalized order statistics, bahadur representations, and sequential nonparametric fixed-width confidence intervals. *J. Stat. Plan. Inference* **19**, 269–282 (1988).
40. B Efron, Bootstrap methods: Another look at the jackknife. *The Annals Stat.* **7**, 1–26 (1979).
41. PJ Bickel, DA Freedman, Some asymptotic theory for the bootstrap. *The annals statistics* **9**, 1196–1217 (1981).
42. R Helmers, P Janssen, N Veraverbeke, *Bootstrapping U-quantiles*. (CWI. Department of Operations Research, Statistics, and System Theory [BS]), (1990).
43. E Joly, G Lugosi, Robust estimation of u-statistics. *Stoch. Process. their Appl.* **126**, 3760–3773 (2016).
44. P Laforgue, S Cléménçon, P Bertail, On medians of (randomized) pairwise means in *International Conference on Machine Learning*. (PMLR), pp. 1272–1281 (2019).
45. P Daniell, Observations weighted according to order. *Am. J. Math.* **42**, 222–236 (1920).
46. F Mosteller, On some useful "inefficient" statistics. *The Annals Math. Stat.* **17**, 377–408 (1946).
47. CR Rao, *Advanced statistical methods in biometric research*. (Wiley), (1952).
48. PJ Bickel, et al., Some contributions to the theory of order statistics in *Proc. Fifth Berkeley Sympos. Math. Statist. and Probability*. Vol. 1, pp. 575–591 (1967).
49. H Chernoff, JL Gastwirth, MV Johns, Asymptotic distribution of linear combinations of functions of order statistics with applications to estimation. *The Annals Math. Stat.* **38**, 52–72 (1967).
50. L LeCam, On the assumptions used to prove asymptotic normality of maximum likelihood estimates. *The Annals Math. Stat.* **41**, 802–828 (1970).
51. P Bickel, E Lehmann, Descriptive statistics for nonparametric models i. introduction in *Selected Works of EL Lehmann*. (Springer), pp. 465–471 (2012).
52. PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models ii. location in *selected works of EL Lehmann*. (Springer), pp. 473–497 (2012).
53. T Lai, H Robbins, K Yu, Adaptive choice of mean or median in estimating the center of a symmetric distribution. *Proc. Natl. Acad. Sci.* **80**, 5803–5806 (1983).
54. PS Laplace, *Theorie analytique des probabilités*. (1812).
55. E Samuel-Cahn, Combining unbiased estimators. *The Am. Stat.* **48**, 34 (1994).
56. Y Chan, X He, A simple and competitive estimator of location. *Stat. & Probab. Lett.* **19**, 137–142 (1994).
57. G Damlano, P Puig, Efficiency of a linear combination of the median and the sample mean: The double truncated normal distribution. *Scand. J. Stat.* **31**, 629–637 (2004).
58. RD Richtmyer, A non-random sampling method, based on congruences, for "monte carlo" problems, (New York Univ., New York. Atomic Energy Commission Computing and Applied ...), Technical report (1958).
59. IM Sobol', On the distribution of points in a cube and the approximate evaluation of integrals. *Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki* **7**, 784–802 (1967).
60. KA Do, P Hall, Quasi-random resampling for the bootstrap. *Stat. Comput.* **1**, 13–22 (1991).
61. H Cramér, *Mathematical methods of statistics*. (Princeton university press) Vol. 43, (1999).
62. FR Hampel, *Contributions to the theory of robust estimation*. (University of California, Berkeley), (1968).
63. MG Genton, A Lucas, Comprehensive definitions of breakdown points for independent and dependent observations. *J. Royal Stat. Soc. Ser. B (Statistical Methodol.)* **65**, 81–94 (2003).
64. SJ Devlin, R Gnanadesikan, JR Kettenring, Robust estimation and outlier detection with correlation coefficients. *Biometrika* **62**, 531–545 (1975).
65. FR Hampel, A general qualitative definition of robustness. *The annals mathematical statistics* **42**, 1887–1896 (1971).
66. FR Hampel, The influence curve and its role in robust estimation. *J. american statistical association* **69**, 383–393 (1974).
67. PJ Rousseeuw, FR Hampel, EM Ronchetti, WA Stahel, *Robust statistics: the approach based on influence functions*. (John Wiley & Sons), (2011).

- 1036 68. AS Nemirovskij, DB Yudin, *Problem complexity and method efficiency in optimization*. (Wiley-  
1037 Interscience), (1983).
- 1038 69. MR Jerrum, LG Valiant, VV Vazirani, Random generation of combinatorial structures from a  
1039 uniform distribution. *Theor. computer science* **43**, 169–188 (1986).
- 1040 70. N Alon, Y Matias, M Szegedy, The space complexity of approximating the frequency moments  
1041 in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*. pp.  
1042 20–29 (1996).
- 1043 71. DL Donoho, PJ Huber, The notion of breakdown point. *A festschrift for Erich L. Lehmann*  
1044 **157184** (1983).
- 1045 72. M Bieniek, Comparison of the bias of trimmed and winsorized means. *Commun. Stat. Methods*  
1046 **45**, 6641–6650 (2016).
- 1047 73. K Danielak, T Rychlik, Theory & methods: Exact bounds for the bias of trimmed means. *Aust.*  
1048 *& New Zealand J. Stat.* **45**, 83–96 (2003).

DRAFT