

Near-consistent robust estimations of moments for unimodal distributions

Tuban Lee^{a,1}

^aInstitute of Biomathematics, Macau SAR 999078, China

This manuscript was compiled on March 16, 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 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), 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 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 the 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, 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 θ 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 log-

arithmic 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 often more elegant. Since the expectation of an L -statistic can often be expressed as an integral of the quantile function, I -statistics are often analytically obtainable. However, the performance of the aforementioned examples is often inferior to 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 research on semiparametric robust mean estimation, the binomial mean (BM_ϵ) is still inconsistent for any skewed distribution if $\epsilon > 0$, although its asymptotic bias is much smaller than that of the trimmed mean (if $\epsilon \rightarrow 0$, $BM \rightarrow \mu$, since the alternating sum of binomial coefficients is zero). Robust location estimators are typically symmetric due to the prevalence 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). Moving from semiparametrics to parametrics, an ideal robust location estimator would have a non-sample-dependent breakdown point (defined in Subsection F) and be consistent for any symmetric distribution and a skewed distribution with finite second moments. This is called an invariant mean. Based on

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.

¹To whom correspondence should be addressed. E-mail: tl@biomathematics.org

the mean-symmetric weighted average-median inequality, the recombined mean is defined as

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

where d is the key factor for bias correction, $\text{SWA}_{\epsilon,n}$ is $\text{BM}_{\epsilon,n}$ in the first three Subsections, but other symmetric weighted averages can also be employed 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 $\alpha \rightarrow \infty$.*

Proof. Finding d and ϵ 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, $rm_{d,\epsilon} = \lim_{c \rightarrow \infty} \left(\frac{(\text{BM}_{\epsilon} + c)^{d+1}}{(\text{median} + c)^d} - c \right) = (d+1)\text{BM}_{\epsilon} - d\text{median} = \mu$. So, $d = \frac{\mu - \text{BM}_{\epsilon}}{\text{BM}_{\epsilon} - \text{median}}$. The quantile function of the exponential distribution is $Q(p) = \ln\left(\frac{1}{1-p}\right)\lambda$. $E[x] = \lambda$. $E[\text{median}] = Q\left(\frac{1}{2}\right) = \ln 2\lambda$. For the exponential distribution, the expectation of $\text{BM}_{\frac{1}{8}}$ is $E\left[\text{BM}_{\frac{1}{8}}\right] = \lambda\left(1 + \ln\left(\frac{46656}{8575\sqrt{35}}\right)\right)$. Obviously, the scale parameter λ 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, $E[\text{BM}_{\epsilon}] = E[\text{median}] = E[X]$. 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 $d_{\text{Pareto}} = \frac{\frac{\alpha x_m}{\alpha - 1} - x_m(1-p_1)^{-\frac{1}{\alpha}}}{x_m(1-p_1)^{-\frac{1}{\alpha}} - x_m(1-p_2)^{-\frac{1}{\alpha}}}$. 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 illustrate 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 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, $F(x) = Q^{-1}(x) \rightarrow x = Q(p) = \lambda Q_0(p) + \mu$. Thus, any weighted average can be expressed as $\lambda \text{WA}_0(\epsilon) + \mu$, where $\text{WA}_0(\epsilon)$ is an integral of $Q_0(p)$ according to the definition of the weighted average. The simultaneous cancellation of μ and λ in $\frac{(\lambda\mu_0 + \mu) - (\lambda\text{BM}_0(\epsilon) + \mu)}{(\lambda\text{BM}_0(\epsilon) + \mu) - (\lambda\text{median}_0 + \mu)}$ assures that d is a constant. Consequently, the roles of BM_{ϵ} 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

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

provided that $\hat{F}_n(\text{SWA}_{\epsilon,n}) \geq \frac{1}{2}$, 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, $\hat{F}_n(x) := \frac{1}{n} \left(\frac{x - X_{sp}}{X_{sp+1} - X_{sp}} + sp \right)$, where $sp = \sum_{i=1}^n \mathbf{1}_{X_i \leq x}$, $\mathbf{1}_A$ is the indicator of event A . The solution of $\hat{F}_n(\text{SWA}_{\epsilon,n}) < \frac{1}{2}$ is reversing the percentile by $1 - \hat{F}_n(\text{SWA}_{\epsilon,n})$, the obtained percentile is also reversed. Without loss of generality, in the following discussion, only the case where $\hat{F}_n(\text{SWA}_{\epsilon,n}) \geq \frac{1}{2}$ is considered. Moreover, in extreme heavy-tailed distributions, the calculated percentile can exceed the breakdown point of SWA_{ϵ} , so the percentile will be modified to $1 - \epsilon$ if this occurs. The quantile mean employs 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, $d = \frac{F(\mu) - F(\text{BM}_{\epsilon})}{F(\text{BM}_{\epsilon}) - \frac{1}{2}}$. The cdf of the exponential distribution is $F(x) = 1 - e^{-\lambda^{-1}x}$, $\lambda \geq 0$, $x \geq 0$, the expectation of BM_{ϵ} can be expressed as $\lambda \text{BM}_0(\epsilon)$, so $F(\text{BM}_{\epsilon})$ is free of λ . 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, $F(E[\text{BM}_{\epsilon}]) = F(\mu) = \frac{1}{2}$. 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

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

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

The definitions of location and scale parameters dictate that they must satisfy $F(x; \lambda, \mu) = F(\frac{x-\mu}{\lambda}; 1, 0)$. Recall that $x = \lambda Q_0(p) + \mu$, so the percentile of any weighted average is free of λ and μ , 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. $qm_{d \approx 0.321, \epsilon = \frac{1}{8}}$ 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 $qm_{d \approx 0.321, \epsilon = \frac{1}{8}}$ are both consistent mean estimators for any symmetric distribution and a skewed distribution with finite second moments. It's evident that the breakdown points 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 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_ϵ , as long as $|u - \frac{1}{2}| \ll 1-2\epsilon$, and other portions of the SQA function satisfy the inequality constraints that define the ν th orderliness on which the SWA_ϵ is based, the mean-SWA-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 is also close to this parameter, yet the mean-BM $_{\frac{1}{8}}$ -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. In the previous article, 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_ϵ -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-median inequality are generally safe.

B. Robust estimations of the central moments. In 1976, Bickel and Lehmann, in their third paper of the landmark series *Descriptive Statistics for Nonparametric Models* (19), generalized a class of estimators called "measures of disperse," which is now often named as Bickel-Lehmann dispersion. As an example,

they proposed a first version of the trimmed standard deviation, $\hat{\tau}^2(F; \epsilon) \equiv \tau^2(F; \epsilon)$, for independent and identically distributed random variables X with a distribution F , where $\tau^2(F; \epsilon) = \frac{1}{1-2\epsilon} \int_{Q(\epsilon)}^{Q(1-\epsilon)} y dG(y)$, Q is the quantile function of G , G is the distribution of $Y = X^2$. Obviously, when $\epsilon = 0$, the result is equivalent to the second raw moment. In 1979, in the same series (20), they explored another 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 (21), was derived in 1993, but the importance of tackling the symmetry assumption has been greatly underestimated. In the final section of that paper (20), they considered another two possible versions of the trimmed standard deviations, which were modified here for comparison,

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

and

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

where $(X - X')_1 \leq \dots \leq (X - X')_{\binom{n}{2}}$ are the order statistics of the "pseudo-sample". 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 $\epsilon \rightarrow 0$, 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 ϵ is 0, [2] is $\sqrt{2}$ times the standard deviation.

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

Nomenclature. Given a robust estimator $\hat{\theta}$. The first part of the name of the robust statistic defined in this paper is a name that indicates the type of estimator, and the second part is the name of the population parameter θ that the estimator is consistent with as $\epsilon \rightarrow 0$. 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 and not a U -statistic, ϵ is indicated in the subscript of the abbreviation of the estimator. For asymmetric estimators, the corresponding γ is also indicated after ϵ . For weighted U -statistics, the breakdown point of the location estimator is indicated, except the median.

In the previous article on semiparametric robust mean estimation, it was shown that the bias of a reasonable robust estimator should be monotonic with respect to the breakdown point in a semiparametric distribution and, naturally, its name should align with the consistent estimator. The trimmed standard deviation following this nomenclature is $Tsd_{\epsilon, \gamma, n} := \left[TM_{\epsilon, \gamma} \left((\psi_2(X_{N_1}, X_{N_2}))_{N=1}^{\binom{n}{2}} \right) \right]^{-\frac{1}{2}}$, where $TM_{\epsilon, \gamma}(Y)$ denotes

the ϵ, γ -trimmed mean with the sequence $(\psi_2(X_{N_1}, X_{N_2}))_{N=1}^{\binom{n}{2}}$ as an input. Removing the square root yields the trimmed variance $(\text{Var}_{\epsilon, \gamma, n})$. 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 γ -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 continuous unimodal distribution is γ -ordered, if $\gamma \geq 1$.*

Proof. Let $Q(p)$, $0 \leq p \leq 1$, denote the quantile of the continuous unimodal distribution $f_X(x)$. The corresponding probability 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) 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 semiparametric robust mean estimation article, it was proven that a right skewed distribution with a monotonic decreasing pdf is always γ -ordered, if $\gamma \geq 1$, 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 obtain 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 ψ_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 ν th ordered, indicating that orderliness does not require unimodality, e.g., a symmetric bimodal distribution is also ordered. The example from the Weibull distribution shows 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 θ 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 (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 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 gradually became clear that the Hodges-Lehmann 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 k th central moment ($k \leq n$) is defined as,

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

where X_{N_1}, \dots, X_{N_k} are the n choose k elements from X , $\psi_k(x_1, \dots, x_k) = \sum_{j=0}^{k-2} (-1)^j \binom{1}{k-j} \sum (x_{i_1}^{k-j} \dots x_{i_{j+1}}) + (-1)^{k-1} (k-1) x_1 \dots x_k$, 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 $(-\frac{k+1}{2})^{-1}(-\Delta)^k, \frac{1}{k}(-\Delta)^k$.*

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 Ξ_k is equivalent.

Ξ_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 Δ is ξ_Δ . The combination of ξ_Δ , i.e., from $\Delta = 0$ to $Q(0) - Q(1)$, is Ξ_k .

The support of ξ_Δ is the extrema of ψ_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 $i \binom{i-1}{g-j-1} \binom{k-i}{j} + (k-i) \binom{k-i-1}{j-k+g-1} \binom{i}{k-j}$ 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 $i \binom{i-1}{g-j-1} \binom{k-i}{j} + (k-i) \binom{k-i-1}{j-k+g-1} \binom{i}{k-j}$ terms having the form $(-1)^{g+1} \frac{1}{k-g+1} x_1^{k-j} x_k^j$. 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 $\sum_{g=i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} (k-i) \binom{k-i-1}{g-i-1} = (k-i) \int_0^1 \sum_{g=i+1}^{k-1} (-1)^{g+1} \binom{k-i-1}{g-i-1} t^{k-g} dt = (k-i) \int_0^1 ((-1)^i (t-1)^{k-i-1} - (-1)^{k+1}) dt = (k-i) \left(\frac{(-1)^k}{i-k} + (-1)^k \right) = (-1)^{k+1} + (k-i)(-1)^k$ and $\sum_{g=k-i+1}^{k-1} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-k+i-1} = \int_0^1 \sum_{g=k-i+1}^{k-1} (-1)^{g+1} i \binom{i-1}{g-k+i-1} t^{k-g} dt = \int_0^1 (i(-1)^{k-i} (t-1)^{i-1} - i(-1)^{k+1}) dt = (-1)^k (i-1)$. 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 $\sum_{g=j+1}^{i+j} (-1)^{g+1} \frac{1}{k-g+1} i \binom{i-1}{g-j-1} \binom{k-i}{j} = i \binom{k-i}{j} \int_0^1 \sum_{g=j+1}^{i+j} (-1)^{g+1} \binom{i-1}{g-j-1} t^{k-g} dt = \binom{k-i}{j} i \int_0^1 ((-1)^j t^{k-j-1} \left(\frac{t}{t-1} \right)^{1-i}) dt = \binom{k-i}{j} i \frac{(-1)^{j+i+1} \Gamma(i) \Gamma(k-j-i+1)}{\Gamma(k-j+1)} = \frac{(-1)^{j+i+1} i! (k-j-i)! (k-i)!}{(k-j)! j! (k-j-i)!} = (-1)^{j+i+1} \frac{i! (k-i)!}{k!} \frac{k!}{(k-j)! j!} = \binom{k}{i}^{-1} (-1)^{1+i} \binom{k}{j} (-1)^j$. The coefficient of $x_1^i x_k^{k-i}$ in $\binom{k}{i}^{-1} (-1)^{1+i} (x_1 - x_k)^k$ is $\binom{k}{i}^{-1} (-1)^{1+i} \binom{k}{i} (-1)^{k-i} = (-1)^{k+1}$, 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 ψ_k follow directly from the properties of the binomial coefficient. \square

ξ_Δ is closely related to $f_\Xi(\Delta)$, which is the pairwise difference distribution, since the probability density of ξ_Δ can be expressed as $f_{\Xi_k}(\bar{\Delta}|\Delta)$ and $\sum_{\bar{\Delta}=-\binom{k}{3+(-1)^k}}^{\binom{k}{-(-1)^k}} (-1)^{\bar{\Delta}} f_{\Xi_k}(\bar{\Delta}|\Delta) = f_\Xi(\Delta)$. 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 ξ_Δ , Ξ_k is monotonic left and

right around zero. In fact, the median of Ξ_k is also close to zero, as it can be cast as a weighted mean of the medians of ξ_Δ . When Δ is small, all values of ξ_Δ are close to zero, resulting in the median of ξ_Δ being close to zero as well. When Δ is large, the median of ξ_Δ depends on its skewness, but the corresponding weight is much smaller, so even if ξ_Δ is highly skewed, the median of Ξ_k will only be slightly shifted from zero. Denote the median of Ξ_k as m_{Ξ_k} , for the five parametric distributions here, $|m_{\Xi_k}|$ s are all $\leq 0.1\sigma$ for Ξ_3 and Ξ_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 Ξ_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, ξ_Δ will be left-skewed, so, for Ξ_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 all of its central moment kernel distributions are also ordered, it is called completely ordered. Although strict complete orderliness is difficult to prove, even 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 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 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,

$$rk m_{d,\epsilon,n} := (d+1) \text{SW} k m_{\epsilon,n} - d m k m_n,$$

where $\text{SW} k m_{\epsilon,n}$ is using the binomial k th central moment ($B k m_{\epsilon,n}$) here, $m k m_n$ is the median k th central moment. Since $\text{SW} k m_{\epsilon,n}$ is an L -statistic, $rk m_{d,\epsilon,n}$ is an arithmetic I -statistic. Similarly, the quantile will not change after scaling. The quantile k th central moment is thus defined as

$$qk m_{d,\epsilon,n} := \hat{Q}_n \left(\left(p \text{SW} k m_{\epsilon,n} - \frac{1}{2} \right) d + p \text{SW} k m_{\epsilon,n} \right),$$

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 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, 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 (2016) (40), 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 Hodges-Lehmann 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 (40, 41), which is more sophisticated and closer to the median H-L 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. 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 based on trimmed mean are completely different from those based on the median, a contradiction arises. However, such contradictions are possible, for example, for the Weibull distribution, just consider the median and mean, $E[m] = \lambda \sqrt[\gamma]{\ln(2)}$, $E[\mu] = \lambda \Gamma(1 + \frac{1}{\alpha})$, 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. To study the conditions that avoid such scenarios, let the quantile average function of a parametric distribution be denoted as $QA(\epsilon, \gamma, \alpha_1, \dots, \alpha_i, \dots, \alpha_k)$, where α_i represent the parameters of the distribution, then, a distribution is γ -congruent if and only if the sign of $\frac{\partial QA}{\partial \alpha_i}$ remains the same for all $0 \leq \epsilon \leq \frac{1}{1+\gamma}$. 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 $0 \leq \epsilon \leq \frac{1}{1+\gamma}$ implies that all γ -weighted averages change in the same direction as the parameters change, as long as they are not undefined. A distribution is completely γ -congruent if and only if it is γ -congruent and all its central moment kernel distributions are also γ -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. Additionally, distributions with infinite moments can be γ -congruent, since the definition is based on the quantile average, not the sample mean.

The following theorems show the conditions that a distribution is congruent or γ -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. \square

Theorem C.2. *A positive define location-scale distribution with a finite second moment is always γ -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 λ or μ are always positive. By application of the definition, the desired outcome is obtained. \square

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

For the Pareto distribution, $\frac{\partial QA(p, \alpha)}{\partial \alpha} = \frac{x_m(1-p)^{-1/\alpha} \ln(1-p)}{\alpha^2}$. Since $\ln(1-p) < 0$ for all $0 < p < 1$, $(1-p)^{-1/\alpha} > 0$ for all $0 < p < 1$ and $\alpha > 0$, so $\frac{\partial QA(p, \alpha)}{\partial \alpha} < 0$, and therefore $\frac{\partial QA(\epsilon, \gamma, \alpha)}{\partial \alpha} < 0$, the Pareto distribution is γ -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, $\text{erfc}^{-1}(2\epsilon) = -\text{erfc}^{-1}(2-2\epsilon)$, $e^{\mu-\sqrt{2}\sigma\text{erfc}^{-1}(2-2\epsilon)} > e^{\mu-\sqrt{2}\sigma\text{erfc}^{-1}(2\epsilon)}$, $\frac{\partial SQA(\epsilon, \sigma)}{\partial \sigma} > 0$, the lognormal distribution is congruent. Theorem C.1 implies that the generalized Gaussian distribution is congruent. For the Weibull distribution, when α changes from 1 to $\frac{1}{2}$, the average probability density on the left side of the median increases, since $\frac{1}{\lambda \ln(2)} < \frac{1}{\lambda \ln^2(2)}$, 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 γ -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 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. 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. A distribution specified by a shape parameter (denoted as α) and a scale parameter (denoted as λ) is often referred to as a two-parameter distribution. Weibull, gamma, Pareto, lognormal, and generalized Gaussian distributions are all two-parameter unimodal distributions. Since α and λ can be converted to skewness and kurtosis, a two-parametric distribution can also be fully specified by these standardized moments. Moreover, if α or skewness or kurtosis is a constant, the two-parameter distribution is reduced to a single-parameter distribution. The above discussion shows that, due to the invariant property, if a single-parameter distribution is chosen as the consistent distribution, the type of invariant moments and their related weighted moments are given, there should exist a unique k -tuple (d_{im}, \dots, d_{ikm}) calibrated by the distribution and the corresponding kernel distributions generated from this distribution. For a right skewed two-parameter 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 specifying d values for a two-parameter distribution requires only skewness or kurtosis, so the other can also be omitted for simplicity. Since many common two-parameter 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 two-parameter distribution as the consistent distribution is to find solutions for the system of equations

$$\begin{cases} rm(SWA, median, D(|rskew|, rkurt, 1)) = \mu \\ rvar(SWvar, mvar, D(|rskew|, rkurt, 2)) = \mu_2 \\ rtm(SWtm, mtm, D(|rskew|, rkurt, 3)) = \mu_3 \\ rfm(SWfm, mfm, D(|rskew|, rkurt, 4)) = \mu_4 \\ rskew = \frac{\mu_3}{\mu_2} \\ rkurt = \frac{\mu_4}{\mu_2^2} \end{cases}, \quad (638)$$

where μ_2 , μ_3 and μ_4 are the population second, third and fourth central moments. $|rskew|$ and $rkurt$ should be the invariant points of the functions $\varsigma(|rskew|) = \left| \frac{rtm(SWtm, mtm, D(|rskew|, 3))}{rvar(SWvar, mvar, D(|rskew|, 2))^{\frac{3}{2}}} \right|$ 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, as the skewness and kurtosis are interrelated for a two-parameter 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).

Algorithm 1 $rkurt$ for a two-parameter distribution

Input: D ; $SWvar$; $SWfm$; $mvar$; mfm ; $maxit$; δ

Output: $rkurt_{i-1}$

```

1:  $i = 0$ 
2:  $rkurt_i \leftarrow \varkappa(kurtosis_{max}) \triangleright$  Using the maximum kurtosis available in  $D$  as an initial guess.
3: repeat
4:    $i = i + 1$ 
5:    $rkurt_{i-1} \leftarrow rkurt_i$ 
6:    $rkurt_i \leftarrow \varkappa(rkurt_{i-1})$ 
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. $|rskew|$ 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 righted skewed two-parameter distribution whose central moment kernel distributions are all congruent, as long as they are within the domain of D (finite), 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. From the definition of D , $\frac{\varkappa(rkurt)}{rkurt} =$

$$\frac{\frac{\mu_{4,cali} - SWfm_{cali}}{SWfm_{cali} - mfm_{cali}} (SWfm - mfm) + SWfm}{\left(\frac{\mu_{2,cali} - SWvar_{cali}}{SWvar_{cali} - mvar_{cali}} (SWvar - mvar) + SWvar \right)^2}, \quad \text{where} \quad (661)$$

the subscript *cali* indicates that the estimates are from the consistent distribution used to calibrate the d values. According to the property of invariance, assuming

$$\left(\frac{\mu_{2,cali} - SWvar_{cali}}{SWvar_{cali} - mvar_{cali}} (SWvar - mvar) + SWvar \right)^2 > 1, \quad (665)$$

$$\text{then, } \frac{\varkappa(rkurt)}{rkurt} < \frac{\frac{\mu_{4,cali} - SWfm_{cali}}{SWfm_{cali} - mfm_{cali}} (SWfm - mfm) + SWfm}{\frac{\mu_{2,cali} - SWvar_{cali}}{SWvar_{cali} - mvar_{cali}} (SWvar - mvar) + SWvar}} = \frac{\frac{\mu_{4,cali} - SWfm_{cali}}{rkurt} (SWfm - mfm) + \frac{SWfm}{rkurt}}{\left(\frac{\mu_{2,cali} - SWvar_{cali}}{rkurt} (SWvar - mvar) + \frac{SWvar}{rkurt} \right)}. \quad (667)$$

$$\begin{aligned}
& \lim_{rkurt \rightarrow \infty} \left(\frac{\frac{\mu_{4cali} - SWfm_{cali}}{rkurt} (SWfm - mfm) + \frac{SWfm}{rkurt}}{\frac{\mu_{4cali} - SWfm_{cali}}{rkurt} (SWfm_{cali} - mfm_{cali})} \right) \\
&= \lim_{rkurt \rightarrow \infty} \left(\left(\frac{\mu_{4cali} - SWfm_{cali}}{rkurt} \right) \frac{SWfm - mfm}{SWfm_{cali} - mfm_{cali}} \right) = \\
& \lim_{rkurt \rightarrow \infty} \left(\left(\frac{(\mu_{4cali} - SWfm_{cali}) \mu_{2cali}^2}{\mu_{4cali}} \right) \frac{SWfm - mfm}{SWfm_{cali} - mfm_{cali}} \right).
\end{aligned}$$

Since $SWfm_{cali}$ and mfm_{cali} are from the same kernel distribution as $\mu_{4cali} = rkurt \mu_{2cali}^2$, so an increase in μ_{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, so, $\lim_{rkurt \rightarrow \infty} \left(\frac{SWfm - mfm}{SWfm_{cali} - mfm_{cali}} \right) < 1$. Also, according to the property of invariance, assuming $\mu_{2cali} < 1$ (the $\mu_{4cali} = rkurt \mu_{2cali}^2$ can be adjusted while $rkurt$ remains unchanged), then $\lim_{rkurt \rightarrow \infty} \left(\frac{(\mu_{4cali} - SWfm_{cali}) \mu_{2cali}^2}{\mu_{4cali}} \right) < 1$.

Therefore, $\lim_{rkurt \rightarrow \infty} \frac{\kappa(rkurt)}{rkurt} < 1$. 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 $SWfm_{cali} > SWfm$, assuming $\mu_{2cali} < 1$, then $(\mu_{4cali} - SWfm_{cali}) \frac{SWfm - mfm}{SWfm_{cali} - mfm_{cali}} + SWfm < \frac{\mu_{4cali}}{\mu_{2cali}^2} = rkurt$, $\frac{\kappa(rkurt)}{rkurt} < 1$, 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 provided that a fix point exists within the domain of D . At this stage, D can only be approximated through a Monte Carlo study. Continuity can be ensured by using linear interpolation. One common encountered problem is that the domain of D depends 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 can return the optimal boundary value. For quantile moments, the logic is similar, if the percentiles do not exceed the breakdown point. If this is the case, 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 Huber M-estimator, but an advantage of this algorithm is that the optimization is solely related to the d value function 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 common unimodal distributions (Table 1, SI Dataset S1), just using the Weibull distribution as the consistent distribution.

E. Variance. As the fundamental theorem in statistics, the central limit theorem declares 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 robust location estimators (2, 36, 43–50) and it was found that the efficiencies of the robust location estimators can differ significantly from that of the arithmetic mean.

Daniell (1920) stated (43) that comparing the efficiencies of various kinds of estimators is useless unless they all tend to coincide asymptotically. Bickel and Lehmann, also in the landmark series (49, 50), argued that meaningful comparisons can be made by studying the standardized variances, 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 comparable 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, $\widehat{\theta}_{r_1}, \widehat{\theta}_{r_2}, \dots, \widehat{\theta}_{r_{j-1}}$. 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 $\bar{\theta}$ 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 (50) proved that a lower bound for the efficiency of TM_ϵ to sample mean is $(1 - 2\epsilon)^2$ and this monotonic bound holds true 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 results were typically as good as the better of the sample mean and median regarding variance (51). It can be interpreted as an attempt to use the variance version of the mean-SWA-median inequality. While they used bootstrap standard error as the criterion, another approach can be dated back to Laplace (1812) (52) 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 (1994), Chan and He (1994), and Damilano and Puig (2004)'s papers (53–55).

Scaled standard error allows for direct comparison of variances for different location estimators in asymmetric distributions. In this study, twelve possible combinations were created using two invariant means and related symmetric weighted averages ($BM_{\frac{1}{8}}, SQM_{\frac{1}{8}}, BM_{\nu=2, \epsilon=\frac{1}{8}}, WM_{\frac{1}{8}}, BWM_{\frac{1}{8}}$, and $TM_{\frac{1}{8}}$ used here). Each combination has a SSE for a single-parameter distribution, which can be inferred through a Monte Carlo study. Then, the combination with the smallest SSE is chosen (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, its computational cost is high. Similar to Subsection D, let $I(|skewness|, kurtosis, k, dtype, n) = ikm_{WA}$ 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 difference is that unlike D , I is defined to be discontinuous but linear interpolation can also ensure continuity. Using this approach, the result is often very close to the optimum choice (SI Datasets S1).

Due to combinatorial explosion, the bootstrap (56), introduced by Efron in 1979, is indispensable for computing invariant moments in practice. In 1981, Bickel and Freed-

Table 1. Evaluation of invariant moments for five common unimodal distributions in comparison with current popular methods

Errors	TM $\frac{1}{8}$	WM $\frac{1}{8}$	H-L	HM	SM $\frac{1}{9}$	rm $\frac{1}{8}$	qm $\frac{1}{8}$	Tsd $\frac{2}{8}$	rvar $\frac{1}{8}$	qvar $\frac{1}{8}$	rtm $\frac{1}{8}$	qtm $\frac{1}{8}$	rfm $\frac{1}{8}$	qfm $\frac{1}{8}$
WAAB	0.128	0.078	0.109	0.102	0.070	0.002	0.004	0.237	0.031	0.013	0.025	0.009	0.071	0.019
WRMSE	0.131	0.082	0.111	0.105	0.074	0.017	0.019	0.235	0.037	0.025	0.030	0.018	0.073	0.027
WAB $_{n=5400}$	0.128	0.078	0.108	0.102	0.070	0.002	0.005	0.235	0.031	0.013	0.025	0.009	0.070	0.019
WSE \vee WSSE	0.014	0.014	0.014	0.014	0.014	0.017	0.018	0.015	0.016	0.018	0.012	0.014	0.012	0.016
WAAB	0.128	0.078	0.109	0.102	0.070	0.003	0.004	0.237	0.007	0.007	0.009	0.007	0.018	0.014
WRMSE	0.131	0.082	0.111	0.105	0.074	0.018	0.019	0.235	0.020	0.021	0.019	0.020	0.033	0.025
WAB $_{n=5400}$	0.128	0.078	0.108	0.102	0.070	0.003	0.004	0.235	0.007	0.007	0.009	0.008	0.018	0.015
WSE \vee WSSE	0.014	0.014	0.014	0.014	0.014	0.017	0.018	0.015	0.018	0.018	0.014	0.017	0.021	0.017

Errors	m	BWM $\frac{1}{8}$	BM $_{\nu=2, \frac{1}{8}}$	SQM $\frac{1}{8}$	BM $\frac{1}{8}$	\bar{x}	im $_{v, \frac{1}{8}}$	var	ivar $_{v, \frac{1}{8}}$	tm	itm $_{v, \frac{1}{8}}$	fm	ifm $_{v, \frac{1}{8}}$
WAAB	0.205	0.104	0.093	0.057	0.057	0.000	0.003	0.000	0.007	0.000	0.009	0.000	0.017
WRMSE	0.208	0.108	0.096	0.061	0.061	0.014	0.017	0.017	0.020	0.021	0.020	0.028	0.031
WAB $_{n=5400}$	0.205	0.104	0.093	0.057	0.057	0.000	0.003	0.000	0.006	0.000	0.010	0.001	0.016
WSE \vee WSSE	0.016	0.014	0.014	0.015	0.015	0.014	0.016	0.017	0.018	0.021	0.014	0.026	0.020

Errors	\bar{x}	rm	im $_v$	var	var $_{bs}$	rvar	ivar $_v$	tm	tm $_{bs}$	rtm	itm $_v$	fm	fm $_{bs}$	rfm	ifm $_v$
RMSE	0.014	0.016	0.016	0.018	0.017	0.019	0.018	0.021	0.019	0.018	0.018	0.027	0.023	0.023	0.022
SE \vee SSE	0.014	0.016	0.016	0.017	0.017	0.019	0.018	0.020	0.019	0.018	0.018	0.025	0.021	0.022	0.021

The first section of the table presents the use of the exponential distribution as the consistent distribution for five common unimodal distributions: Weibull, gamma, Pareto, lognormal and generalized Gaussian distributions. The second section uses the Weibull distribution as the consistent distribution. The third section uses the Weibull distribution plus optimization (ikm_v is invariant k th moment, variance-optimized). The fourth section uses the Weibull distribution for the Weibull distribution; the breakdown points are all $\frac{1}{8}$ (not indicated). BM $\frac{1}{8}$ is the weighted average used in recombined/quantile moments. The table includes the average asymptotic bias (AAB, as $n \rightarrow \infty$), root mean square error (RMSE, at $n = 5400$), average bias (AB, at $n = 5400$) and variance (SE \vee SSE, at $n = 5400$) of these estimators, all reported in the units of the standard deviations of the distribution or corresponding kernel distributions. The notation bs indicates the quasi-bootstrap central moments. W means that the results were weighted by the number of Google Scholar search results on May 30, 2022 (including synonyms). The calibrations of d values and the computations of AAB, AB, and SSE were described in Subsection E, F and SI Methods. Detailed results and related codes are available in SI Dataset S1 and [GitHub](#).

man (57) 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 were proven by Helmers, Janssen, and Veraverbeke (1990) (58). In the previous article, the advantages of quasi-bootstrap were discussed (59–61). 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 complex deterministic estimator that is not only computationally efficient but also statistical efficient. The only drawback of quasi-bootstrap compared to non-bootstrap is that a small bootstrap size can produce additional finite sample bias but this can be corrected by recalibrating the d values (SI Text). The default bootstrap size is set as 18 thousand here, as it balances computational cost and finite sample bias, except for the asymptotic value calculation. In general, the variances of invariant central moments are much smaller than those of corresponding unbiased sample central moments (deduced by Cramér (62)), except that of the corresponding second central moment (Table 1).

F. Robustness. The measure of robustness to gross errors used in this paper is the breakdown point proposed by Hampel (63) in 1968. Previous work has shown that the median of means (MoM) is asymptotically equivalent to the median Hodge-Lehmann mean. Therefore it is also biased for any asymmetric distribution. Nevertheless, the concentration bound of MoM depends on $\sqrt{\frac{1}{n}}$ (64), so it is quite natural to deduce that it is

a consistent robust estimator. The concept, sample-dependent breakdown point, is defined to avoid ambiguity.

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, ζ 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 (65, 66).

Bear in mind that it differs from the “infinitesimal robustness” defined by Hampel, which is related to whether the asymptotic influence function is bounded (67–69). 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 break down this estimator.

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, which extends the method in Donoho and Huber (1983)’s proof of the breakdown point of the Hodges-Lehmann estimator (70).

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

that as $n \rightarrow \infty$, $k \ll n$, the asymptotic breakdown point of the weighted U -statistic is $1 - (1 - \epsilon)^{\frac{1}{k}}$, where ϵ is the breakdown point of the weighted average.

Proof. According to the definition of ϵ -contamination (70), suppose m contaminants are added to the sample. The fraction of bad values in the sample can be represented as $\epsilon_U = \frac{m}{n+m}$, where n denotes the original number of data points that remain unaffected. In the kernel distribution, $\binom{n}{k}$ out of a total of $\binom{n+m}{k}$ points are not corrupted. Then, the breakdown can be avoided 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 ϵ is the breakdown point of the weighted average, $\frac{1}{2} \geq \epsilon \geq 0$,

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

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 asymptotic 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 ϵ 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 if the population mean is finite. As a result, the sample mean not only has zero robustness to gross errors, but also has zero robustness to departures from the regularity conditions. To compare the performance of estimators under departures from assumptions, it is necessary to impose constraints on these departures, since if departures are unlimited, any estimators can be broken.

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 (71, 72). 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 corresponding semiparametric models. An estimator based on a smaller model is naturally more robust to distributional shift within that model. While this comparison is limited to the smaller semiparametric model and is not universal, quite surprisingly, the results coincide with those obtained from bias bound analysis. Bias bounds are more universal since they can be deduced for distributions with finite moments without assuming unimodality (71, 72). However, 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 greater than that of $SQM_{\frac{1}{8}}$ in the gamma distribution, but it has much smaller average biases (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 asymptotic bias $|\hat{\theta} - \theta|$, where $\hat{\theta}$ is the estimation of θ . 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

$$AAB_{\hat{\theta}} := \frac{1}{C} \sum_{\substack{\delta + \tilde{\mu}_4 \leq \mu_4 \leq C\delta + \tilde{\mu}_4 \\ \tilde{\mu}_4 \text{ is a multiple of } \delta}} E_{\hat{\theta}|\mu_4} [|\hat{\theta} - \theta|]$$

where $\tilde{\mu}_4$ is the kurtosis specifying the two-parameter distribution, $E_{\hat{\theta}|\mu_4}$ denotes the expected value given that the $\tilde{\mu}_4$ is fixed.

Standardization plays a crucial role in comparing the performance 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 preferred because of its central role in standard error estimation. During standardization, all absolute differences in AAB are divided by the standard deviations of the corresponding distributions. 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 weighted central moments. The only difference is that the population standard deviation is replaced by the asymptotic standard deviation of the kernel distribution (σ_{km}).

In Table 1, $\delta = 0.1$, $C = 120$. 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. This setting may seem 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 when 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 come 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. Interestingly,

the asymptotic biases of $TM_{\frac{1}{8}}$ and $WM_{\frac{1}{8}}$, after averaging and weighting, are 0.128σ and 0.078σ , 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), 0.173σ and 0.126σ , for distributions with finite moments without assuming unimodality (71, 72).

Discussion

Moments, including raw moments, central moments, and standardized moments, are key parameters that determine probability distributions. Central moments are much more popular than raw moments because they are invariant to translation. In 1947, Hsu and Robbins proved that the arithmetic mean converges completely to the population mean provided the second moment is finite (73). The strong law of large numbers (proven by Kolmogorov in 1933) (74) implies that the k th sample central moment is asymptotically unbiased. Recently, fascinating statistical phenomena regarding Taylor's law for distributions with infinite moments have been discovered by Drton and Xiao (2016) (75), Pillai and Meng (2016) (76), Cohen, Davis, and Samorodnitsky (2020) (77), and Brown, Cohen, Tang, and Yam (2021) (78). In their inspiring comment to Brown et al's paper (78), Lindquist and Rachev (2021) raised a critical question: "What are the proper measures for the location, spread, asymmetry, and dependence (association) for random samples with infinite mean?" (79). They suggested using median, interquartile range, and medcouple (80) as the robust versions of the first three standardized moments (81–83). This is not the focus of this paper, but it is almost sure that the estimators proposed here will have a place. Since the efficiency of an L -statistic to the sample mean is generally monotonic with respect to the breakdown point (50), and the estimation of central moments can be transformed into a location estimation problem, similar monotonic relations can be expected. For distributions with infinite moments, the desired measures should be as robust as possible. Clearly now, if one wants to preserve the original relationship between each moment while ensuring maximum robustness, the natural choices are median, median variance, and median skewness. Similar to the most robust version of L -moment (84) being trimmed L -moment (85), moments now also have their standard most robust version based on the complete congruence of the underlying distribution.

More generally, parametrics, nonparametrics, and semiparametrics are the current three main branches of statistics. Consistent robust estimation is impossible without specific parametric assumptions. Maximum likelihood was first introduced by Fisher in 1922 (86) in a multinomial model and later generalized by Cramér (1946), Hájek (1970), and Le Cam (1972) (48, 62, 87). Besides Newcomb (2, 3), the general robust estimation of parametric models dates back to 1939, when Wald (88) suggested the use of minimax estimates to solve such problems. Hodges and Lehmann in 1950 (89) expanded upon this concept and obtained minimax estimates for a series of important problems. It was soon clear that a minimax estimator should be a Bayes estimator with regard to the least favorable prior distribution of θ as a minimax estimator is the best in the worst case scenario. Following Huber's seminal work (4), M -statistics have dominated the field of parametric robust statistics for over half a century. In 1984, Bickel ad-

ressed the challenge of robustly estimating the parameters of a linear model while acknowledging the possibility that the model may be invalid but still within the confines of a larger model (90). As suggested by the title *Parametric Robustness: Small Biases can be Worthwhile*, biases exists, but by carefully designing the estimators, they can be very small. The study of semiparametric models was initiated by Stein (91) in 1956. Estimation of the center of symmetry for an unknown symmetric distribution is an important example in his paper. The adaptive estimation of the center of symmetry was studied by van Eeden (1970) and Takeuchi (1971) (92, 93). Bickel, in 1982, simplified the general heuristic necessary condition proposed by Stein (91) (1956) and derived sufficient conditions for this type of problem, adaptive estimation (94). As pointed out by Begun, Hall, Huang, and Wellner (1983) (95), the two problems, semiparametrics (or adaptive estimation) and parametric robustness, are closely related but different. The paradigm shift here opens up the possibility that by defining a large semiparametric model and constructing estimators simultaneously for two or more very different semiparametric/parametric models within the large semiparametric model, then even for a seemingly "wrong" parametric model belongs to the large semiparametric model but not to the semiparametric/parametric models used for calibration, their performance might still be near-optimal due to the common nature shared by the models used by the estimators. The models can be directly expanded and not limited to a single parametric form. Maybe it can be named as comparametrics. Closely related topics are "mixture model" and "constraint defined model" generalized in Bickel, Klaassen, Ritov, and Wellner's classic semiparametric textbook (1993) (96) and the method of sieves, introduced by Grenander in 1981 (97). As building blocks of statistics, invariant moments provide an option for estimating distribution parameters robustly and near-consistently with moderate variances under mild assumptions. This can improve the consistency of statistical results across studies, particularly when heavy-tailed distributions may be present (98–102).

Data Availability. Data for Table 1 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. P.J 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. C Bernard, R Kazzi, S Vanduffel, Range value-at-risk bounds for unimodal distributions under partial information. *Insur. Math. Econ.* **94**, 9–24 (2020).
16. WR van Zwet, *Convex Transformations of Random Variables: Nebst Stellingen.* (1964).
17. 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).
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 Dharmadhikari, 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.* **14**, 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. E Joly, G Lugosi, Robust estimation of u-statistics. *Stoch. Process. their Appl.* **126**, 3760–3773 (2016).
41. P Laforgue, S Cléménçon, P Bertail, On medians of (randomized) pairwise means in *International Conference on Machine Learning.* (PMLR), pp. 1272–1281 (2019).
42. K Pearson, Contributions to the mathematical theory of evolution. *Philos. Transactions Royal Soc. London.* **A 185**, 71–110 (1894).
43. P Daniell, Observations weighted according to order. *Am. J. Math.* **42**, 222–236 (1920).
44. F Mosteller, On some useful “inefficient” statistics. *The Annals Math. Stat.* **17**, 377–408 (1946).
45. CR Rao, *Advanced statistical methods in biometric research.* (Wiley), (1952).
46. 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).
47. 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).
48. L LeCam, On the assumptions used to prove asymptotic normality of maximum likelihood estimates. *The Annals Math. Stat.* **41**, 802–828 (1970).
49. P Bickel, E Lehmann, Descriptive statistics for nonparametric models i. introduction in *Selected Works of EL Lehmann.* (Springer), pp. 465–471 (2012).
50. PJ Bickel, EL Lehmann, Descriptive statistics for nonparametric models ii. location in *selected works of EL Lehmann.* (Springer), pp. 473–497 (2012).
51. 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).
52. PS Laplace, *Théorie analytique des probabilités.* (1812).
53. E Samuel-Cahn, Combining unbiased estimators. *The Am. Stat.* **48**, 34 (1994).
54. Y Chan, X He, A simple and competitive estimator of location. *Stat. & Probab. Lett.* **19**, 137–142 (1994).
55. G Damilano, 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).
56. B Efron, Bootstrap methods: Another look at the jackknife. *The Annals Stat.* **7**, 1–26 (1979).
57. PJ Bickel, DA Freedman, Some asymptotic theory for the bootstrap. *The annals statistics* **9**, 1196–1217 (1981).
58. R Helmers, P Janssen, N Veraverbeke, *Bootstrapping U-quantiles.* (CWI. Department of Operations Research, Statistics, and System Theory [BS]), (1990).
59. 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).
60. 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).
61. KA Do, P Hall, Quasi-random resampling for the bootstrap. *Stat. Comput.* **1**, 13–22 (1991).
62. H Cramér, *Mathematical methods of statistics.* (Princeton university press) Vol. 43, (1999).
63. FR Hampel, *Contributions to the theory of robust estimation.* (University of California, Berkeley), (1968).
64. L Devroye, M Lerasle, G Lugosi, RI Oliveira, Sub-gaussian mean estimators. *The Annals Stat.* **44**, 2695–2725 (2016).
65. 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).
66. SJ Devlin, R Gnanadesikan, JR Kettenring, Robust estimation and outlier detection with correlation coefficients. *Biometrika* **62**, 531–545 (1975).
67. FR Hampel, A general qualitative definition of robustness. *The annals mathematical statistics* **42**, 1887–1896 (1971).
68. FR Hampel, The influence curve and its role in robust estimation. *J. american statistical association* **69**, 383–393 (1974).
69. PJ Rousseeuw, FR Hampel, EM Ronchetti, WA Stahel, *Robust statistics: the approach based on influence functions.* (John Wiley & Sons), (2011).
70. DL Donoho, PJ Huber, The notion of breakdown point. *A festschrift for Erich L. Lehmann* **157184** (1983).
71. M Bieniek, Comparison of the bias of trimmed and winsorized means. *Commun. Stat. Methods* **45**, 6641–6650 (2016).
72. K Danielak, T Rychlik, Theory & methods: Exact bounds for the bias of trimmed means. *Aust. & New Zealand J. Stat.* **45**, 83–96 (2003).
73. PL Hsu, H Robbins, Complete convergence and the law of large numbers. *Proc. national academy sciences* **33**, 25–31 (1947).
74. A Kolmogorov, Sulla determinazione empirica di una legge di distribuzione. *Inst. Ital. Attuari, Giorn.* **4**, 83–91 (1933).
75. M Drton, H Xiao, Wald tests of singular hypotheses. *Bernoulli* **22**, 38–59 (2016).
76. NS Pillai, XL Meng, An unexpected encounter with cauchy and lévy. *The Annals Stat.* **44**, 2089–2097 (2016).
77. JE Cohen, RA Davis, G Samorodnitsky, Heavy-tailed distributions, correlations, kurtosis and Taylor’s law of fluctuation scaling. *Proc. Royal Soc. A* **476**, 20200610 (2020).
78. M Brown, JE Cohen, CF Tang, SCP Yam, Taylor’s law of fluctuation scaling for semivariates and higher moments of heavy-tailed data. *Proc. Natl. Acad. Sci.* **118**, e2108031118 (2021).
79. WB Lindquist, ST Rachev, Taylor’s law and heavy-tailed distributions. *Proc. Natl. Acad. Sci.* **118**, e2118893118 (2021).
80. G Brys, M Hubert, A Struyf, A robust measure of skewness. *J. Comput. Graph. Stat.* **13**, 996–1017 (2004).
81. DC Hoaglin, F Mosteller, JW Tukey, *Exploring data tables, trends, and shapes.* (John Wiley & Sons), (2011).
82. PJ Huber, *Robust statistics.* (Wiley), pp. 309–312 (1981).
83. RA Maronna, RD Martin, VJ Yohai, M Salibián-Barrera, *Robust statistics: theory and methods (with R).* (John Wiley & Sons), (2019).
84. JR Hosking, L-moments: Analysis and estimation of distributions using linear combinations of order statistics. *J. Royal Stat. Soc. Ser. B (Methodological)* **52**, 105–124 (1990).
85. EA Elamir, AH Seheult, Trimmed l-moments. *Comput. Stat. & Data Analysis* **43**, 299–314 (2003).
86. RA Fisher, On the mathematical foundations of theoretical statistics. *Philos. transactions Royal Soc. London. Ser. A, containing papers a mathematical or physical character* **222**, 309–368 (1922).
87. J Hájek, Local asymptotic minimax and admissibility in estimation in *Proceedings of the sixth Berkeley symposium on mathematical statistics and probability.* Vol. 1, pp. 175–194 (1972).
88. A Wald, Contributions to the theory of statistical estimation and testing hypotheses. *The Annals Math. Stat.* **10**, 299–326 (1939).
89. J Hodges, EL Lehmann, Some problems in minimax point estimation in *Selected Works of EL Lehmann.* (Springer), pp. 15–30 (2012).
90. P Bickel, Parametric robustness: small biases can be worthwhile. *The Annals Stat.* **12**, 864–879 (1984).
91. C Stein, et al., Efficient nonparametric testing and estimation in *Proceedings of the third Berkeley symposium on mathematical statistics and probability.* Vol. 1, pp. 187–195 (1956).
92. C Van Eeden, Efficiency-robust estimation of location. *The Annals Math. Stat.* **41**, 172–181 (1970).
93. K Takeuchi, A uniformly asymptotically efficient estimator of a location parameter. *J. Am. Stat. Assoc.* **66**, 292–301 (1971).
94. PJ Bickel, On adaptive estimation. *The Annals Stat.* **10**, 647–671 (1982).
95. JM Begun, WJ Hall, WM Huang, JA Wellner, Information and asymptotic efficiency in parametric-nonparametric models. *The Annals Stat.* **11**, 432–452 (1983).
96. P Bickel, CA Klaassen, Y Ritov, JA Wellner, *Efficient and adaptive estimation for semiparametric models.* (Springer) Vol. 4, (1993).
97. U Grenander, *Abstract Inference.* (1981).
98. JT Leek, RD Peng, Reproducible research can still be wrong: adopting a prevention approach. *Proc. Natl. Acad. Sci.* **112**, 1645–1646 (2015).
99. B Baribault, et al., Metastudies for robust tests of theory. *Proc. Natl. Acad. Sci.* **115**, 2607–2612 (2018).
100. MJ Schuemie, G Hripcsak, PB Ryan, D Madigan, MA Suchard, Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data. *Proc. Natl. Acad. Sci.* **115**, 2571–2577 (2018).
101. P Patil, G Parmigiani, Training replicable predictors in multiple studies. *Proc. Natl. Acad. Sci.* **115**, 2578–2583 (2018).
102. E National Academies of Sciences, et al., *Reproducibility and Replicability in Science.* (National Academies Press), (2019).