

Semiparametric mean estimations based on the orderliness of quantile averages

Tuban Lee

This manuscript was compiled on May 17, 2023

As one of the most fundamental problems in statistics, robust location estimation has many prominent solutions, such as the symmetric trimmed mean, symmetric Winsorized mean, Hodges–Lehmann estimator, Huber M -estimator, and median of means. Recent studies suggest that their maximum biases concerning the mean can be quite different in asymmetric distributions, but the underlying mechanisms and average performance remain largely unclear. In this article, similar to the mean-median-mode inequality, it is proven that in the context of nearly all common unimodal distributions, there exists an orderliness of symmetric quantile averages with different breakdown points. Further deductions explain why the Winsorized mean and median of means typically have smaller biases compared to the trimmed mean. Building on the U -orderliness, the superiority of the median Hodges–Lehmann mean is discussed.

semiparametric | mean-median-mode inequality | asymptotic | unimodal
| Hodges–Lehmann estimator

In 1823, Gauss (1) proved that for any unimodal distribution with a finite second moment, $|m - \mu| \leq \sqrt{\frac{3}{4}}\omega$, where μ is the population mean, m is the population median, and ω is the root mean square deviation from the mode, M . This pioneering work revealed that despite potential bias with respect to the mean in robust estimates, the deviation remains bounded in unit of a scale parameter under certain assumptions. Bernard, Kazzi, and Vanduffel (2020) (2) further derived asymptotic bias bounds of any quantile for unimodal distributions by reducing this optimization problem to a parametric one, which can be solved analytically. They showed that the population median, m , has the smallest maximum distance to the population mean, μ , among all symmetric quantile averages (SQA _{ϵ}). Daniell, in 1920, (3) analyzed a class of estimators, linear combinations of order statistics, and identified that ϵ -symmetric trimmed mean (STM _{ϵ}) belongs to this class. Another popular choice, the ϵ -symmetric Winsorized mean (SWM _{ϵ}), named after Winsor and introduced by Tukey (4) and Dixon (5) in 1960, is also an L -estimator. Bieniek (2016) derived exact bias upper bounds of the Winsorized mean based on Danielak and Rychlik's work (2003) on the trimmed mean for any distribution with a finite second moment and confirmed that the former is smaller than the latter (6, 7). In 1963, Hodges and Lehmann (8) proposed a class of nonparametric location estimators based on rank tests and, from the Wilcoxon signed-rank statistic (9), deduced the median of pairwise means as a robust location estimator for a symmetric population. Both L -statistics and R -statistics achieve robustness essentially by removing a certain proportion of extreme values. In 1964, Huber (10) generalized maximum likelihood estimation to the minimization of the sum of a specific loss function, which measures the residuals between the data points and the model's parameters. Some L -estimators are also M -estimators, e.g., the sample mean is an M -estimator with a squared error loss function, while the sample median is an

M -estimator with an absolute error loss function (10). The Huber M -estimator is obtained by applying the Huber loss function that combines elements of both squared error and absolute error to achieve robustness against gross errors and high efficiency for contaminated Gaussian distributions (10). Sun, Zhou, and Fan (2020) examined the concentration bounds of Huber M -estimator (11). Mathieu (2022) (12) further derived the concentration bounds of M -estimators and demonstrated that, by selecting the tuning parameter which depends on the variance, Huber M -estimator can also be a sub-Gaussian estimator. The concept of median of means (MoM _{$k, b = \frac{n}{k}$} , k is the number of size in each block, b is the number of blocks) was implicitly introduced several times in Nemirovsky and Yudin (1983) (13), Jerrum, Valiant, and Vazirani (1986), (14) and Alon, Matias and Szegedy (1996) (15)'s works. Given its good performance even for distributions with infinite second moments, MoM has received increasing attention over the past decade (16–23). Devroye, Lerasle, Lugosi, and Oliveira (2016) showed that MoM nears the optimum of sub-Gaussian mean estimation with regards to concentration bounds when the distribution has a heavy tail (21). For a comparison of concentration bounds of trimmed mean, Huber M -estimator, median of means and other relevant estimators, readers are directed to Gobet, Lerasle, and Métivier's paper (2022) (24). Laforgue, Clemencon, and Bertail (2019) proposed the median of randomized means (MoRM _{k, b}) (23), wherein, rather than partitioning, an arbitrary number, b , of blocks are built independently from the sample, and showed that MoRM has better non-asymptotic sub-Gaussian property compared to MoM. In fact, asymptotically, the Hodges–Lehmann (H-L) estimator is equivalent to MoM _{$k=2, b = \frac{n}{k}$} and MoRM _{$k=2, b$} , and they can be seen as the pairwise mean distribution is approximated by the sampling without replacement and bootstrap, respectively. For the asymptotic validity, readers are referred to the foundational works of Efron (1979) (25), Bickel and Freedman (1981,

Significance Statement

In 1964, van Zwet introduced the convex transformation order for comparing the skewness of two distributions. This paradigm shift played a fundamental role in defining robust measures of distributions, from spread to kurtosis. Here, rather than the stochastic ordering between two distributions, the orderliness of quantile averages within a distribution is investigated. By classifying distributions through the signs of derivatives, a series of sophisticated robust mean estimators are deduced. Nearly all common nonparametric robust location estimators are found to be special cases thereof.

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

1984) (26, 27), and Helmers, Janssen, and Veraverbeke (1990) (28).

Here, the ϵ, b -stratified mean is defined as

$$SM_{\epsilon, b, n} := \frac{b}{n} \left(\sum_{j=1}^{\frac{b-1}{2b\epsilon}} \sum_{i_j=\frac{(2bj-b-1)n\epsilon}{b-1}+1}^{\frac{(2bj-b+1)n\epsilon}{b-1}} X_{i_j} \right),$$

where $X_1 \leq \dots \leq X_n$ denote the order statistics of a sample of n independent and identically distributed random variables X_1, \dots, X_n . $b \in \mathbb{N}$, $b \geq 3$. The definition was further refined to guarantee the continuity of the breakdown point by incorporating an additional block in the center when $\lfloor \frac{b-1}{2b\epsilon} \rfloor \bmod 2 = 0$, or by adjusting the central block when $\lfloor \frac{b-1}{2b\epsilon} \rfloor \bmod 2 = 1$ (SI Text). If the subscript n is omitted, only the asymptotic behavior is considered. If b is omitted, $b = 3$ is assumed. $SM_{\epsilon, b=3}$ is equal to STM_{ϵ} , when $\epsilon > \frac{1}{6}$. The basic idea of the stratified mean, when $\frac{b-1}{2\epsilon} \in \mathbb{N}$, $b \bmod 2 = 1$ is to distribute the data into $\frac{b-1}{2\epsilon}$ equal-sized non-overlapping blocks according to their order, then further sequentially group these blocks into b equal-sized strata and compute the mean of the middle stratum, which is the median of means of each stratum. In situations where $i \bmod 1 \neq 0$, a potential solution is to generate multiple smaller samples that satisfy the equality by sampling without replacement, and subsequently calculate the mean of all estimations, the details of determining the sample size and sampling times are included in the SI Text. Although the principle is similar to that of the median of means, without the random shift, the result is different from $MoM_{k=\frac{n}{b}, b}$. Additionally, the stratified mean differs from the mean of the sample obtained through stratified sampling methods, introduced by Neymean (1934) (29) or ranked set sampling (30), introduced by McIntyre in 1952, as these sampling methods are designed to obtain more representative samples or improve the efficiency of sample estimates, but the sample mean based on them are not robust. When $b \bmod 2 = 1$, the stratified mean can be regarded as replacing the other equal-sized strata with the middle stratum, which, in principle, is analogous to the Winsorized mean that replaces extreme values with less extreme percentiles. Furthermore, while the bounds confirm that the Winsorized mean and median of means outperform the trimmed mean (6, 7, 21, 24) in worst-case performance, the complexity of bound analysis makes it difficult to achieve a complete and intuitive understanding of these results. Also, a clear explanation for the average performance of them remains elusive. The aim of this paper is to define a series of semi-parametric models using the signs of derivatives, reveal their elegant interrelations and connections to parametric models, and show that by exploiting these models, a set of sophisticated robust mean estimators can be deduced, which have strong robustness to departures from assumptions.

Quantile average and weighted average

The symmetric trimmed mean, symmetric Winsorized mean, and stratified mean are all L -estimators. More specifically, they are symmetric weighted averages, which are defined as

$$SWA_{\epsilon, n} := \frac{\sum_{i=1}^{\lceil \frac{n}{2} \rceil} \frac{X_i + X_{n-i+1}}{2} w_i}{\sum_{i=1}^{\frac{n}{2}} w_i},$$

where w_i s are the weights applied to the symmetric quantile averages according to the definition of the corresponding L -estimators. For example, for the ϵ -symmetric trimmed mean, $w_i = \begin{cases} 0, & i < n\epsilon \\ 1, & i \geq n\epsilon \end{cases}$, provided that $n\epsilon \in \mathbb{N}$. The mean and median are indeed two special cases of the symmetric trimmed mean.

To extend the symmetric quantile average to the asymmetric case, there are two possible definitions for the ϵ, γ -quantile average $QA(\epsilon, \gamma, n)$, i.e.,

$$\frac{1}{2}(\hat{Q}_n(\gamma\epsilon) + \hat{Q}_n(1-\epsilon)), \quad [1]$$

and

$$\frac{1}{2}(\hat{Q}_n(\epsilon) + \hat{Q}_n(1-\gamma\epsilon)), \quad [2]$$

where $\gamma \geq 0$ and $0 \leq \epsilon \leq \frac{1}{1+\gamma}$, $\hat{Q}_n(p)$ is the empirical quantile function. For trimming from both sides, [1] and [2] are equivalent. [1] is assumed in this article unless otherwise specified, since many common asymmetric distributions are right skewed, and [1] allows trimming only from the right side by setting $\gamma = 0$.

Analogously, the weighted average can be defined as

$$WA_{\epsilon, \gamma} := \frac{\int_{\epsilon_0=0}^{\frac{1}{1+\gamma}} QA(\epsilon_0, \gamma) w_{\epsilon_0}}{\int_{\epsilon_0=0}^{\frac{1}{1+\gamma}} w_{\epsilon_0}}.$$

For instance, the ϵ, γ -trimmed mean ($TM_{\epsilon, \gamma}$) is a weighted average with a left trim size of $\gamma\epsilon n$ and a right trim size of ϵn , where $w_{\epsilon_0} = \begin{cases} 0, & \epsilon_0 < \epsilon \\ 1, & \epsilon_0 \geq \epsilon \end{cases}$.

Classifying distributions by the signs of derivatives

Let \mathcal{P}_k denote the set of all distributions over \mathbb{R} whose moments, from the first to the k th, are all finite. Without loss of generality, all classes discussed in the following are subclasses of the nonparametric class of distributions $\mathcal{P}_k^k := \{\text{All continuous distribution } P \in \mathcal{P}_k\}$. Besides fully and smoothly parameterizing by a Euclidean parameter or just assuming regularity conditions, there are many ways to classify distributions. In 1956, Stein initiated the problem of estimating parameters in the presence of an infinite dimensional nuisance shape parameter (31). A notable example discussed in his groundbreaking work was the estimation of the center of symmetry for an unknown symmetric distribution. In 1993, Bickel, Klaassen, Ritov, and Wellner published an influential semiparametrics textbook (32) and systematically classified many common models into three classes: parametric, nonparametric, and semiparametric. However, there is another old and commonly encountered class of distributions that receives little attention in semiparametric literature: the unimodal distribution. It is a very unique semiparametric model because its definition is based on the signs of derivatives, i.e., assuming P is continuous, $(f'(x) > 0 \text{ for } x \leq M) \wedge (f'(x) < 0 \text{ for } x \geq M)$. Let \mathcal{P}_U denote the set of all unimodal distributions. Five parametric distributions in \mathcal{P}_U are detailed as examples here: Weibull, gamma, Pareto, lognormal and generalized Gaussian.

Consider the sign of the derivative of the quantile average with respect to the breakdown point, a right-skewed distribution is called γ -ordered, if and only if

$$\forall 0 \leq \epsilon \leq \frac{1}{1+\gamma}, \frac{\partial \text{QA}_{\epsilon,\gamma}}{\partial \epsilon} \leq 0.$$

The left-skewed case can be obtained by reversing the inequality $\frac{\partial \text{QA}_{\epsilon,\gamma}}{\partial \epsilon} \leq 0$ to $\frac{\partial \text{QA}_{\epsilon,\gamma}}{\partial \epsilon} \geq 0$ and employing the second definition of QA, as given in [2]; for simplicity, it will be omitted in the following discussion. If $\gamma = 1$, the γ -ordered distribution is referred to as ordered.

Furthermore, many common right-skewed distributions are partial bounded, indicating a convex behavior of the QA function when $\epsilon \rightarrow 0$. If assuming convexity further, the second γ -orderliness can be defined as the following for a right-skewed distribution,

$$\forall 0 \leq \epsilon \leq \frac{1}{1+\gamma}, \frac{\partial^2 \text{QA}_{\epsilon,\gamma}}{\partial \epsilon^2} \geq 0 \wedge \frac{\partial \text{QA}_{\epsilon,\gamma}}{\partial \epsilon} \leq 0.$$

Analogously, the ν th γ -orderliness of a right-skewed distribution can be defined as $(-1)^\nu \frac{\partial^\nu \text{QA}_{\epsilon,\gamma}}{\partial \epsilon^\nu} \geq 0 \wedge \dots \wedge -\frac{\partial \text{QA}_{\epsilon,\gamma}}{\partial \epsilon} \geq 0$. If $\gamma = 1$, the ν th γ -orderliness is referred as ν th orderliness. Let \mathcal{P}_O denote the set of all distributions that are ordered and let \mathcal{P}_{O_ν} and $\mathcal{P}_{\gamma O_\nu}$ denote the sets of all distributions that are ν th ordered and ν th γ -ordered, respectively. The following theorems can be used to quickly identify parametric distributions in \mathcal{P}_O , \mathcal{P}_{O_ν} , and $\mathcal{P}_{\gamma O_\nu}$.

Theorem .1. *For any random variable X whose probability distribution function belongs to a location-scale family, the distribution is ν th γ -ordered if and only if the family of probability distributions is ν th γ -ordered.*

Proof. Let Q_0 denote the quantile function of the standard distribution without any shifts or scaling, then, after a location-scale transformation, the quantile function is $Q(p) = \lambda Q_0(p) + \mu$, where λ is the scale parameter, μ is the location parameter. According to the definition of the ν th γ -orderliness, the signs of derivatives of the QA function remain the same after this transformation. Since the location-scale transformation can also be performed inversely, the proof is complete. \square

Theorem .1 shows that in the analytical proof of the ν th γ -orderliness, both the location and scale parameters can be regarded as constants. It is also instrumental in proving other theorems, as illustrated below.

Theorem .2. *Any symmetric distribution with a finite second moment is ν th ordered.*

Proof. Without loss of generality, assuming continuity and $m = 0$, a symmetric distribution is a probability distribution such that for all x , $f(x) = f(-x)$. The cdf of it satisfies $F(x) = 1 - F(-x)$. Let $x = Q(p)$, by definition, $F(Q(p)) = p$, then, $F(Q(p)) = p = 1 - F(-Q(p))$ and $F(Q(1-p)) = 1-p \Leftrightarrow p = 1 - F(Q(1-p))$. Therefore, $F(-Q(p)) = F(Q(1-p))$. Since the cdf is monotonic, $-Q(p) = Q(1-p) \Leftrightarrow Q(p) + Q(1-p) = 0$. As a result, the SQA function is always a horizontal line; the ν th order derivative is zero. The case of $m \neq 0$ follows directly from Theorem .1. \square

As a consequence of Theorem .2 and the fact that generalized Gaussian distribution is symmetric around the median, it is ν th ordered.

Theorem .3. *Any continuous right skewed distribution whose quantile function Q satisfies $Q^{(\nu)}(p) \geq 0 \wedge \dots \wedge Q^{(i)}(p) \geq 0 \wedge Q^{(2)}(p) \geq 0$, $i \bmod 2 = 0$, is ν th γ -ordered, provided that $0 \leq \gamma \leq 1$.*

Proof. Since $(-1)^i \frac{\partial^i \text{QA}_{\epsilon,\gamma}}{\partial \epsilon^i} = \frac{1}{2}((- \gamma)^i Q^i(\gamma \epsilon) + Q^i(1 - \epsilon))$, $0 \leq \epsilon \leq \frac{1}{1+\gamma}$, $1 \leq i \leq \nu$, when $i \bmod 2 = 0$, $(-1)^i \frac{\partial^i \text{QA}_{\epsilon,\gamma}}{\partial \epsilon^i} \geq 0$ for all $\gamma \geq 0$. When $i \bmod 2 = 1$, if further assuming $0 \leq \gamma \leq 1$, $(-1)^i \frac{\partial^i \text{QA}_{\epsilon,\gamma}}{\partial \epsilon^i} \geq 0$, since $Q^{(i+1)}(\epsilon) \geq 0$. \square

It is now trivial to prove that the Pareto distribution follows the ν th γ -orderliness, provided that $0 \leq \gamma \leq 1$, since the quantile function of the Pareto distribution is $Q(p) = x_m(1-p)^{-\frac{1}{\alpha}}$, where $x_m > 0$, $\alpha > 0$, so $Q^{(\nu)}(p) \geq 0$ according to the chain rule.

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

ACKNOWLEDGMENTS. I sincerely acknowledge the insightful comments from the editor which considerably elevated the lucidity and merit of this paper.

1. CF Gauss, *Theoria combinationis observationum erroribus minimis obnoxiae*. (Henricus Dieterich), (1823).
2. C Bernard, R Kazzi, S Vanduffel, Range value-at-risk bounds for unimodal distributions under partial information. *Insur. Math. Econ.* **94**, 9–24 (2020).
3. P Daniell, Observations weighted according to order. *Am. J. Math.* **42**, 222–236 (1920).
4. JW Tukey, A survey of sampling from contaminated distributions in *Contributions to probability and statistics*. (Stanford University Press), pp. 448–485 (1960).
5. WJ Dixon, Simplified Estimation from Censored Normal Samples. *The Annals Math. Stat.* **31**, 385–391 (1960).
6. M Bieniek, Comparison of the bias of trimmed and winsorized means. *Commun. Stat. Methods* **45**, 6641–6650 (2016).
7. K Danielak, T Rychlik, Theory & methods: Exact bounds for the bias of trimmed means. *Aust. & New Zealand J. Stat.* **45**, 83–96 (2003).
8. J Hodges Jr, E Lehmann, Estimates of location based on rank tests. *The Annals Math. Stat.* **34**, 598–611 (1963).
9. F Wilcoxon, Individual comparisons by ranking methods. *Biom. Bull.* **1**, 80–83 (1945).
10. PJ Huber, Robust estimation of a location parameter. *Ann. Math. Stat.* **35**, 73–101 (1964).
11. Q Sun, WX Zhou, J Fan, Adaptive huber regression. *J. Am. Stat. Assoc.* **115**, 254–265 (2020).
12. T Mathieu, Concentration study of m-estimators using the influence function. *Electron. J. Stat.* **16**, 3695–3750 (2022).
13. AS Nemirovskij, DB Yudin, *Problem complexity and method efficiency in optimization*. (Wiley-Interscience), (1983).
14. MR Jerrum, LG Valiant, VV Vazirani, Random generation of combinatorial structures from a uniform distribution. *Theor. computer science* **43**, 169–188 (1986).
15. N Alon, Y Matias, M Szegedy, The space complexity of approximating the frequency moments in *Proceedings of the twenty-eighth annual ACM symposium on Theory of computing*. pp. 20–29 (1996).
16. PL Bühlmann, Bagging, subbagging and bragging for improving some prediction algorithms in *Research report/Seminar für Statistik, Eidgenössische Technische Hochschule (ETH)*. (Seminar für Statistik, Eidgenössische Technische Hochschule (ETH), Zürich), Vol. 113, (2003).
17. JY Audibert, O Catoni, Robust linear least squares regression. *The Annals Stat.* **39**, 2766–2794 (2011).
18. D Hsu, S Sabato, Heavy-tailed regression with a generalized median-of-means in *International Conference on Machine Learning*. (PMLR), pp. 37–45 (2014).
19. S Minsker, Geometric median and robust estimation in banach spaces. *Bernoulli* **21**, 2308–2335 (2015).
20. C Brownlees, E Joly, G Lugosi, Empirical risk minimization for heavy-tailed losses. *The Annals Stat.* **43**, 2507–2536 (2015).
21. L Devroye, M Lerasle, G Lugosi, RI Oliveira, Sub-gaussian mean estimators. *The Annals Stat.* **44**, 2695–2725 (2016).
22. E Joly, G Lugosi, Robust estimation of u-statistics. *Stoch. Process. their Appl.* **126**, 3760–3773 (2016).
23. P Laforgue, S Cléménçon, P Bertail, On medians of (randomized) pairwise means in *International Conference on Machine Learning*. (PMLR), pp. 1272–1281 (2019).
24. E Gobet, M Lerasle, D Métivier, Mean estimation for Randomized Quasi Monte Carlo method. working paper or preprint (2022).
25. B Efron, Bootstrap methods: Another look at the jackknife. *The Annals Stat.* **7**, 1–26 (1979).

- 275 26. PJ Bickel, DA Freedman, Some asymptotic theory for the bootstrap. *The annals statistics* **9**,
276 1196–1217 (1981).
- 277 27. PJ Bickel, DA Freedman, Asymptotic normality and the bootstrap in stratified sampling. *The*
278 *annals statistics* **12**, 470–482 (1984).
- 279 28. R Helmers, P Janssen, N Veraverbeke, *Bootstrapping U-quantiles*. (CWI. Department of
280 Operations Research, Statistics, and System Theory [BS]), (1990).
- 281 29. J Neyman, On the two different aspects of the representative method: The method of stratified
282 sampling and the method of purposive selection. *J. Royal Stat. Soc.* **97**, 558–606 (1934).
- 283 30. G McIntyre, A method for unbiased selective sampling, using ranked sets. *Aust. journal*
284 *agricultural research* **3**, 385–390 (1952).
- 285 31. C Stein, , et al., Efficient nonparametric testing and estimation in *Proceedings of the third*
286 *Berkeley symposium on mathematical statistics and probability*. Vol. 1, pp. 187–195 (1956).
- 287 32. P Bickel, CA Klaassen, Y Ritov, JA Wellner, *Efficient and adaptive estimation for semiparamet-*
288 *ric models*. (Springer) Vol. 4, (1993).

DRAFT