# Semiparametric robust mean estimation based on the orderliness of quantile averages

**Tuban Lee**[a,1]

[a]Institute of Biomathematics, Macau SAR 999078, China

**As arguably the most fundamental problem in statistics, nonparametric robust location estimation has many prominent solutions, such as the trimmed mean, Winsorized mean, Hodges–Lehmann estimator, and median of means. Recent research suggests that their biases with respect to mean can be quite different in asymmetric distributions. Here, 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 generally have smaller biases compared to the trimmed mean. Building on the $U$-orderliness, the supremacy of weighted 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 $\mu$ is the population mean, $m$ is the population median, $\omega$ is the root mean square deviation from the mode, $M$. Bernard, Kazzi, and Vanduffel (2020) [2] derived bias bounds for the $\epsilon$-symmetric quantile average (SQA$_\epsilon$) for unimodal distributions, building on the works of Karlin and Novikoff (1963) and Li, Shao, Wang, and Yang (2018) [3, 4]. They showed that the $m$ has the smallest maximum distance to the $\mu$ among all symmetric quantile averages. Daniell, in 1920, [5] analyzed a class of estimators, linear combinations of order statistics, and identified that $\epsilon$-symmetric trimmed mean (TM$_\epsilon$) belongs to this class. Another popular choice, the $\epsilon$-symmetric Winsorized mean (WM$_\epsilon$), which was named after Winsor and introduced by Tukey [6] and Dixon [7] in 1960, is also an $L$-statistic. Without assuming unimodality, Bieniek (2016) derived exact bias upper bounds of the Winsorized mean based on Danielak and Rychlik's work (2003) on the trimmed mean and confirmed that the former is smaller than the latter [8, 9]. In 1963, Hodges and Lehmann [10] proposed a class of nonparametric location estimators based on rank tests and, from the Wilcoxon signed-rank statistic [11], deduced the median of pairwise means as a robust location estimator for a symmetric population. The concept of median of means (MoM$_{k,b}$, $k$ is the number of size in each block, $b$ is the number of blocks) was implicit several times in Nemirovsky and Yudin (1983) [12], Jerrum, Valiant, and Vazirani (1986), [13] and Alon, Matias and Szegedy (1996) [14]'s works. Having good performance even for distributions with infinite second moments, the advantages of MoM have received increasing attention over the past decade [15–22]. Devroye, Lerasle, Lugosi, and Oliveira (2016) showed that MoM nears the optimum of nonparametric mean estimation with regards to concentration bounds when the distribution has a heavy tail [20]. In fact, asymptotically, the Hodges-Lehmann (H-L) estimator is equivalent to MoM$_{k=2, b=\frac{n}{k}}$, and it can be seen as the pairwise mean

distribution is approximated by the bootstrap and sampling without replacement, respectively (for the asymptotic validity, the reader is referred to the foundational works of Efron (1979) [23], Bickel and Freedman [24, 25], and Helmers, Janssen, and Veraverbeke (1990) [26]).

Here, the $\epsilon$,$b$-stratified mean is defined as

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

where $X_1 \leq ... \leq X_n$ denote the order statistics of a sample of $n$ independent and identically distributed random variables $X_1, ..., X_n$, $\epsilon \bmod \frac{2}{b-1} = 0$, $\frac{1}{\epsilon} \geq 9$. If the subscript $n$ is omitted, only the asymptotic behavior is considered. If $b$ is omitted, $b = 3$ is assumed. The basic idea is to distribute the random variables into $b$ blocks according to their order, and then compute the mean of the middle block, which is the median of all $b$ blocks. Although the principle is similar to the median of means, without the random shift, the result is different from $\mathrm{MoM}_{k=\frac{n}{b},b}$. The exact solution for $n \bmod \frac{b-1}{2\epsilon} \neq 0$ is imputing the remaining values with multiple hot deck imputation (proposed by Little and Rubin in 1986) [27], since it preserves the original distribution (proven by Reilly in 1991) [28]. If $n \bmod \frac{b-1}{2\epsilon} = \varrho$, the algorithm should run $\binom{n}{\varrho}$ times. An approximation solution is randomly imputing the remaining values several times and then computing the mean of all estimations. The stratified mean is a type of stratum mean which is related to the stratified sampling. The most similar version was proposed by Takahasi and Wakimoto in 1968 [29], which is stratifying order statistics into several non-overlapping blocks and then computing the mean of one block. The median of means and stratified mean are consistent

**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 inequalities, a series of sophisticated robust mean estimators are deduced. Nearly all common nonparametric robust location estimators are found to be special cases thereof.

mean estimators if their asymptotic breakdown points are zero. However, if $\epsilon = \frac{1}{9}$, the biases of the $\text{SM}_{\frac{1}{9}}$ are nearly identical to those of the $\text{WM}_{\frac{1}{9}}$ in asymmetric distributions (Figure **??**, if no other subscripts, $\epsilon$ is omitted for simplicity), i.e., their robustness to departures from the symmetry assumption is similar in practice. More importantly, the bounds confirm that the worst-case performances of $\text{WM}_\epsilon$ are better than those of $\text{TM}_\epsilon$ in terms of bias, but due to the complexity, any extensions are extremely difficult. The aim of this paper is to define a series of semiparametric models using inequalities, demonstrate their elegant interrelations and connections to parametric models, and deduce a set of sophisticated robust mean estimators.

## Quantile average and weighted average

$\epsilon$-symmetric trimmed mean, $\epsilon$-symmetric Winsorized mean, and $\epsilon$-stratified mean are all $L$-statistics. More specifically, they are symmetric weighted averages, which is defined as

$$\text{SWA}_{\epsilon,n} := \frac{\sum_{i=1}^{\frac{n}{2}} \frac{X_i + X_{n-i+1}}{2} w_i}{\sum_{i=1}^{\frac{n}{2}} w_i},$$

where $w_i$s are the weights applied to the symmetric quantile average according to the definition of the corresponding $L$-statistic. For example, for the $\epsilon$-symmetric trimmed mean, $w_i = \begin{cases} 0, & i < n\epsilon \\ 1, & i \geq n\epsilon \end{cases}$. Mean $(\lim_{\epsilon \to 0} \text{TM}_\epsilon)$ and median $(\text{TM}_{\frac{1}{2}})$ are two special cases of symmetric trimmed mean. In 1974, Hogg investigated asymmetric trimmed mean and found its advantages for some special applications ([30]). To extend to the asymmetric case, the quantile average can be defined as

$$\text{QA}(\epsilon, \gamma, n) := \frac{1}{2}(\hat{Q}_n(\epsilon) + \hat{Q}_n(1 - \gamma\epsilon)).$$

where $\gamma > 0$ and $\epsilon \leq \frac{1}{1+\gamma}$, $\hat{Q}_n(p)$ is the empirical quantile function. For example, $\text{QA}(\epsilon = 0.2, \gamma = 2, n) = \frac{1}{2}(\hat{Q}_n(0.2) + \hat{Q}_n(0.6))$. Symmetric quantile average is a special case of quantile average when $\gamma = 1$.

Analogously, weighted average can be defined as

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

Converting this asymptotic definition to finite sample definition requires rounding the $n\epsilon_0$, for simplicity, only asymptotic definition is considered here. For example, the $\epsilon, \gamma$-asymmetric trimmed mean $(\text{TM}_{\epsilon,\gamma})$ is a weighted average that trims the left side $\epsilon$ and trims the right side $\gamma\epsilon$, where $w_{\epsilon_0} = \begin{cases} 0, & \epsilon_0 < \epsilon \\ 1, & \epsilon_0 \geq \epsilon \end{cases}$. Noted that a weighted average is an $L$-statistic, but an $L$-statistic might not be a weighted average, because all quantile averages have the same $\gamma$ in a weighted average. For the sake of brevity, in the following text, if $\gamma$ is not indicated, symmetry will be assumed.

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

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. S Karlin, A Novikoff, Generalized convex inequalities. *Pac. J. Math.* **13**, 1251–1279 (1963).
4. L Li, H Shao, R Wang, J Yang, Worst-case range value-at-risk with partial information. *SIAM J. on Financial Math.* **9**, 190–218 (2018).
5. P Daniell, Observations weighted according to order. *Am. J. Math.* **42**, 222–236 (1920).
6. JW Tukey, A survey of sampling from contaminated distributions in *Contributions to probability and statistics.* (Stanford University Press), pp. 448–485 (1960).
7. WJ Dixon, Simplified Estimation from Censored Normal Samples. *The Annals Math. Stat.* **31**, 385 – 391 (1960).
8. M Bieniek, Comparison of the bias of trimmed and winsorized means. *Commun. Stat. Methods* **45**, 6641–6650 (2016).
9. K Danielak, T Rychlik, Theory & methods: Exact bounds for the bias of trimmed means. *Aust. & New Zealand J. Stat.* **45**, 83–96 (2003).
10. J Hodges Jr, E Lehmann, Estimates of location based on rank tests. *The Annals Math. Stat.* **34**, 598–611 (1963).
11. F Wilcoxon, Individual comparisons by ranking methods. *Biom. Bull.* **1**, 80–83 (1945).
12. AS Nemirovskij, DB Yudin, *Problem complexity and method efficiency in optimization.* (Wiley-Interscience), (1983).
13. MR Jerrum, LG Valiant, VV Vazirani, Random generation of combinatorial structures from a uniform distribution. *Theor. computer science* **43**, 169–188 (1986).
14. 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).
15. PL Bühlmann, Bagging, subagging 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).
16. JY Audibert, O Catoni, Robust linear least squares regression. *The Annals Stat.* **39**, 2766–2794 (2011).
17. D Hsu, S Sabato, Heavy-tailed regression with a generalized median-of-means in *International Conference on Machine Learning.* (PMLR), pp. 37–45 (2014).
18. S Minsker, Geometric median and robust estimation in banach spaces. *Bernoulli* **21**, 2308–2335 (2015).
19. C Brownlees, E Joly, G Lugosi, Empirical risk minimization for heavy-tailed losses. *The Annals Stat.* **43**, 2507–2536 (2015).
20. L Devroye, M Lerasle, G Lugosi, RI Oliveira, Sub-gaussian mean estimators. *The Annals Stat.* **44**, 2695–2725 (2016).
21. E Joly, G Lugosi, Robust estimation of u-statistics. *Stoch. Process. their Appl.* **126**, 3760–3773 (2016).
22. P Laforgue, S Clémençon, P Bertail, On medians of (randomized) pairwise means in *International Conference on Machine Learning.* (PMLR), pp. 1272–1281 (2019).
23. B Efron, Bootstrap methods: Another look at the jackknife. *The Annals Stat.* **7**, 1–26 (1979).
24. PJ Bickel, DA Freedman, Some asymptotic theory for the bootstrap. *The annals statistics* **9**, 1196–1217 (1981).
25. PJ Bickel, DA Freedman, Asymptotic normality and the bootstrap in stratified sampling. *The annals statistics* **12**, 470–482 (1984).
26. R Helmers, P Janssen, N Veraverbeke, *Bootstrapping U-quantiles.* (CWI. Department of Operations Research, Statistics, and System Theory [BS]), (1990).
27. RJ Little, DB Rubin, *Statistical analysis with missing data.* (John Wiley & Sons) Vol. 793, (2019).
28. M Reilly, *Semi-parametric methods of dealing with missing or surrogate covariate data.* (University of Washington), (1991).
29. K Takahasi, K Wakimoto, On unbiased estimates of the population mean based on the sample stratified by means of ordering. *Annals institute statistical mathematics* **20**, 1–31 (1968).
30. RV Hogg, Adaptive robust procedures: A partial review and some suggestions for future applications and theory. *J. Am. Stat. Assoc.* **69**, 909–923 (1974).