Semiparametric robust mean estimations based on the orderliness of quantile averages

Tuban Lee

12

20

21

22

23

27

31

This manuscript was compiled on June 4, 2023

As one of the most fundamental problems in statistics, the 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. This study establishes several forms of orderliness among quantile averages, similar to the mean-median-mode inequality, within a wide range of semi-parametric distributions, particularly highlighting the unique role of γ -symmetric distributions. From this, a sequence of advanced robust mean estimators emerges, which also explains why the Winsorized mean and median of means typically have smaller biases compared to the trimmed mean. Building on the γ -U-orderliness, the superiority of the median Hodges–Lehmann mean is discussed.

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

n 1823, Gauss (1) proved that for any unimodal distribution, $|m-\mu| \leq \sqrt{\frac{3}{4}}\omega$ and $\sigma \leq \omega \leq 2\sigma$, where μ is the population mean, m is the population median, ω is the root mean square deviation from the mode, and σ is the population standard deviation. This pioneering work revealed that despite potential bias in robust mean estimates, the deviation remains bounded in units of a scale parameter under certain assumptions. Bernard, Kazzi, and Vanduffel (2020) (2) further derived asymptotic bias bounds of any quantile for unimodal distributions with finite second moments, by reducing this optimization problem to a parametric one, which can be solved analytically. They showed that m has the smallest maximum distance to μ among all symmetric quantile averages (SQA_c). Daniell, in 1920, (3) analyzed a class of estimators, linear combinations of order statistics, and identified that the ϵ -symmetric trimmed mean (STM_{ϵ}) belongs to this class. Another popular choice, the ϵ -symmetric Winsorized mean (SWM $_{\epsilon}$), 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 signedrank 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, 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 the 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, the Huber M-estimator can also be a sub-Gaussian estimator. The concept of the median of means $(MoM_{k,b=\frac{n}{k},n})$ was first introduced by Nemirovsky and Yudin (1983) in their work on stochastic optimization (13). Given its good performance even for distributions with infinite second moments, the MoM has received increasing attention over the past decade (14-17). Devroye, Lerasle, Lugosi, and Oliveira (2016) showed that $MoM_{k,b=\frac{n}{h},n}$ nears the optimum of sub-Gaussian mean estimation with regards to concentration bounds when the distribution has a heavy tail (15). Laforgue, Clemencon, and Bertail (2019) proposed the median of randomized means $(MoRM_{k,b,n})$ (16), wherein, rather than partitioning, an arbitrary number, b, of blocks are built independently from the sample, and showed that $MoRM_{k,b,n}$ has a better nonasymptotic sub-Gaussian property compared to $MoM_{k,b=\frac{n}{r},n}$. In fact, asymptotically, the Hodges-Lehmann (H-L) estimator is equivalent to $MoM_{k=2,b=\frac{n}{h}}$ 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. When $k \ll n$, the difference between sampling with replacement and without replacement is negligible. For the asymptotic validity, readers are referred to the foundational works of Efron (1979) (18), Bickel and Freedman (1981, 1984) (19, 20), and Helmers, Janssen, and Veraverbeke (1990) (21).

35

36

37

41

42

43

44

45

46

47

50

51

52

53

57

59

60

61

62

64

65

66

67

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, instead of examining 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 is 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.

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

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

69

70

71

73

75

77

81

82

83

84

85

86

87

88

89

90

91

92

93 94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

111

112

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

where $X_1 \leq \ldots \leq X_n$ denote the order statistics of a sample of n independent and identically distributed random variables X_1, \ldots, 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 $\lfloor \frac{b-1}{2b\epsilon} \rfloor \mod 2 = 0$, or by adjusting the central block when $\lfloor \frac{b-1}{2b\epsilon} \rfloor \mod 2 = 1$ (SI Text). If the subscript n is omitted, only the asymptotic behavior is considered. If b is omitted, b = 3 is assumed. $\mathrm{SM}_{\epsilon,b=3}$ is equivalent to STM_{ϵ} , when $\epsilon > \frac{1}{6}$. When $\frac{b-1}{2\epsilon} \in \mathbb{N}$ and $b \mod 2 = 1$, the basic idea of the stratified mean 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 $i \mod 1 \neq 0$, 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 smaller sample size and the number of sampling times are provided in the SI Text. Although the principle resembles that of the median of means, $SM_{\epsilon,b,n}$ is different from $MoM_{k=\frac{n}{L},b,n}$ as it does not include the random shift. Additionally, the stratified mean differs from the mean of the sample obtained through stratified sampling methods, introduced by Nevman (1934) (22) or ranked set sampling (23), introduced by McIntyre in 1952, as these sampling methods aim to obtain more representative samples or improve the efficiency of sample estimates, but the sample means based on them are not robust. When $b \mod 2 = 1$, 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, 15) 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 semiparametric 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 mean estimators can be deduced, which exhibit 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

$$\mathrm{SWA}_{\epsilon,n} \coloneqq \frac{\sum_{i=1}^{\lceil \frac{n}{2} \rceil} \frac{X_i + X_{n-i+1}}{2} w_i}{\sum_{i=1}^{\lceil \frac{n}{2} \rceil} w_i},$$

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,

 $w_i = \begin{cases} 0, & i < n\epsilon \\ 1, & i \geq n\epsilon \end{cases} \text{ when } n\epsilon \in \mathbb{N}. \text{ The mean and median are}$

indeed two special cases of the symmetric trimmed mean.

To extend the symmetric quantile average to the asymmetric case, two definitions for the ϵ, γ -quantile average (QA_{ϵ, γ, n}) are proposed. The first definition is:

$$\frac{1}{2}(\hat{Q}_n(\gamma\epsilon) + \hat{Q}_n(1-\epsilon)), \qquad [1] \quad {}_{122}$$

118

121

127

128

131

132

134

137

139

140

141

142

143

144

145

146

147

148

149

150

152

153

156 157

159

160

162 163

164 165

166

167

168

169

170

and the second definition is:

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

where $\hat{Q}_n(p)$ is the empirical quantile function; γ is used to adjust the degree of asymmetry, $\gamma \geq 0$; and $0 \leq \epsilon \leq \frac{1}{1+\gamma}$. For trimming from both sides, [1] and [2] are essentially equivalent. The first definition along with $\gamma \geq 0$ and $0 \leq \epsilon \leq \frac{1}{1+\gamma}$ are assumed in the rest of 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

$$\mathrm{WA}_{\epsilon,\gamma,n} \coloneqq \frac{\int_0^{\frac{1}{1+\gamma}} \mathrm{QA}\left(\epsilon_0,\gamma,n\right) w(\epsilon_0) d\epsilon_0}{\int_0^{\frac{1}{1+\gamma}} w(\epsilon_0) d\epsilon_0}.$$

For any weighted average, if γ is omitted, it is assumed to be 1. The ϵ, γ -trimmed mean $(TM_{\epsilon, \gamma, n})$ is a weighted average with a left trim size of $n\gamma\epsilon$ and a right trim size of $n\epsilon$.

where
$$w(\epsilon_0) = \begin{cases} 0, & \epsilon_0 < \epsilon \\ 1, & \epsilon_0 \ge \epsilon \end{cases}$$
. Using this definition, regard-

less of whether $n\gamma\epsilon \notin \mathbb{N}$ or $n\epsilon \notin \mathbb{N}$, the TM computation remains the same, since this definition is based on the empirical quantile function. However, in this article, considering the computational cost in practice, non-asymptotic definitions of various types of weighted averages are primarily based on order statistics. Unless stated otherwise, the solution to their decimal issue is the same as that in SM.

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.

- CF Gauss, Theoria combinationis observationum erroribus minimis obnoxiae. (Henricus Dieterich), (1823).
- 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).
- JW Tukey, A survey of sampling from contaminated distributions in Contributions to probability and statistics. (Stanford University Press), pp. 448–485 (1960).
- WJ Dixon, Simplified Estimation from Censored Normal Samples. The Annals Math. Stat. 31 385 – 301 (1960)
- K Danielak, T Rychlik, Theory & methods: Exact bounds for the bias of trimmed means. Aus. & New Zealand J. Stat. 45, 83–96 (2003).
- M Bieniek, Comparison of the bias of trimmed and winsorized means. Commun. Stat. Methods 45, 6641–6650 (2016).
- 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)
- T Mathieu, Concentration study of m-estimators using the influence function. *Electron. J. Stat.* 16, 2605–2750 (2022).
- AS Nemirovskij, DB Yudin, Problem complexity and method efficiency in optimization. (Wiley Interscience), (1983).

2 | Lee

- 14. D Hsu, S Sabato, Heavy-tailed regression with a generalized median-of-means in *International* 171 Conference on Machine Learning. (PMLR), pp. 37–45 (2014).

 15. L Devroye, M Lerasle, G Lugosi, RI Oliveira, Sub-gaussian mean estimators. The Annals Stat. 172
- 173 **44**, 2695–2725 (2016). 174
- 16. P Laforgue, S Clémençon, P Bertail, On medians of (randomized) pairwise means in Interna-175 16. P Latorgue, S Dietriengon, P Detrain, On medians of gardenized, partition of median management tional Conference on Machine Learning. (PMLR), pp. 1272–1281 (2019).
 17. G LECUÉ, M LERASLE, Robust machine learning by median-of-means: Theory and practice. 176

177

178

- The Annals Stat. 48, 906-931 (2020).
- 18. B Efron, Bootstrap methods: Another look at the jackknife. The Annals Stat. 7, 1–26 (1979). 179
- 19. PJ Bickel, DA Freedman, Some asymptotic theory for the bootstrap. The annals statistics 9, 180 1196-1217 (1981). 181
- 182 $20. \ \ \, \text{PJ Bickel}, \, \text{DA Freedman}, \, \text{Asymptotic normality and the bootstrap in stratified sampling}. \, \, \textit{The} \,$ 183 annals statistics 12, 470-482 (1984).
- 21. R Helmers, P Janssen, N Veraverbeke, Bootstrapping U-quantiles. (CWI. Department of 184 Operations Research, Statistics, and System Theory [BS]), (1990). 185
- 186 22. J Neyman, On the two different aspects of the representative method: The method of stratified sampling and the method of purposive selection. J. Royal Stat. Soc. 97, 558-606 (1934). 187
- 188 23. G McIntyre, A method for unbiased selective sampling, using ranked sets. Aust. journal 189 agricultural research 3, 385-390 (1952).