

Comparing Estimators of Quartiles Under Various Models

Steven T. Garren

Abstract

Some authors define the first quartile of a data set to be the median of the ordered data *strictly* to the left of the overall median. Other authors define the first quartile to be the median of the ordered data to the left of and *including* the overall median. We propose a third estimator (of the population first quartile) which has minimum mean squared error under normality, among weighted averages of two order statistics, and robustness against nonnormal distributions is examined. Likewise, estimators of the third quartile may be defined. The three estimators of first and third quartiles are compared in terms of mean squared error under normality, a t -distribution with 3 degrees of freedom, a uniform distribution, and an exponential distribution. The preferred estimator depends on the distribution and the sample size modulus 4, for sample sizes no larger than 30.

Keywords: Exponential distribution, Normal distribution, t -distribution, Uniform distribution.

2000 Mathematics Subject Classification: 62F10.

1 Introduction

Let x_i be the i th order statistic, for $i = 1, \dots, n$. If $n = 5$, then the sample median is x_3 . However, multiple definitions exist for the first and third quartiles of a sample. When the population is symmetric, we restrict attention to the first quartile without loss of generality. Some authors (c.f., Devore 2004, p. 41) define the sample first quartile, or sample 25th

percentile, to be x_2 when $n = 5$. This definition is reasonable in the sense that the five ordered observations divide the data set into four non-overlapping regions. However, if the five observations were sampled independently from a normal distribution, then x_2 is positively biased. Therefore, defining the sample first quartile to be $(x_1 + x_2)/2$ (c.f., Moore and McCabe 2003, p. 42) reduces the magnitude of the bias somewhat, but only for some sample sizes including $n = 5$.

Herein, we discuss estimating the population first quartile, θ , for general $n \geq 3$. A new estimator, which is a function of two of the order statistics, is proposed, such that this estimator has smallest mean squared error (MSE) under normality among all weighted averages of these two order statistics. The three estimators are compared under squared error loss for the normal distribution, a t -distribution with 3 degrees of freedom, a uniform(0, 1) distribution, and an exponential distribution with mean one. Since the exponential distribution is not symmetric, estimating the third quartile is also discussed. However, since the normal, t , and uniform distributions are symmetric, then results pertaining to the first quartile may also be applied to the third quartile. These four distributions were selected due to the varying heaviness of their tails.

The results presented herein are based mainly on simulations, which may be performed for estimating any percentile. To ensure a high level of precision, 5 million simulations were used to produce each number in the tables, except for the numbers which can be calculated analytically.

2 Defining Sample Quartiles

Sample median typically is defined to be the middle of the ordered observations when n is odd, and

Steven T. Garren is Associate Professor, Department of Mathematics and Statistics, James Madison University, Harrisonburg, Virginia 22807, USA. Email: garrenst@jmu.edu

the average of the two middle observations when n is even. However, two different definitions of the sample first quartile are commonly used by different authors. One definition (c.f., Moore and McCabe 2003, p. 42) is the following: Let $\hat{\theta}_1$ be the median of the ordered observations *strictly* to the left of the position of the overall median. The other definition (c.f., Devore 2004, p. 41) is the following: Let $\hat{\theta}_2$ be the median of the ordered observations to the left of and *including* the position of the overall median. When n is even, $\hat{\theta}_1 = \hat{\theta}_2$. However, when n is odd, then $\hat{\theta}_1 < \hat{\theta}_2$ *a.s.* for absolutely continuous distributions. An estimator $\hat{\theta}_3$, with a reduced MSE under normality, will be proposed in Section 4. Since $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$ are asymptotically unbiased as $n \rightarrow \infty$ under general regularity conditions (c.f., Bain and Engelhardt 1992, p. 244), we restrict attention to small sample sizes.

3 Mean Order Statistics Under Normality

The mean order statistics under a standard normal distribution for some sample sizes n are given in Table 1. Since the normal distribution is symmetric, only the smallest half of the order statistics are listed in the tables (without the negative sign, to save space). These mean order statistics are estimated based on 5,000,000 simulations. The standard errors, not shown herein, for the results in Table 1 are all smaller than 0.00035.

4 An Example: Determining $\hat{\theta}_3$

The following example shows how to compute the proposed estimator, $\hat{\theta}_3$ (with reduced MSE under normality), based on $n = 17$ observations. When sampling $n = 17$ observations from a $N(\mu, \sigma^2)$ distribution, the population first quartile is $\theta = \mu - 0.67449\sigma$. From Table 1, under a $N(0, 1)$ model, the two mean order statistics which flank -0.67449 are the means of x_4 and x_5 . When estimating θ based on any $N(\mu, \sigma^2)$ distribution, we consider all estimators of the form $wx_4 + (1 - w)x_5$, where w is a weight between 0 and 1. The MSE is mini-

mized when $w = [E(x_5 - x_4)(x_5 - \theta)]/E(x_5 - x_4)^2$. The weight w can be precisely estimated using the simulated mean of $(x_5 - x_4)(x_5 - \theta)$ and the simulated mean of $(x_5 - x_4)^2$. Hence, when $n = 17$, the proposed estimator $\hat{\theta}_3$ is $0.308x_4 + 0.692x_5$. For comparison, note that $\hat{\theta}_1 = (x_4 + x_5)/2$ and $\hat{\theta}_2 = x_5$.

In general, the appropriate weights for minimizing MSE when estimating θ under normality for $n = 3, \dots, 30$ are listed in Table 2, and are grouped according to q_n , where

$$q_n = n \text{ modulus } 4.$$

This grouping is reasonable since a seasonal trend of the weights exists, where the trend is repeated as the sample size increases by four. This repeating trend of four units is not too surprising, since the sample median has two definitions (one for when n is even, and one for when n is odd), and similarly a sample quartile has two definitions, depending on n . Seasonal trends will also be noted when determining which estimator has smallest MSE, in Section 5.

The weights from Table 2 are graphed in Figure 3. To analyze robustness, these weights also will be used when data are generated from t_3 , uniform, and exponential distributions.

The two flanking order statistics used to compute $\hat{\theta}_3$ are the same ones used for computing $\hat{\theta}_1$ and $\hat{\theta}_2$, for $n = 3, \dots, 30$. The weights given to the flanking marker further from the sample median for computing $(\hat{\theta}_1, \hat{\theta}_2)$ are (0.5, 0.5) for $q_n = 0$; (0.5, 0) for $q_n = 1$; (0, 0) for $q_n = 2$; and (1, 0.5) for $q_n = 3$.

5 Comparing Estimators

The population first quartiles for the $N(0, 1)$, t_3 , uniform(0, 1), and exponential(1) distributions are -0.67449 , -0.76489 , 0.25 , and $-\log(0.75) = 0.28768$, respectively. The population third quartile for the exponential(1) distribution is $-\log(0.25) = 1.38629$.

The estimators $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$ are compared according to MSE, which is estimated from 5,000,000 simulations. We skip determining MSE for $n = 2$, since one cannot reasonably estimate the first quartile under such a small sample. As sample sizes get

larger than 30, the differences among the three estimators diminish. The ratios of $n \times \text{MSE}$ to variance of the original distribution are shown in Tables 4 and 5, for sample sizes 3 to 30, and results discussed below are based on these sample sizes.

5.1 Normal Data

For $N(\mu, \sigma^2)$ data, when comparing $\hat{\theta}_1$ with $\hat{\theta}_2$, the preferred estimator is $\hat{\theta}_1$ when $q_n = 1$ but is $\hat{\theta}_2$ when $q_n = 3$, according to Table 4. Note that when $q_n = 1$, the estimator $\hat{\theta}_1$ is based on averaging two order statistics but $\hat{\theta}_2$ is based on just one order statistic. Likewise, when $q_n = 3$, the estimator $\hat{\theta}_2$ is based on averaging two order statistics but $\hat{\theta}_1$ is based on just one order statistic. If q_n is 0 or 2, then n is even and $\hat{\theta}_1 = \hat{\theta}_2$. Hence, when $\hat{\theta}_1$ and $\hat{\theta}_2$ differ, the preferred estimator between those two is the one based on two order statistics rather than just one order statistic. When comparing all three estimators, $\hat{\theta}_3$ has the smallest MSE, since it was defined to minimize MSE among all weighted averages of the two flanking order statistics.

5.2 t_3 -distributed Data

Next, we consider a sample from a t -distribution with three degrees of freedom. The variance of this heavy-tailed distribution is 3, and the results in Table 4 consist of $n \times \text{MSE}/3$. When comparing the two estimators, $\hat{\theta}_1$ and $\hat{\theta}_2$, the smaller MSE for all odd values of n from 3 to 30 is produced by $\hat{\theta}_2$. Of course $\hat{\theta}_1 = \hat{\theta}_2$ when n is even. For all values of n from 3 to 30, the estimator $\hat{\theta}_2$ has smaller MSE than $\hat{\theta}_3$ for $q_n = 1, 2, 3$, but larger MSE for $q_n = 0$.

5.3 Uniform Data

Now, we consider data from a uniform distribution. Since the variance of a uniform(0,1) distribution is $1/12$, in Table 4 are listed $12n \times \text{MSE}$. Since the order statistics of a uniform(0,1) random variable are beta-distributed, the MSE of $\hat{\theta}_1$ and $\hat{\theta}_2$ can be computed analytically when these estimators are based on just one order statistic. For example, when $q_n = 1$, $\text{MSE}(\hat{\theta}_2) = (3n + 1)/[16(n + 1)(n + 2)]$. Also, when $q_n = 2$, $\text{MSE}(\hat{\theta}_1) = \text{MSE}(\hat{\theta}_2) = 3/[16(n + 1)]$.

Furthermore, when $q_n = 3$, $\text{MSE}(\hat{\theta}_1) = 3/[16(n + 2)]$. Moreover, $\hat{\theta}_1$ is unbiased when n is odd.

The estimator $\hat{\theta}_1$ always has smaller MSE than $\hat{\theta}_2$ for odd n from 3 to 30. Furthermore, $\hat{\theta}_3$ has smaller MSE than $\hat{\theta}_1$ only when $q_n = 2$.

5.4 Exponential Data: Estimating the First Quartile

Next, we consider estimating the population first quartile when data are sampled from an exponential distribution. Table 5 shows that the estimator $\hat{\theta}_1$ always has smaller MSE than $\hat{\theta}_2$ for odd n from 3 to 30, similar to the uniform case in Section 5.3. This similarity is not too surprising since both the exponential and uniform distributions have truncated left tails. Also, as in the uniform case, the estimator $\hat{\theta}_3$ has smaller MSE than $\hat{\theta}_1$ only when $q_n = 2$.

5.5 Exponential Data: Estimating the Third Quartile

Finally, we consider estimating the population third quartile when data are sampled from an exponential distribution. Hence, we are focusing on a moderately heavy right tail. When estimating the third quartile, we redefine $\hat{\theta}_1$ to be the median of the ordered observations *strictly* to the right of the position of the overall median, and $\hat{\theta}_2$ to be the median of the ordered observations to the right of and *including* the position of the overall median. Hence, $\hat{\theta}_1 > \hat{\theta}_2$ *a.s.* for odd n , and $\hat{\theta}_1 = \hat{\theta}_2$ for even n .

These results are similar to those in Section 5.2 regarding the heavy-tailed t_3 -distribution. The estimator $\hat{\theta}_2$ has smaller MSE than $\hat{\theta}_1$ for odd n between 3 and 30. Furthermore, $\hat{\theta}_2$ has smaller MSE than $\hat{\theta}_3$ except when $q_n = 0$.

6 Summary

Which of the three estimators proposed herein is most preferable depends on the distribution, so constructing quantile-quantile plots may be beneficial. The overall results are summarized in Table 6, in terms of which estimator has smaller MSE among

$\hat{\theta}_1$ and $\hat{\theta}_2$, and which estimator has smallest MSE among $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$.

6.1 Comparing $\hat{\theta}_1$ With $\hat{\theta}_2$

Commonly-used estimators of a population quartile are $\hat{\theta}_1$ and $\hat{\theta}_2$, and are also employed for constructing boxplots. First, we summarize comparisons of $\hat{\theta}_1$ with $\hat{\theta}_2$. When the distribution is normal, the estimator based on two order statistics is preferable to the estimator based on just one order statistic.

REMARK 6.1 *Under normality for $3 \leq n \leq 30$, the estimator $\hat{\theta}_1$ is preferable to $\hat{\theta}_2$ when $q_n = 1$, but $\hat{\theta}_2$ is preferable to $\hat{\theta}_1$ when $q_n = 3$.*

We now consider the other distributions.

REMARK 6.2 *In our examples for odd $n \leq 30$ where the tail is truncated (i.e., when estimating the population first quartile of a uniform or an exponential distribution), the preferred estimator is $\hat{\theta}_1$, which is further away from the sample median than θ_2 is. However, in our examples where the tail is somewhat heavy (i.e., when estimating the population third quartile of an exponential or a t_3 -distribution), the preferred estimator is $\hat{\theta}_2$, which is closer to the sample median than $\hat{\theta}_1$ is.*

6.2 Comparing $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$

We now summarize our results when determining which estimator among $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$ has the smallest MSE. Under normality, $\hat{\theta}_3$, by its definition, has the smallest MSE. Consider the four possible values of q_n .

Case where n is odd: When n is odd, q_n is 1 or 3. When the tail is truncated (i.e., when estimating the population first quartile of a uniform or an exponential distribution), $\hat{\theta}_1$ has the smallest MSE. However, when the tail is heavy (i.e., when estimating the population third quartile of an exponential or a t_3 distribution), $\hat{\theta}_2$ has the smallest MSE.

Case where $q_n = 0$: In our examples where the tail is truncated (i.e., when estimating the population

first quartile of a uniform or an exponential distribution), the estimator $\hat{\theta}_1$ ($= \hat{\theta}_2$) has the smallest MSE. However, in our examples where the tail is heavy (i.e., when estimating the population third quartile of an exponential or a t_3 -distribution), the estimator $\hat{\theta}_3$ has the smallest MSE.

Case where $q_n = 2$: When the tail is truncated (i.e., when estimating the population first quartile of a uniform or an exponential distribution), $\hat{\theta}_3$ has the smallest MSE. When the tail is heavy (i.e., when estimating the population third quartile of an exponential or a t_3 distribution), $\hat{\theta}_1$ ($= \hat{\theta}_2$) has the smallest MSE.

References

- Bain, L. J., and Engelhardt, M. (1992), *Introduction to Probability and Mathematical Statistics, 2nd ed.*, Pacific Grove, CA, Brooks/Cole.
- Devore, J. (2004), *Probability and Statistics for Engineering and the Sciences, 6th ed.*, Belmont, CA: Brooks/Cole.
- Moore, D. S., and McCabe, G. P. (2003), *Introduction to the Practice of Statistics, 4th ed.*, New York: W. H. Freeman.

Table 1: *Normal Distribution*: Negative Mean Order Statistics, Estimated From 5,000,000 Simulations.

n	1	2	3	4	5	6	7	8	9	10	11	12	13
3	0.846	0											
4	1.030	0.297											
5	1.163	0.495	0										
6	1.267	0.642	0.202										
7	1.352	0.758	0.353	0									
8	1.424	0.853	0.473	0.153									
9	1.485	0.933	0.572	0.275	0								
10	1.539	1.002	0.656	0.376	0.123								
11	1.587	1.062	0.729	0.463	0.225	0							
12	1.629	1.116	0.793	0.537	0.313	0.103							
13	1.668	1.164	0.850	0.603	0.389	0.191	0						
14	1.704	1.208	0.901	0.662	0.456	0.268	0.088						
15	1.736	1.248	0.948	0.715	0.516	0.336	0.166	0					
16	1.766	1.285	0.990	0.764	0.571	0.397	0.234	0.078					
17	1.794	1.319	1.030	0.808	0.620	0.452	0.296	0.146	0				
18	1.820	1.351	1.066	0.848	0.665	0.502	0.352	0.208	0.069				
19	1.845	1.380	1.100	0.886	0.707	0.548	0.402	0.265	0.131	0			
20	1.868	1.408	1.131	0.921	0.746	0.591	0.449	0.316	0.188	0.062			
21	1.889	1.434	1.160	0.954	0.782	0.630	0.492	0.363	0.239	0.119	0		
22	1.910	1.458	1.188	0.985	0.816	0.667	0.532	0.406	0.287	0.171	0.057		
23	1.929	1.481	1.214	1.014	0.847	0.702	0.570	0.447	0.331	0.218	0.109	0	
24	1.948	1.503	1.239	1.041	0.877	0.734	0.604	0.485	0.371	0.263	0.156	0.052	
25	1.966	1.524	1.263	1.067	0.905	0.764	0.637	0.520	0.410	0.304	0.201	0.100	0
26	1.982	1.544	1.285	1.091	0.932	0.793	0.668	0.553	0.445	0.342	0.242	0.144	0.048

Table 2: Weights for Computing $\hat{\theta}_3$.

$q_n = 0$	n	4	8	12	16	20	24	28
	weight	0.436	0.442	0.444	0.445	0.446	0.445	0.445
$q_n = 1$	n	5	9	13	17	21	25	29
	weight	0.273	0.295	0.304	0.308	0.311	0.312	0.314
$q_n = 2$	n	6	10	14	18	22	26	30
	weight	0.136	0.160	0.170	0.176	0.180	0.182	0.184
$q_n = 3$	n	3	7	11	15	19	23	27
	weight	0.634	0.602	0.593	0.588	0.584	0.583	0.580

The *weight* is given to the order statistic which is flanking -0.67449 and is further from the sample median, and $(1 - \textit{weight})$ is given to the order statistic which is flanking -0.67449 and is closer to the sample median.

Figure 3: Graph of Table 2

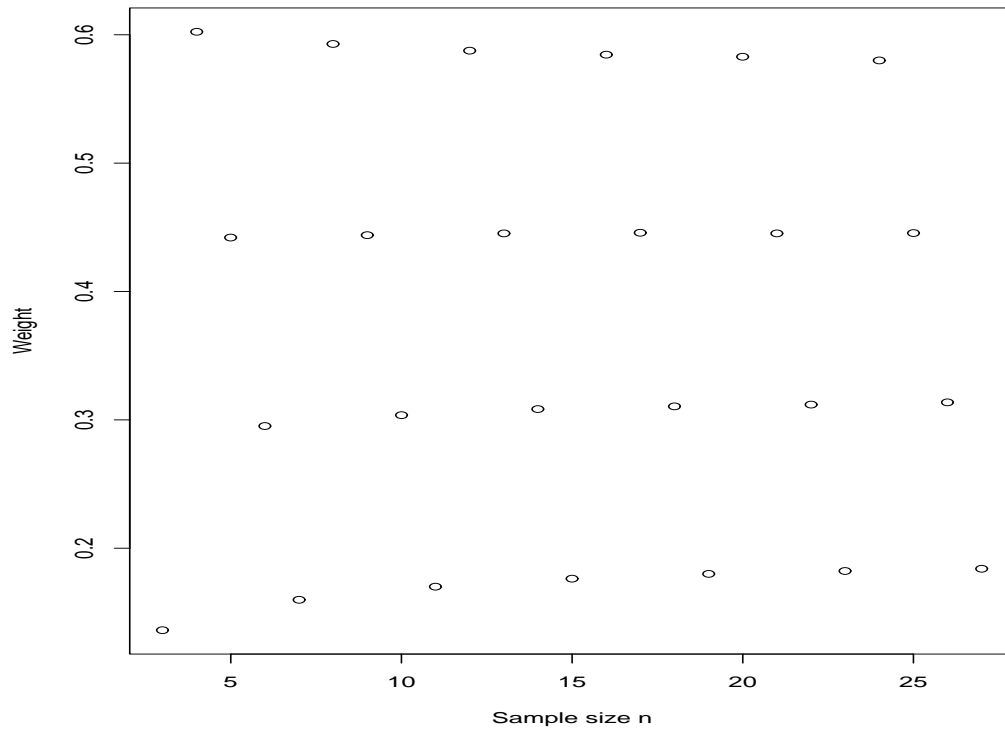


Table 4: $n \times$ Ratio of MSE of Estimators to Variance of Distribution, Based on 5,000,000 Simulations

	n	Normal			t_3			Uniform		
		$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_3$	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_3$	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_3$
$q_n = 0$	4	1.344	1.344	1.330	1.597	1.597	1.382	1.400	1.400	1.515
	8	1.534	1.534	1.527	1.011	1.011	0.966	1.732	1.732	1.801
	12	1.619	1.619	1.615	0.960	0.960	0.936	1.880	1.880	1.930
	16	1.667	1.667	1.664	0.946	0.946	0.930	1.960	1.960	1.999
	20	1.698	1.698	1.696	0.942	0.942	0.930	2.013	2.013	2.045
	24	1.721	1.721	1.719	0.939	0.939	0.929	2.048	2.048	2.075
	28	1.738	1.738	1.736	0.938	0.938	0.930	2.077	2.077	2.101
$q_n = 1$	5	1.631	1.718	1.436	2.380	0.967	1.337	1.250	2.321	1.559
	9	1.674	1.769	1.586	1.234	0.956	1.020	1.596	2.332	1.778
	13	1.711	1.790	1.654	1.095	0.953	0.981	1.765	2.321	1.896
	17	1.738	1.803	1.696	1.044	0.950	0.967	1.864	2.311	1.965
	21	1.756	1.812	1.723	1.018	0.949	0.960	1.931	2.303	2.014
	25	1.766	1.814	1.739	1.002	0.947	0.955	1.979	2.297	2.049
	29	1.778	1.820	1.755	0.992	0.946	0.953	2.012	2.292	2.073
$q_n = 2$	6	1.684	1.684	1.610	1.113	1.113	1.334	1.929	1.929	1.626
	10	1.751	1.751	1.696	1.040	1.040	1.103	2.045	2.045	1.802
	14	1.776	1.776	1.733	1.012	1.012	1.048	2.100	2.100	1.905
	18	1.795	1.795	1.760	0.996	0.996	1.020	2.132	2.132	1.966
	22	1.805	1.805	1.775	0.986	0.986	1.005	2.152	2.152	2.012
	26	1.811	1.811	1.785	0.979	0.979	0.994	2.167	2.167	2.044
	30	1.815	1.815	1.793	0.973	0.973	0.985	2.177	2.177	2.070
$q_n = 3$	3	1.767	1.359	1.295	2.842	1.139	1.413	1.350	1.800	1.502
	7	1.844	1.543	1.521	1.399	0.888	0.938	1.750	2.042	1.888
	11	1.853	1.627	1.615	1.199	0.886	0.912	1.904	2.116	2.013
	15	1.854	1.674	1.665	1.121	0.894	0.910	1.985	2.152	2.075
	19	1.853	1.703	1.697	1.079	0.900	0.912	2.036	2.170	2.109
	23	1.853	1.726	1.720	1.053	0.904	0.914	2.070	2.186	2.135
	27	1.854	1.741	1.737	1.035	0.908	0.917	2.095	2.196	2.153

Table 5: $n \times$ Ratio of MSE of Estimators to Variance of the Exponential Distribution, Based on 5,000,000 Simulations

	n	1st quartile			3rd quartile		
		$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_3$	$\hat{\theta}_1$	$\hat{\theta}_2$	$\hat{\theta}_3$
$q_n = 0$	4	0.428	0.428	0.482	2.850	2.850	2.526
	8	0.376	0.376	0.400	2.770	2.770	2.630
	12	0.361	0.361	0.376	2.810	2.810	2.721
	16	0.353	0.353	0.364	2.841	2.841	2.777
	20	0.349	0.349	0.358	2.863	2.863	2.812
	24	0.346	0.346	0.354	2.881	2.881	2.839
	28	0.345	0.345	0.351	2.896	2.896	2.860
	$q_n = 1$	5	0.285	0.645	0.409	4.354	2.367
9		0.302	0.511	0.365	3.503	2.639	2.877
13		0.310	0.457	0.352	3.299	2.747	2.910
17		0.314	0.428	0.346	3.210	2.807	2.930
21		0.317	0.410	0.343	3.158	2.840	2.939
25		0.320	0.398	0.341	3.130	2.867	2.950
29		0.322	0.389	0.340	3.108	2.885	2.956
$q_n = 2$		6	0.444	0.444	0.362	2.974	2.974
	10	0.403	0.403	0.342	3.016	3.016	3.212
	14	0.384	0.384	0.336	3.022	3.022	3.163
	18	0.373	0.373	0.334	3.021	3.021	3.132
	22	0.366	0.366	0.333	3.016	3.016	3.107
	26	0.361	0.361	0.333	3.016	3.016	3.092
	30	0.358	0.358	0.333	3.012	3.012	3.078
	$q_n = 3$	3	0.340	0.783	0.590	4.681	1.841
7		0.341	0.512	0.454	3.881	2.283	2.468
11		0.339	0.445	0.412	3.585	2.484	2.596
15		0.338	0.414	0.392	3.433	2.595	2.674
19		0.337	0.397	0.380	3.344	2.669	2.730
23		0.336	0.385	0.372	3.287	2.720	2.770
27		0.336	0.378	0.367	3.241	2.754	2.795

Table 6: Best Estimator in Terms of MSE, When Comparing $\hat{\theta}_1$ and $\hat{\theta}_2$, and When Comparing $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$, for $n = 3, \dots, 30$. (Only the subscript on $\hat{\theta}_j$ is listed for $j = 1, 2, 3$. The distributions are listed from lightest to heaviest tail.)

Distribution	Comparing $\hat{\theta}_1$ and $\hat{\theta}_2$		Comparing $\hat{\theta}_1$, $\hat{\theta}_2$, and $\hat{\theta}_3$			
	$q_n = 1$	$q_n = 3$	$q_n = 0$	$q_n = 1$	$q_n = 2$	$q_n = 3$
Uniform	1	1	1&2	1	3	1
Exponential (1st quartile)	1	1	1&2	1	3	1
Normal	1	2	3	3	3	3
Exponential (3rd quartile)	2	2	3	2	1&2	2
t_3	2	2	3	2	1&2	2