

Performance of Some Robust Estimators for the Weibull Distribution

N. GÜNDÜZ*, E. BAŞAR AND C. AYDIN

Gazi University, Faculty of Science, Statistics Department, Ankara, Turkey

In this study, we consider some of univariate quantile-based robust estimators. We focus on the estimators such as median, interquartile range, quartile and octile skewness for the Weibull distribution which is one of the most widely applied probability function because of its versatility and relative simplicity. It is important to use robust estimators as a measure of distribution properties for analyzing data in the case of contamination with outliers. For small data sets, it is reported that by introducing kernel estimation for smoothing empirical distribution function, a reduction in mean square error of estimator is achieved by Fernholz (1997) and Hubert et al. (2013). In kernel estimation, it is well known that bandwidth selection is more important than selection of kernel density since bandwidth controls the smoothness of the estimated distribution function. Using simulation studies, we examine some quantile-based estimators for the Weibull distribution with various sample size. The performance of estimators is measured by mean squared error under Different outlier contaminated data. We applied this idea in the case of real data.

DOI: [10.12693/APhysPolA.128.B-203](https://doi.org/10.12693/APhysPolA.128.B-203)

PACS: 02.70.Rr, 02.50.-r

1. Introduction

In this study, quantile-based robust estimators are investigated for the Weibull distribution. It is known that the Weibull distribution has wide spread application in medicine, biology engineering, and as a probability model it is very common in modeling the problems of the area of survival and reliability analysis. For this purpose, kernel estimation is applied to random samples which are taken from the Weibull distribution for varied parameters in order to obtain smoothed distribution function. After that, the mean square error (MSE) of robust estimators and variances are obtained.

As a quantile-based robust statistics, the median (med), interquartile range (IQR), quartile skewness (QS), octile skewness (OS) is considered. We examine the reduction in MSE for estimators. In the case of contamination with Different proportion, the behavior of related statistics is investigated by a simulation study for random sample of Different size.

The first proposals about kernel smoothing for distribution functions estimates has been made by Nadaraya [1] and Azzalini [2]. In the following years, for small data sets and in the case of outlier, Fernholz [3] proved that the MSE of the estimators obtained from the smoothed distribution function is less than the MSE of estimators obtained from the empirical distribution function.

For estimators which have discontinuous influence function, it is declared that kernel smoothing is especially useful [4]. Hubert et al. [5] stated the results of a simulation study conducted for the kernel smoothing to random

samples that are taken from a gamma distribution with various parameters including the cases of contamination. Additionally, they reported that a considerable decrease in the MSE of the quantile based estimators has occurred.

Quantile function and the definition of the estimators considered took place in the second section of the study. In Sect. 3, the smoothing procedure of the empirical distribution function is explained. In this study, for the Weibull distribution a simulation study is constructed to determine bandwidth by minimizing integrated mean square error ($IMSE$) of smoothed distribution function. It is explained in Sect. 4. In Sect. 4.1 the algorithm of simulation study is given. In Sect. 5 the results of estimators such as MSE , variance, and bias of the related statistics are tabulated for non-contamination and contamination cases separately. In case of real data bandwidth selection procedure is explained and illustrated to lifetime data in Sect. 6, and finally conclusion takes place in the last section.

The written algorithm which is used for the Weibull distribution in this study has a general structure, so it can be easily applied to other distributions. We want to use the benefits of the algorithm in order to make similar studies in general for any distribution.

2. Quantile function

Let $\{x_1, x_2, \dots, x_n\}$ be an independent and identically distributed random sample drawn from an absolutely continuous distribution function $F(x)$ with probability density function $f(x)$. We use the conventionally-established quantile function

$$Q(p) = \inf \{x : p \leq F(x)\}, \quad 0 \leq p \leq 1. \quad (1)$$

Every member of the real line is connected with one quantile function value, since quantile function is the inverse

*corresponding author; e-mail: ngunduz@gazi.edu.tr

of the distribution function. Empirical distribution function is

$$F_n(x) = \widehat{F}(x) = \frac{1}{n} \sum_{i=1}^n I_{(-\infty, x]}(x_i),$$

$$I_{(-\infty, x]}(u) = \begin{cases} 1, & u \in (-\infty, x], u \leq x, \\ 0, & u \notin (-\infty, x], u > x. \end{cases} \quad (2)$$

Accordingly, empirical quantile function is defined as:

$$Q_n(p) = \inf \{x : p \leq F_n(x)\}, \quad 0 \leq p \leq 1. \quad (3)$$

When kernel smoothing is applied, robust quantile-based estimators with high breakdown point can be achieved more accurately from smoothed quantile function, since we can create much more quantile values that is provided by the order statistics. In this study, some robust quantile-based estimators such as, median as a location estimator, *IQR* as a scale estimator, *QS* and *OS* as a measure of skewness are investigated for the Weibull distribution [5, 6].

3. Kernel smoothing

An empirical distribution function estimates the distribution function of a random variable by assigning equal probability to each observation in a sample. It is discontinuous at many points. Kernel smoothing is applied to achieve a smoother empirical distribution function, so we have a continuous estimate of distribution function, which makes it possible to estimate the density of a random variable based on an observed sample.

Kernel-based estimator of a distribution function is given as follows [1]:

$$\tilde{F}_{n,h}(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right), \quad (4)$$

where $K(\cdot)$ is distribution function of the Epanechnikov kernel and h is smoothing parameter called the bandwidth [7]. If a random variable X has a distribution function $F(x)$, that is Differentiable twice and has continuous second derivative, while $n \rightarrow \infty$, $h \rightarrow 0$, $nh \rightarrow \infty$, the expected value and the variance of the smoothed distribution function estimator $\tilde{F}_{n,h}(x)$ is given as follows:

$$E\left(\tilde{F}_{n,h}(x)\right) = F(x) + \frac{1}{2}h^2 f'(x) \mu_2(k) + O(h^2), \quad (5)$$

$$V\left(\tilde{F}_{n,h}(x)\right) = \frac{F(x)(1 - F(x))}{n} - \frac{2hf(x)c}{n} + O\left(\frac{h}{n}\right), \quad (6)$$

where $\mu_2(k) = \int_{-\infty}^{\infty} t^2 k(t) dt$ and $c = \int_{-\infty}^{\infty} tk(t) K(t) dt$.

For the Epanechnikov kernel we have $\mu_2(k) = 1$, the constant c is 0.2875 [2].

4. Bandwidth determination using simulation

It is known that the choice of the kernel function $K(\cdot)$ is less important than the choice of bandwidth in kernel estimation [8]. It is common practice to use *MSE*

for the measure of performance of estimator. *IMSE* of smoothed distribution function estimate $\tilde{F}_{n,h}(x)$ is given below

$$IMSE\left(\tilde{F}_{n,h}(x)\right) = E \int_{-\infty}^{\infty} \left[\tilde{F}_{n,h}(x) - F(x)\right]^2 dx. \quad (7)$$

We designed a simulation study, and obtained the sampling distribution of the estimator $\tilde{F}_{n,h}(x)$. Then we looked for bandwidth which minimizes the *IMSE*. Specifically, we drew 5000 random samples of the size 40 for each selected h value. We calculated *IMSE* over that sampling distribution and plotted it against h . For *Weibull*(1.5, 1) and *Weibull*(4, 1), Fig. 1 gives plots of *IMSE* estimate versus h . *IMSE* is minimum when $h = 0.6$ and $h = 0.25$ for *Weibull*(1.5, 1) and *Weibull*(4, 1), respectively.

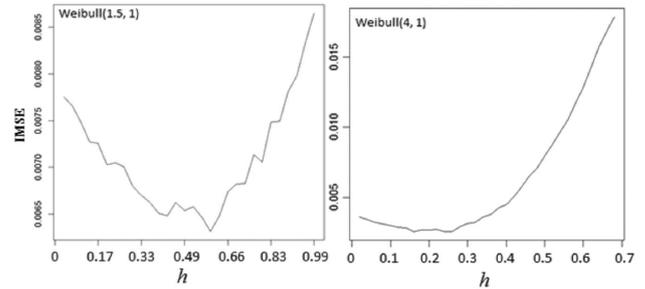


Fig. 1. *IMSE* estimates versus h values.

These particular bandwidth values are used for the simulation study to investigate the reduction in *MSE* and variance of quantile-based estimators.

4.1. Algorithm

<p>Data producing</p> <p>Random samples are drawn and summary statistics including order statistics are obtained. The working interval that covers the range is determined. Working interval is partitioned such that satisfying the following requirement:</p> <p><i>length of working interval</i> \leq <i>partition number</i> * <i>step</i>.</p> <p>This partitioning makes possible to combine results of simulation study. Distribution function values are calculated for the selected grid of working interval.</p>
<p>↓</p>
<p>Kernel smoothing</p> <p>For a drawn sample $\{x_1, x_2, \dots, x_n\}$, smoothed distribution function estimate $\tilde{F}_{n,h}(x)$ is obtained by</p> $\tilde{F}_{n,h}(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$ <p>$x = \text{initial value} + k * \text{step};$ $K = \text{partitioning number}$ $k = 0, 1, 2, \dots, K,$ where x_i is a value of random sample and x is member of working interval.</p>
<p>↓</p>

Bandwidth selection

Draw 5000 random samples of size 40 for every selected h value in the interval $0 \leq h \leq 2$. $h_k = 0 + l * (0.01)$, $l = 0, 1, 2, \dots, 200$. We obtained the sampling distribution of $\tilde{F}_{n,h}(x)$ for each selected h value. We estimate the *IMSE* over the sampling distribution of $\tilde{F}_{n,h}(x)$, in discrete terms, as follows:

IMSE = Diff square + Bias square.

$$Diff\ square = \frac{1}{r} \sum_{i=1}^r \sum_{k=1}^K \left(\tilde{F}_{(n,h),i}(x_k) - \bar{F}_{n,h}(x_k) \right)^2 *$$

step,
 $r = replication\ number$
 $x_k = initial\ point + k * step$
 $k = 0, 1, 2, \dots, K$.

$$Bias\ square = \sum_{k=1}^K \left(\bar{F}_{n,h}(x_k) - \bar{F}_n(x_k) \right)^2 * step,$$

$x_k = x + k * step;$
 $k = 0, 1, 2, \dots, K$.
 $\bar{F}_{n,h}(x_k)$ is the average of repeated smoothed distribution function and $\bar{F}_n(x_k)$ is the average of repeated empirical distribution function. $\bar{F}_n(x_k)$ is used instead of $F(x)$.

↓

Smoothed estimators

As the smoothed distribution function value for each sample is known, by controlling, expression $\left| \tilde{F}_{(n,h),i}(x_k) - \frac{1}{2} \right| \leq \varepsilon$ (ε is chosen to be 1/1000), we can state the smoothed median for i -th sample, as $m\bar{e}d(i) = x_k$. In similar way, by controlling the expression $\left| \bar{F}_{(n,h),i}(x_l) - \frac{1}{2} \right| \leq \varepsilon$ (e.g. $\varepsilon = 1/1000$), we determine the average of smoothed median as $m\bar{e}d = x_l$. Then, we have the *MSE* for smoothed median, $m\bar{e}d$, $MSE(m\bar{e}d) = \frac{1}{r} \sum_{i=1}^r (m\bar{e}d(i) - m\bar{e}d)^2 + (m\bar{e}d - med(Weibull))^2$.

The same approach is used for smoothed quantile estimates of IQR_n , QS_n and OS_n .

5. Simulation results

Simulation results about quantile-based estimators are obtained by applying in each case the given algorithm for two Different Weibull distributions with parameters $Weibull(1.5, 1)$ and $Weibull(4, 1)$.

We specified $h = 0.6$ for $Weibull(1.5, 1)$ and $h = 0.25$ for $Weibull(4, 1)$ which are reported in the first part of simulation study. We applied the simulation algorithm in case of 0%, 5%, and 10% contamination. For contamination cases, data were drawn from normal distribution $N(\mu = 20, \sigma^2 = 1)$. We chose normal distribution to represent contamination structure and we believe that it is good enough to keep track of the performance of estimators for these contamination cases. We draw 1000 samples for various sizes for each estimator. Here we only report simulation results of sample size 40 in Tables I-IV.

TABLE I

MSE, variance and bias of med_n and $\tilde{Q}_n(0.5)$ under $Weibull(1.5, 1)$ and $Weibull(4, 1)$ distributions for 0%, 5% and 10% contamination, respectively.

Cont.*	Est.**	MSE	Variance	Bias
<i>Weibull(1.5, 1)</i>				
no	med_n	0.01362	0.01356	0.00723
	$\tilde{Q}_n(0.5)$	0.01081	0.00985	0.03103
5%	med_n	0.01775	0.01561	0.04625
	$\tilde{Q}_n(0.5)$	0.01691	0.01170	0.07216
10%	med_n	0.02715	0.01829	0.09408
	$\tilde{Q}_n(0.5)$	0.02999	0.01450	0.12445
<i>Weibull(4, 1)</i>				
no	med_n	0.00256	0.00256	0.00070
	$\tilde{Q}_n(0.5)$	0.00200	0.00200	0.00073
5%	med_n	0.00307	0.00277	0.01723
	$\tilde{Q}_n(0.5)$	0.00256	0.00221	0.01861
10%	med_n	0.00441	0.00306	0.03684
	$\tilde{Q}_n(0.5)$	0.00421	0.00255	0.04074

*contamination, **estimator

Table I shows 20% reduction in *MSE* of smoothed median for uncontaminated $Weibull(1.5, 1)$ i.e. $(1 - (MSE(\tilde{Q}(0.5))/MSE(med_n))) = 1 - 0.01081/0.01362 = 1 - 0.79 = 0.21$. For $Weibull(4, 1)$ this reduction is also 20%. 5% contamination, *MSE* of smoothed median has decreased 5% and 16% for $Weibull(1.5, 1)$ and $Weibull(4, 1)$, respectively. However, we noted an increase of 10% for $Weibull(1.5, 1)$ under 10% contamination and 5% decrease for $Weibull(4, 1)$. For $Weibull(4, 1)$ we see that consistently similar bias occurs across all scenarios.

TABLE II

MSE, variance and bias of IQR_n and the $I\tilde{Q}R_n$ under $Weibull(1.5, 1)$ and $Weibull(4, 1)$ distributions for 0%, 5% and 10% contamination, respectively.

Cont.	Est.	MSE	Variance	Bias
<i>Weibull(1.5, 1)</i>				
no	IQR_n	0.02380	0.02368	0.01100
	$I\tilde{Q}R_n$	0.01880	0.01276	0.07772
5%	IQR_n	0.04900	0.03716	0.10877
	$I\tilde{Q}R_n$	0.05217	0.02366	0.16884
10%	IQR_n	0.37550	0.30405	0.26731
	$I\tilde{Q}R_n$	0.14045	0.04940	0.30175
<i>Weibull(4, 1)</i>				
no	IQR_n	0.32374	0.00392	-0.56553
	$I\tilde{Q}R_n$	0.28681	0.00202	-0.53365
5%	IQR_n	0.29604	0.00491	-0.53956
	$I\tilde{Q}R_n$	0.26003	0.00286	-0.50713
10%	IQR_n	0.46765	0.23573	-0.48158
	$I\tilde{Q}R_n$	0.22941	0.01102	-0.46733

Results in Table II shows decrease of 21% in *MSE* of smoothed *IQR* for uncontaminated $Weibull(1.5, 1)$ i.e. $(1 - (MSE(I\tilde{Q}R)/MSE(IQR_n))) = 1 - 0.01880/$

0.02380 = 1 - 0.79 = 0.21). For *Weibull*(4,1) this reduction is 11%. However, we noted an increase of 6% for *Weibull*(1.5,1) with 5% contamination and 10% decrease for *Weibull*(4,1). 10% contamination leads to 62% and 51% decrease in *MSE* of smoothed *IQR* for *Weibull*(1.5,1) and *Weibull*(4,1), respectively. For *Weibull*(4,1) we see that consistently similar bias occurs across all scenarios.

TABLE III

MSE, variance and bias of QS_n and the $\tilde{Q}S_n$ under *Weibull*(1.5,1) and *Weibull*(4,1) distributions for 0%, 5% and 10% contamination, respectively.

Cont.	Est.	<i>MSE</i>	Variance	Bias
<i>Weibull</i> (1.5,1)				
no	QS_n	0.04149	0.04132	-0.01277
	$\tilde{Q}S_n$	0.00517	0.00394	-0.03508
5%	QS_n	0.04058	0.04001	0.02398
	$\tilde{Q}S_n$	0.00655	0.00655	0.00005
10%	QS_n	0.04734	0.04186	0.07399
	$\tilde{Q}S_n$	0.01191	0.00960	0.04803
<i>Weibull</i> (4,1)				
no	QS_n	0.04130	0.04128	-0.00439
	$\tilde{Q}S_n$	0.00423	0.00418	0.00744
5%	QS_n	0.04042	0.03980	0.02501
	$\tilde{Q}S_n$	0.00664	0.00551	0.03358
10%	QS_n	0.04798	0.04330	0.06845
	$\tilde{Q}S_n$	0.01442	0.00886	0.07459

From Table III the reduction *MSE* of $\tilde{Q}S_n$ has decreased at least 70% in each case.

TABLE IV

MSE, variance and bias of OS_n and the $\tilde{O}S_n$ under *Weibull*(1.5,1) and *Weibull*(4,1) distributions for 0%, 5% and 10% contamination, respectively.

Cont.	Est.	<i>MSE</i>	Variance	Bias
<i>Weibull</i> (1.5,1)				
no	OS_n	0.03014	0.02738	-0.05255
	$\tilde{O}S_n$	0.01185	0.00718	-0.06833
5%	OS_n	0.03907	0.03883	0.01557
	$\tilde{O}S_n$	0.01142	0.01141	0.00360
10%	OS_n	0.13353	0.08885	0.21139
	$\tilde{O}S_n$	0.01391	0.00975	0.06448
<i>Weibull</i> (4,1)				
no	OS_n	0.02916	0.02694	-0.04716
	$\tilde{O}S_n$	0.00741	0.00727	0.01190
5%	OS_n	0.04626	0.04602	0.01522
	$\tilde{O}S_n$	0.04233	0.03167	0.10324
10%	OS_n	0.24224	0.17299	0.26317
	$\tilde{O}S_n$	0.20298	0.08901	0.33760

In Table IV, when there is contamination the reduction in *MSE* of octile skewness has decreased by at most 15% for *Weibull*(4,1). In other cases, the reduction in *MSE* has decreased by at least 60%.

6. Real data example

6.1. Bandwidth selection for real data

In this part of the study, the structures which are constructed in simulation study and summarized in Sect. 4 and 5, are used for the application of real data problem. If $n \rightarrow \infty$, $h \rightarrow 0$, $nh \rightarrow \infty$, from Eqs. (5) and (6), asymptotic integrated mean square error (*AIMSE*) is defined as:

$$AIMSE(h) = \frac{\int_{-\infty}^{\infty} F(x)[1-F(x)] dx}{n} - \frac{2hc}{n} + \frac{h^4 R}{4} + O(h^4), \tag{8}$$

where R is the roughness of $f(x)$, $R = \int_{-\infty}^{\infty} (f'(x))^2 dx$.

By equating the first derivative of (8) to zero, we see that *AIMSE* is minimized at

$$h_0 = \left(\frac{2c}{R}\right)^{\frac{1}{3}} n^{-\frac{1}{3}}. \tag{9}$$

Here, most striking point is that optimal bandwidth is inverse proportional to the roughness R . As $f(x)$ is estimated by kernel estimation, also the estimate of $f''(x)$ can be achieved by kernel estimation by Epanechnikov kernel with

$$\tilde{f}''(x) = \frac{1}{h_d^3 n} \sum_{i=1}^n k''\left(\frac{x-x_i}{h_d}\right). \tag{10}$$

Since $R = -E(f''(X))$, we can estimate roughness of $f(x)$:

$$\hat{R} = -\frac{1}{n} \sum_{i=1}^n \tilde{f}''(x_i). \tag{11}$$

In general, a plug-in bandwidth determination rule which is given by Silverman [9] is $h_d = 2.34 \min\left(\hat{\sigma}_n, \frac{IQR_n}{1.349}\right) n^{-\frac{1}{5}}$. In this study, we used bandwidth determination rule as follows: $h_d = 2.34 \min\left(\hat{\sigma}_n, \frac{IQR_n}{1.349}, 2.219Q_n, 1.192S_n\right) n^{-\frac{1}{5}}$, where Q_n and S_n are robust scale estimators alternative to median absolute deviation as a scale parameter [10].

Using the bandwidth h_d we estimate the roughness from (10) and (11), then we calculate optimal bandwidth h_0 (9) in order to estimate the smoothed distribution function. Consequently, smoothed quantile-based estimators are obtained.

6.2. Application to real data

This procedure is applied to a real data collected from register of patients admitted to Başkent University Hospital between January 1, 1990 and November 30, 1992 [11]. In Ref. [11] it was reported that the lifetimes of transplanted kidneys in months and 34 failures had been observed during the study period. Here failure means that the transplanted kidney was not compatible.

Figure 2 depicts the histogram and density estimation of kidney lifetime. It appears that the empirical distribution is right skewed and there might be an outlier.

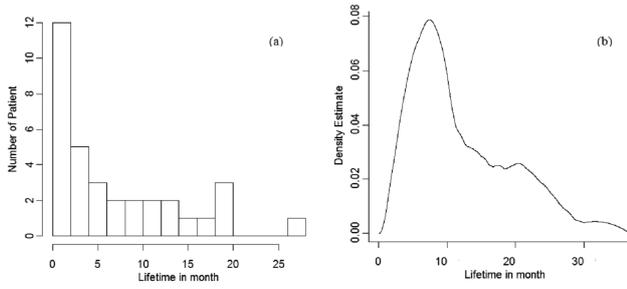


Fig. 2. (a) Histogram and (b) density estimate for data of lifetimes.

However, further investigation of the underlying boxplot fails to confirm this.

In Figure 3, both empirical and smoothed distribution functions are graphically represented together.

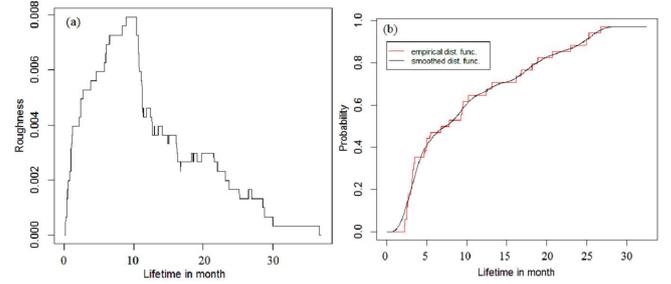


Fig. 3. (a) Estimation of roughness, (b) empirical and smoothed distribution function of data.

TABLE V

Empirical and smoothed quantile-based estimates for the lifetime data of kidney transplant patients.

estimators	$Q_n(0.125)$	$Q_n(0.25)$	$Q_n(0.50)$	$Q_n(0.75)$	$Q_n(0.875)$	IQR_n	QS_n	OS_n
empirical	0.52	1.02	4.30	11.65	17.65	0.63	0.38	0.56
smoothed	0.44	1.20	4.32	12.08	17.30	0.88	0.43	0.54

Bandwidth which is used in estimation of roughness is obtained as $h_d = 5.11$, then we got the estimation of roughness as $\hat{R} = 0.005$. After that we estimated smoothed distribution function with optimal bandwidth $h_0 = 1.49$ according to Epanechnikov kernel density. Both empirical and smoothed estimators that are obtained from empirical and smoothed distribution function are tabulated in Table V.

We have estimated asymptotic variance of first, second and third quantile estimates by obtaining estimate of density function for the real data. For $0 < p < 1$ values $SE(\widehat{Q}_p) = \sqrt{Var(Q_p)} = \sqrt{p(1-p)/(n(f(Q_p))^2)}$.

We have obtained $SE(\widehat{Q}_{0.25}) = 1.019$, $SE(\widehat{Q}_{0.50}) = 1.463$ and $SE(\widehat{Q}_{0.75}) = 3.027$. Since there is no outlier the empirical and smoothed quartile estimates are close to each other for this real data.

7. Conclusion

For small data sets, when kernel estimation is used, as it is expected, a great reduction in MSE of estimators is achieved. This has occurred in each case for median, IQR_n , QS_n , OS_n . A considerable reduction in MSE has appeared for both skewness measures QS_n , OS_n . For median, relatively small changes in MSE have occurred. In case when MSE decreases, we see that, almost for all, bias has increased strongly. So, in practice, when kernel smoothing is used, it is necessary to give more attention to control bias and to introduce bias reducing methods.

We think that we got reasonable results by applying the related kernel estimation procedure for the real data

problem. On the other hand, the programs for simulation studies and application are coded in software *R* without using any robust packages.

References

- [1] E.A. Nadaraya, *Theory Prob. Appl.* **9**, 497 (1964).
- [2] A. Azzalini, *Biometrika* **68**, 326 (1981).
- [3] L.T. Fernholz, *J. Stat. Plan Infer.* **57**, 29 (1997).
- [4] G. Brys, M. Hubert, A. Struyf, *J. Comput. Graph Stat.* **13**, 996 (2004).
- [5] M. Hubert, I. Gijbels, D. Vanpaemel, *Test* **22**, 448 (2013).
- [6] G. Brys, M. Hubert, A. Struyf, in: *Developments in Robust Statistics: Int. Conf. on Robust Statistics*, 2001, Eds. R. Dutter, P. Filzmoser, U. Gather, P.J. Rousseeuw, Vol. 114, Physika Verlag, Heidelberg 2003, p. 98.
- [7] V. Epanechnikov, *Theory Probab. Its Appl.* **14**, 53 (1969).
- [8] M.P. Wand, M.C. Jones, *Kernel Smoothing, CRC Monographs on Statistics and Applied Probability*, Chapman and Hall, London 1994.
- [9] B.W. Silverman, *Density Estimation for Statistics and Data Analysis*, Chapman and Hall, London 1986.
- [10] P.J. Rousseeuw, C. Croux, *J. Am. Stat. Assoc.* **88**, 1273 (1993).
- [11] E. Başar, Ph.D. Thesis, Science Institute of Hacettepe University, Ankara 1993.