

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/358045701>

Outlier Detection Technique for Univariate Normal Datasets

Article in *American Journal of Theoretical and Applied Statistics* · January 2022

DOI: 10.11648/j.ajtas.20221101.11

CITATIONS

0

READS

431

3 authors:



[Ooko Silas Owuor](#)

Maasai Mara University

1 PUBLICATION 0 CITATIONS

SEE PROFILE



[Benedict Troon](#)

Maasai Mara University

27 PUBLICATIONS 19 CITATIONS

SEE PROFILE



[Kevin Okumu](#)

Kenyatta University

4 PUBLICATIONS 0 CITATIONS

SEE PROFILE

Outlier Detection Technique for Univariate Normal Datasets

Ooko Silas Owuor¹, Troon John Benedict², Otieno Okumu Kevin²

¹Department of Mathematics and Physical Sciences, Maasai Mara University, Narok, Kenya

²Department of Economics, Maasai Mara University, Narok, Kenya

Email address:

sylusooko7@gmail.com (O. S. Owuor), troon@mmarau.ac.ke (T. J. Benedict), kevinotieno15@gmail.com (O. O. Kevin)

To cite this article:

Ooko Silas Owuor, Troon John Benedict, Otieno Okumu Kevin. Outlier Detection Technique for Univariate Normal Datasets. *American Journal of Theoretical and Applied Statistics*. Vol. 11, No. 1, 2022, pp. 1-12. doi: 10.11648/j.ajtas.20221101.11

Received: December 19, 2021; **Accepted:** January 8, 2022; **Published:** January 21, 2022

Abstract: This paper presents an outlier detection technique for univariate normal datasets. Outliers are observations that lie at an abnormal distance from the mean. Outlier detection is a useful technique in such areas as fraud detection, financial analysis, health monitoring and Statistical modelling. Many recent approaches detect outliers according to reasonable, pre-defined concepts of an outlier. Methods of outlier detection such as Gaussian method of outlier detection have been widely used in the detection of outliers for univariate data-sets, however, such methods use measure of central tendency and dispersion that are affected by outliers hence making the method to be less robust towards detection of outliers. The study aimed at providing an alternative method that can be used in outlier detection for univariate normal data sets by deploying the measures of variation and central tendency that are least affected by the outliers (median and the geometric measure of variation). The study formulated an outlier detection formula using median and geometric measure of variation and then applied the formulation on randomly simulated normal dataset with outliers and recorded the number of outliers detected by the method in comparison to the other two existing best methods of outlier detection. The study then compared the sensitivity of the three methods in outlier detection. The simulation was done in two different ways, the first considered the variation in mean with a constant standard deviation while the second test held the mean constant while varying the standard deviation. The formulated outlier detection technique performed the best, eliminating the most required number of outliers compared to other two Gaussian outlier detection techniques when there was variation in mean. The study also established that the formulated method of outlier detection was stricter when the standard deviation was varied but still stands out to be the best as an outlier is defined relative to the mean and not the standard deviation. The study established that the formulated method is more sensitive than the Gaussian Method of outlier detection but performed as well as the best existing outlier detection technique. In conclusion, the study established that the formulated method could be employed in outlier detections for univariate normal data-sets as it performed almost the same to the best existing method of outlier detection for univariate data-sets.

Keywords: Outlier, Anomaly, Outlier Detection, Gaussian

1. Introduction

Outlier detection; also known as anomaly detection this process is the identification of rare items [1, 2] events and observations which arise and are significantly different from the other observations in the data [3, 4]. Identification of these events (outliers) is very important given they may lead to bad data and this may lead to poor running of the experiment for they may hide very essential information about the data. If it can be determined earlier that a point is outlying then it can be worth ejecting it for the purposes of better results. Secondly, in some cases, it may not be possible to deter-

mine if an outlying point is bad data since Outliers may be due to random variation or may indicate something scientifically interesting. In any event, we typically do not want to simply delete the outlying observation. However, if the data contains significant outliers, we may need to consider the use of robust statistical techniques. [5-7] Before application of these techniques we have to determine whether the outlier is univariate or multivariate. Univariate outliers can be found when looking at a distribution of values in a single feature space. Multivariate outliers can be found in an n-dimensional space (of n-features). Looking at distributions in n-dimensional spaces can be very difficult for the human brain, that is why we need to train a model to do it for us [8,

9]. Outlier detection is an important research problem in data mining that aims to find objects that are considerably dissimilar, exceptional and inconsistent with respect to the majority data in an input database [10]. The following are the existing outlier detection techniques that the study focused on.

1.1. Gaussian Model

Estimation of mean and standard deviation is done in training stage using the maximum likelihood estimates (MLE). A wide range, nearly 100 of outlier tests has been put in place in different ways depending on the data set and the parameters like mean and variance and the expected values of the outliers. To ensure the test carried are optima or close to optima statistical discordancy tests are usually carried out in the test stage [11-13]. The usually used outlier test for normal distribution is the mean-variance and Boxplot tests. In the mean variance test for Gaussian distribution $N(\mu, \sigma)$, where the population has mean and variance σ . Outlier is considered to be a point that lie 3 or more standard deviation i.e. $> 3\sigma$ away from the mean.

Similarly, the tests can be applied to some other distribution like t-distribution and the Poisson distribution with the former featuring a latter tail and the latter a longer right tail than a normal distribution.

The box plot test also gives a profound test by deployment of 5 major attributes i.e. smallest non-outlier observation [min], lower quartile [Q1], upper quartile [Q3], medium and the largest non-outlier observation [max]. The quantity (Q3 - Q1) is called the interquartile range (IQR). This helps use clearly define the boundary beyond which the data will be considered an outlier. A point X_1 is labeled or referred to as an outlier if, $X_i > Q3 + k(IQR)$ or $X_i < Q1 - k(IQR)$.

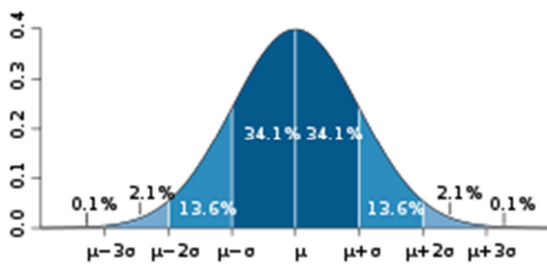


Figure 1. Outliers are points $> |\mu + 3\sigma|$, for some $k=1.5$.

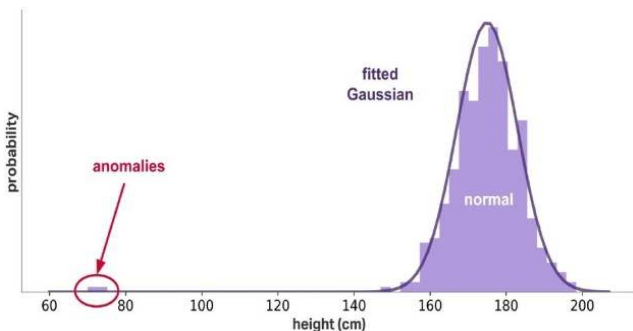


Figure 2. Outlier detection by fitted Gaussian model.

Basing our argument on low dimensional outlier detection technique, we settle with the Gaussian model which defines an outlier as a point $X > |\mu + 3\sigma|$ as the best existing detection technique as it takes into consideration both the probabilistic and normal distributions.

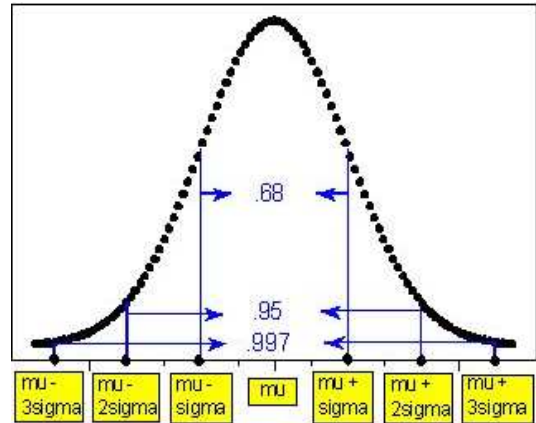


Figure 3. Outliers are points $> |\mu + 3\sigma|$.

However, this method has a number of shortcomings since by Central Limit Theorem (Which states; If you have small, independent random variables, then their sum is distributed approximately a bell curve [14-16]). By so doing, if an outlier occurs at some point away from the normal curve, then the normal curve will shift towards the outlier.

The anomaly/outlier towards the left as shown in Figure 5 will shift the normal curve towards the left as illustrated below;

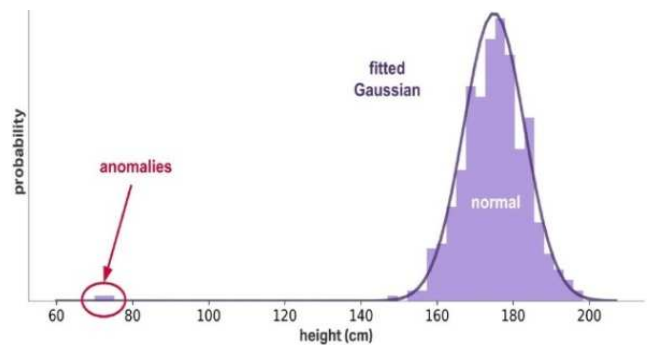


Figure 4. Positioning of an outlier towards the left of the curve.



Figure 5. Shift of the normal curve towards the outlier changing the mean but keeping the standard deviation constant.

In an event outliers occur both side of the curve then it's likely to spread the normal curve having an effect on the

standard deviation but keeping the mean constant. This is illustrated in the Figure 7 below;

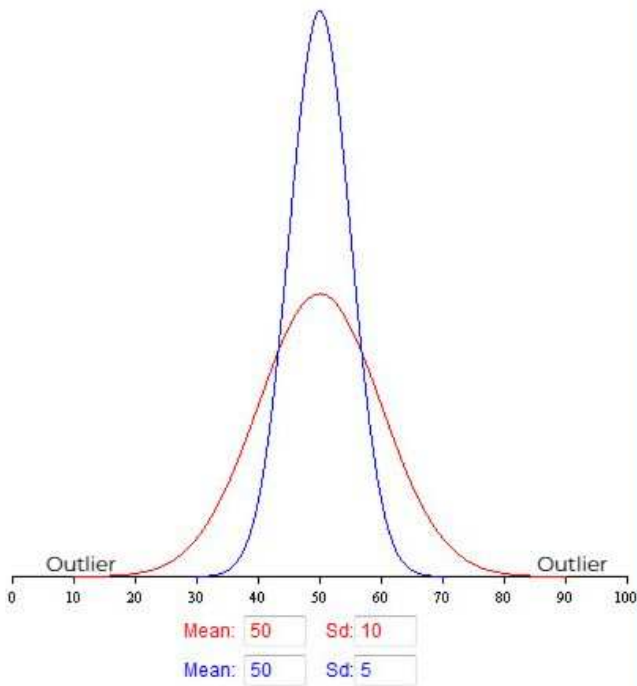


Figure 6. Stretching of the curve when outliers occur at both sides of the curve.

Since the existing Gaussian model detection technique uses parameters that are affected by the outliers as illustrated above, the study would like to come up with a technique that does not rely on these parameters μ and σ . In this regard, the study will replace the former with Median since the Median is least affected by the outliers [17, 18] and the latter with the Geometric measure because it takes into account the compounding that occurs from period to period. Because of this, investors usually consider the geometric mean a more accurate measure of returns than the arithmetic mean [19].

In so doing we will have our new detection technique define outlier as a point $X > |Med + 3g|$

Where;

Med is the median and g is the geometric measure. The main objectives of this study is Outlier detection for univariate data set using geometric technique.

1.2. Regression Model

A regression model is also used to detect the outliers. In this scenario, an outlier is considered to be an observation for which the residual is larger compared to other observations in the data-set. Such observations are imputed accordingly for higher accuracy in statistical findings. This study however is going to focus on the Gaussian detection techniques.

2. Methods

For normal observations, the outlier detection technique by Gaussian model stipulates that an outlier is given by a

point $X > |\bar{X} \pm 3\sigma|$. Since the arithmetic mean as a measure of central tendency that is affected by the outliers, and we know that in a perfectly symmetric data, the mean, the mode and the median are the same [20, 21], the study replaced the mean with the median since the median is a measure of central tendency that is not affected by the existence of the outliers in the dataset [22]. This lead to the same equation given as;

$$X > |Med \pm 3\sigma| \tag{1}$$

The study expects the formula to be better than the Gaussian outlier detection method given we have done away with a measure of central tendency that is affected by the outliers. Moreover, since the standard deviation is a measure of variation that is affected by the existence of the outliers in the dataset and definitely will affect the accuracy of the detection, the study therefore found it necessary to replace the standard deviation with a geometric mean. The geometric mean however was calculated around the median, a measure of central tendency that is not affected by the outliers [22, 23], so as to come up with a geometric mean that is also not affected by the existence of the outliers in the dataset. The study expects this to make our outlier detection formula even better. The formula will therefore be given as;

$$X > |Med \pm 3G| \tag{2}$$

The next task now is to calculate the geometric averages with respect to the mean, this is given as, the study borrowed a concept from [24]

$$G = \sqrt[n]{\prod_{i=1}^n (xi - Med)}$$

While formulating the G, the study established that most important is the deviation from the median can either take a positive, a negative or a zero value, making the formula not applicable in an event we get a negative value since we cannot get a real root of a negative number. In response to this shortcoming, the study took the absolute of the deviations given the rule of geometric averaging holds that most important is the magnitude of the deviation and not the direction [25, 26]. By doing so, we obtain the equation as;

$$G = \sqrt[n]{\prod_{i=1}^n |xi - Med|}$$

The next challenge comes when some data points are equal to the median leading to zero deviation, thus, the product of the deviation from the median will eventually be zero leading to indefinite root. In response to this problem, the study added an arbitrary constant k. The study derived constant K by identifying the best constant that best detects outliers in low dimensional data-set and as well have the least effect to the deviation, the constant was obtained by plotting these constants against outliers removed in given sets of data. This is as shown in Figure 7 below;

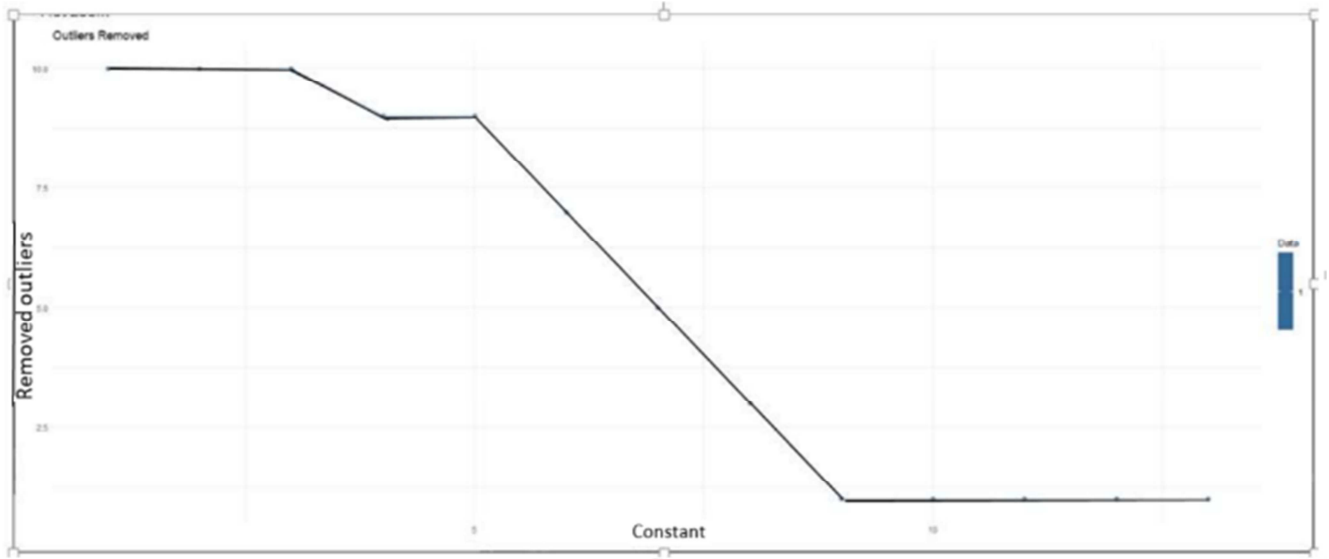


Figure 7. Curve used to derive the most appropriate k-value.

The curve flattened at the 9th constant which was 0.1, even though the other constants removed the same number of outliers after flattening, the study considered the least of these constants (0.1) which will have the least effect on the deviation from the median given the study doesn't want to interfere much with the deviation from the median. The formula therefore becomes;

$$G = \sqrt[n]{\prod_{i=1}^n |di| + k, i = 1 \dots n}$$

$xi = Med$

Where $di=|xi-Med|$ and $k=0.1$

The study further introduced logarithms in order to help in eliminating infinite number that are likely to arise when the population size is large as the geometric measure of variation from the median is being calculated. This gives;

$$G = \begin{cases} \exp\left(\frac{\sum_{i=1}^n \log(|xi-Med|+k)}{n}\right), & x1 \neq x2 \neq \dots \neq Med \\ 0, & x1 = x2 = \dots = Med \end{cases}$$

Where $k=0.1$

Therefore the new formula for outlier detection will state that an outlier is any point given as;

$$X > |Med \pm 3G|$$

3. Results

To test for the effectiveness of the newly invented formula, the study examined and compared the sensitivity of the two Gaussian detection techniques ($Xi > Q3+k(IQR)$ or $Xi < Q1-k(IQR)$).

The study formulated random normal data-sets to help in randomly obtaining the data for simulation. This was done by combining two randomly formulated data-sets with varying mean and standard deviation to help examine the effect of changing either the mean or the standard deviation to the sensitivity of the model. Given the Gaussian formula under scrutiny ($X > |\mu + 3\sigma|$) uses two measure of central tendency that are affected by the outliers, this process was done in two ways;

1. Combining two data-sets with constant standard deviation but varying mean

The first simulation is summarized in the table below;

Table 1. Outliers detected by all formulas with first set of data.

Technique	New formula	Gaussian 1st equation	Gaussian 2 nd Equation
Outliers available	5	5	5
Outliers removed	7	6	4

From table 1, the first simulation, the study combined the first data set ($X1; n=50, \mu = 4$ and $\sigma = 2$) with ($X2; n=5, \mu = 30$ and $\sigma = 2$), with ($X2; n=5, \mu = 30$ and $\sigma = 2$). The expected 5 outliers from the combined datasets were then subjected to the three detection techniques. The study simulated the dataset in the first Gaussian detection technique ($Xi > Q3 +k(IQR)$) or $Xi < Q1 -k(IQR)$), 6 outliers were detected as shown in the table above (Refer to Appendix 1: Figure 8). When the same dataset was simulated against the

second Gaussian equation ($X > |\mu + 3\sigma|$), 4 outliers were detected (Refer to Appendix 1: Figure 9).

When the outliers were simulated in the new equation, 7 outliers were detected (refer to Appendix 1: Figure 10).

From the results, the new formula eliminated the most number of outliers (7), the number exceeds the expected number of outliers, 5, since when two datasets of different means are combined they form a new mean, therefore the outliers are likely to be more or less than were suppose.

The Gaussian second equation performed second best with 6 outliers eliminated. The Gaussian equation under scrutiny managed to detect only 4 outliers, this may be attributed to by the fact that it use the measures of central tendency that is

affected by the outliers.

The study examined another pair of datasets, (X3; n=250, $\mu = 15$ and $\sigma = 5$) and (X4; n=10, $\mu = 80$ and $\sigma = 5$).

The table below summarizes the simulations.

Table 2. Outliers detected by all formulas with second set of data.

Technique	New formula	Gaussian 1st equation	Gaussian 2 nd Equation
Outliers available	10	10	10
Outliers removed	22	12	10

From table 2, when the data was simulated using the first Gaussian formula (refer to Appendix 2: Figure 11), 10 outliers were detected. When the data was simulated against the second Gaussian formula (refer to Appendix 2: Figure 12), 10 outliers were detected.

Finally, when the dataset was simulated against the new detection technique and 22 outliers were detected.

The new formula eliminated the more outliers compared to the rest which may be contributed to by the fact that it uses measures that are least affected by the outliers. The Gaussian equation with the interquartile range once again performed better, eliminating 12 outliers. The equation under scrutiny

(the second Gaussian equation) eliminated the least number of outliers.

The number of outliers removed differ and may be even more than expected because the moment two datasets are combined, they form a new mean interfering with the number of outliers as outliers by definition, are observations that lip an abnormal distance from the mean.

The study carried another sensitivity test on another sets of datasets by combining (X5; n=150, $\mu = 90$ and $\sigma = 5$) with (X6; n=25, $\mu = 200$ and $\sigma = 5$).

The summary is as shown in the table below;

Table 3. Outliers detected by all formulas with third set of data.

Technique	New formula	Gaussian 1st equation	Gaussian 2 nd Equation
Outliers available	10	10	10
Outliers removed	27	26	26

When the outliers were simulated using the first Gaussian equation ($X_i > Q_3 + k(IQR)$ or $X_i < Q_1 - k(IQR)$), 26 outliers were eliminated (Refer to Appendix 3: Figure 14).

When the outliers were simulated using the Gaussian second equation, 27 outliers (refer to Appendix 3: Figure 15).

Finally, when the data was simulated against the new detection technique, 26 outliers were detected (Refer to Appendix 3: Figure 16).

The three techniques performed relatively the same even though the new technique eliminated one more outlier than the rest.

2. Combining two data-sets with constant mean but varying standard deviation

The study examined the sensitivity by combining (X7; n=250, $\mu = 40$ and $\sigma = 45$) with (X8; n=15, $\mu = 40$ and $\sigma = 5$).

The finding are summarized in the table as shown;

Table 4. Outliers detected by all formulas with first set of data.

Technique	New formula	Gaussian 1st equation	Gaussian 2 nd Equation
Outliers available	15	15	15
Outliers removed	27	4	1

Simulating the outliers using the Gaussian first equation, 4 outliers were removed. The Gaussian second equation eliminated 1 outlier and the new detection technique removed 27 outliers.

The study also examined the following datasets;

(X9; n=500, $\mu = 20$ and $\sigma = 10$) and (X10; n=55, $\mu = 20$ and $\sigma = 5$). The summary is as shown in the table below;

Table 5. Outliers detected by all formulas with second set of data..

Technique	New formula	Gaussian 1st equation	Gaussian 2 nd Equation
Outliers available	55	55	55
Outliers removed	72	8	3

Simulating the outliers using the Gaussian first equation, 4 outliers were ejected from the data-set, the Gaussian second equation detected 1 outlier while the new equation once again removed the most, 27. This may be as a result the use of measures of central tendencies that are least affected by

the outliers. Lastly, the study examined the sensitivity in one more pair of data-sets; (X9; n=500, $\mu = 20$ and $\sigma = 10$) with. When the outliers were detected, the Gaussian first equation ejected 8 outliers, 3 by the second equation and 72 by the new equation. The new equation is stricter because the

measures it uses are not affected by the outliers.

4. Summary

The study sought to determine the Outlier detection for univariate data set using geometric technique. The study sought to empirically detect outliers using the univariate normal outlier detection technique in simulated data. The study also aim to measure precision of the univariate outlier detection model in comparison to the Gaussian outlier detection models. The data used in this study was randomly generated from a normal distribution. This chapter gives a summary of the findings, makes conclusions and recommendations based on the findings.

1) Summary of findings

The study sought to establish an outlier detection technique for univariate normal datasets. The measures of central tendency in the Gaussian equation ($X > \mu + 3\sigma$) which are highly affected by the outliers were replaced by those that are least affected by the outliers to form a new equation which defined an outlier as a point $X > Med + 3G$.

The study sought to empirically detect outliers using univariate normal outlier detection technique in simulated data. The normal data sets were randomly generated and

simulation done using the new formula, the study noted that the formula was able to detect the outliers in the data-set. Univariate normal outlier detection model in comparison to the Gaussian outlier detection model. The study formulated same sets of datasets and observed the sensitivity of the models.

2) Conclusion

After conducting sensitivity test on several sets of datasets, the study established that the new formula is the best in outlier detection, in all cases examined it performed better than the Gaussian detection model ($X > \mu + 3\sigma$). The second Gaussian equation ($X_i > Q_3 + k(IQR)$ or $X_i < Q_1 - k(IQR)$) performed as well as the new formula ($X > Med \pm 3G$).

When there was variation in mean but constant standard deviation. This, however, was not the case when standard deviation was varied with constant mean, in this case the new model detected more outliers than any of the two Gaussian equations. This is due to the fact that the new equation used on the measures that are least affected by the outliers.

The new equation proved more sensitive and precise in outlier detection. Even though the new technique was stricter when the standard deviation was varied, it still stands as the best technique according to the study as an outlier is defined relative to the mean and not the standard deviation.

Appendix

Appendix 1. Outlier Detection Simulation for Small Data Sets with Mean Variation

```
> set.seed(45)
> x1<-rnorm(50,mean=4,sd=2)
> x2<-rnorm(5,mean=30,sd=2)
> v=c(x1,x2)
> v[v<quantile(v,.25)-1.5*IQR(v) | v>quantile(v,.75)+1.5*IQR(v)] <- NA #Gaussian Eqn 1
v
[1] 4.681599 2.593319 3.240925 2.507905 2.203785 3.330412 2.997244 3.650929 7.618075 3.539790 1.739164 4.431978 6.464475 7.218717
[15] 4.803101 3.454032 3.927695 3.699378 NA 0.695008 1.729710 4.455340 3.633363 3.172963 3.124809 3.947631 2.280332 4.333089
[29] 6.950981 4.390846 4.318844 2.559613 2.128995 4.570865 2.521530 4.858298 9.467968 1.333193 7.720191 4.491940 2.508021 1.031724
[43] 4.444097 4.955655 5.462397 4.344210 6.373384 3.299181 6.295195 6.700168 NA NA NA NA NA
> sum(is.na(v))
[1] 6
```

Figure 8. Sensitivity of Gaussian 1st equation with variation in mean.

```
> set.seed(45)
> x1<-rnorm(50,mean=4,sd=2)
> x2<-rnorm(5,mean=30,sd=2)
> v=c(x1,x2)
> v[v< mean(v)-3*sd(v) | v> mean(v)+3*sd(v)] <- NA #Gaussian Eqn 2
v
[1] 4.681599 2.593319 3.240925 2.507905 2.203785 3.330412 2.997244 3.650929 7.618075 3.539790 1.739164 4.431978
[13] 6.464475 7.218717 4.803101 3.454032 3.927695 3.699378 11.537621 0.695008 1.729710 4.455340 3.633363 3.172963
[25] 3.124809 3.947631 2.280332 4.333089 6.950981 4.390846 4.318844 2.559613 2.128995 4.570865 2.521530 4.858298
[37] 9.467968 1.333193 7.720191 4.491940 2.508021 1.031724 4.444097 4.955655 5.462397 4.344210 6.373384 3.299181
[49] 6.295195 6.700168 NA NA NA 26.944019 NA
> sum(is.na(v))
[1] 4
```

Figure 9. Sensitivity of Gaussian 2nd equation with variation in mean.

```
> set.seed(45)
> x1<-rnorm(50,mean=4,sd=2)
> x2<-rnorm(5,mean=30,sd=2)
> v=c(x1,x2)
> v[v<quantile(v,.50)-3*exp(sum(Log(abs(v-quantile(v,.50))+1))/length(v)) |
+ v>quantile(v,.50)+3*exp(sum(Log(abs(v-quantile(v,.50))+1))/length(v))] <- NA #New Eqn
v
[1] 4.681599 2.593319 3.240925 2.507905 2.203785 3.330412 2.997244 3.650929 7.618075 3.539790 1.739164 4.431978 6.464475 7.218717
[15] 4.803101 3.454032 3.927695 3.699378 NA 0.695008 1.729710 4.455340 3.633363 3.172963 3.124809 3.947631 2.280332 4.333089
[29] 6.950981 4.390846 4.318844 2.559613 2.128995 4.570865 2.521530 4.858298 NA 1.333193 7.720191 4.491940 2.508021 1.031724
[43] 4.444097 4.955655 5.462397 4.344210 6.373384 3.299181 6.295195 6.700168 NA NA NA NA NA
> sum(is.na(v))
[1] 7
```

Figure 10. Sensitivity of the New Equation with variation in mean.

Appendix 3. Outlier Detection Simulation for Large Data Sets with Mean Variation

```

> set.seed(45)
> x7<-rnorm(150,mean=90,sd=5)
> x8<-rnorm(25,mean=200,sd=5)
> v=c(x7,x8)
> v[v<quantile(v,.25)-1.5*IQR(v) | v>quantile(v,.75)+1.5*IQR(v)] <- NA #Gaussian Eqn 1
> v
[1] 91.70400 86.48330 88.10231 86.26976 85.50946 88.32603 87.49311 89.12732 99.04519 88.84948 84.34791 91.07994
[13] 96.16119 98.04679 92.00775 88.63508 89.81924 89.24844 NA 81.73752 84.32427 91.13835 89.08341 87.93241
[25] 87.81202 89.86908 85.70083 90.83272 97.37745 90.97711 90.79711 86.39903 85.32249 91.42716 86.30382 92.14574
[37] 103.66992 83.33298 99.30048 91.22985 86.27005 82.57931 91.11024 92.38914 93.65599 90.86053 95.93346 88.24795
[49] 95.73799 96.75042 95.58075 91.01173 90.55745 82.36005 90.28754 100.93730 89.64523 96.46498 91.88438 85.91745
[61] 89.64153 79.20747 100.09597 89.74935 88.59893 91.58919 85.15548 87.49829 85.19538 82.04977 91.87867 95.14599
[73] 96.13116 95.24197 90.83707 95.68999 95.28363 86.59921 83.31305 80.98449 93.50489 79.77955 85.05972 93.13691
[85] 88.08129 108.01761 91.66105 83.77471 91.23683 82.97986 94.53591 86.48023 82.15114 91.72656 93.39969 95.71359
[97] 92.02958 85.40680 101.48157 86.55151 80.60851 95.32172 95.69164 93.89804 86.83979 90.94214 85.45155 89.31991
[109] 91.14678 86.37434 93.80565 89.60261 98.63746 92.02601 82.98189 90.33236 91.52715 94.36589 89.85526 90.65314
[121] 90.37835 93.93083 81.86099 90.39073 86.92468 90.46637 90.82302 84.63411 88.12258 89.90127 93.04207 92.37104
[133] 81.54538 94.60573 92.60130 92.17799 88.80717 84.58242 93.90847 86.33207 90.91520 92.11036 89.86157 92.27670
[145] 96.28178 83.34470 85.56218 89.47476 89.54752 80.99463 NA NA NA NA NA NA
[157] NA NA NA NA NA NA NA NA NA NA NA NA
[169] NA NA NA NA NA NA NA NA NA NA NA NA
> sum(is.na(v))
[1] 26
    
```

Figure 14. Sensitivity of Gaussian 1st equation with variation in mean.

```

> set.seed(45)
> x7<-rnorm(150,mean=90,sd=5)
> x8<-rnorm(25,mean=200,sd=5)
> v=c(x7,x8)
> v[v< mean(v)-3*sd(v) | v> mean(v)+3*sd(v)] <- NA #Gaussian Eqn 2
> v
[1] 91.70400 86.48330 88.10231 86.26976 85.50946 88.32603 87.49311 89.12732 99.04519 88.84948 84.34791 91.07994
[13] 96.16119 98.04679 92.00775 88.63508 89.81924 89.24844 108.84405 81.73752 84.32427 91.13835 89.08341 87.93241
[25] 87.81202 89.86908 85.70083 90.83272 97.37745 90.97711 90.79711 86.39903 85.32249 91.42716 86.30382 92.14574
[37] 103.66992 83.33298 99.30048 91.22985 86.27005 82.57931 91.11024 92.38914 93.65599 90.86053 95.93346 88.24795
[49] 95.73799 96.75042 95.58075 91.01173 90.55745 82.36005 90.28754 100.93730 89.64523 96.46498 91.88438 85.91745
[61] 89.64153 79.20747 100.09597 89.74935 88.59893 91.58919 85.15548 87.49829 85.19538 82.04977 91.87867 95.14599
[73] 96.13116 95.24197 90.83707 95.68999 95.28363 86.59921 83.31305 80.98449 93.50489 79.77955 85.05972 93.13691
[85] 88.08129 108.01761 91.66105 83.77471 91.23683 82.97986 94.53591 86.48023 82.15114 91.72656 93.39969 95.71359
[97] 92.02958 85.40680 101.48157 86.55151 80.60851 95.32172 95.69164 93.89804 86.83979 90.94214 85.45155 89.31991
[109] 91.14678 86.37434 93.80565 89.60261 98.63746 92.02601 82.98189 90.33236 91.52715 94.36589 89.85526 90.65314
[121] 90.37835 93.93083 81.86099 90.39073 86.92468 90.46637 90.82302 84.63411 88.12258 89.90127 93.04207 92.37104
[133] 81.54538 94.60573 92.60130 92.17799 88.80717 84.58242 93.90847 86.33207 90.91520 92.11036 89.86157 92.27670
[145] 96.28178 83.34470 85.56218 89.47476 89.54752 80.99463 203.55123 200.20106 201.49178 206.21652 198.23154 185.74605
[157] 204.39366 189.76887 208.22311 197.90699 204.58194 201.80422 211.02346 204.32687 207.48455 203.35738 193.95942 197.29920
[169] 195.43880 194.67938 194.23424 201.54845 200.59899 203.08030 195.52335
> sum(is.na(v))
[1] 0
    
```

Figure 15. Sensitivity of Gaussian 2nd equation with variation in mean.

```

> set.seed(45)
> x7<-rnorm(150,mean=90,sd=5)
> x8<-rnorm(25,mean=200,sd=5)
> v=c(x7,x8)
> v[v<quantile(v,.50)-3*exp(sum(log(abs(v-quantile(v,.50))+.1))/length(v)) |
+ v>quantile(v,.50)+3*exp(sum(log(abs(v-quantile(v,.50))+.1))/length(v))] <- NA #New Eqn
> v
[1] 91.70400 86.48330 88.10231 86.26976 85.50946 88.32603 87.49311 89.12732 99.04519 88.84948 84.34791 91.07994
[13] 96.16119 98.04679 92.00775 88.63508 89.81924 89.24844 NA 81.73752 84.32427 91.13835 89.08341 87.93241
[25] 87.81202 89.86908 85.70083 90.83272 97.37745 90.97711 90.79711 86.39903 85.32249 91.42716 86.30382 92.14574
[37] 103.66992 83.33298 99.30048 91.22985 86.27005 82.57931 91.11024 92.38914 93.65599 90.86053 95.93346 88.24795
[49] 95.73799 96.75042 95.58075 91.01173 90.55745 82.36005 90.28754 100.93730 89.64523 96.46498 91.88438 85.91745
[61] 89.64153 79.20747 100.09597 89.74935 88.59893 91.58919 85.15548 87.49829 85.19538 82.04977 91.87867 95.14599
[73] 96.13116 95.24197 90.83707 95.68999 95.28363 86.59921 83.31305 80.98449 93.50489 79.77955 85.05972 93.13691
[85] 88.08129 NA 91.66105 83.77471 91.23683 82.97986 94.53591 86.48023 82.15114 91.72656 93.39969 95.71359
[97] 92.02958 85.40680 101.48157 86.55151 80.60851 95.32172 95.69164 93.89804 86.83979 90.94214 85.45155 89.31991
[109] 91.14678 86.37434 93.80565 89.60261 98.63746 92.02601 82.98189 90.33236 91.52715 94.36589 89.85526 90.65314
[121] 90.37835 93.93083 81.86099 90.39073 86.92468 90.46637 90.82302 84.63411 88.12258 89.90127 93.04207 92.37104
[133] 81.54538 94.60573 92.60130 92.17799 88.80717 84.58242 93.90847 86.33207 90.91520 92.11036 89.86157 92.27670
[145] 96.28178 83.34470 85.56218 89.47476 89.54752 80.99463 NA NA NA NA NA NA
[157] NA NA NA NA NA NA NA NA NA NA NA NA
[169] NA NA NA NA NA NA NA NA NA NA NA NA
> sum(is.na(v))
[1] 27
    
```

Figure 16. Sensitivity of the New Equation with variation in mean.

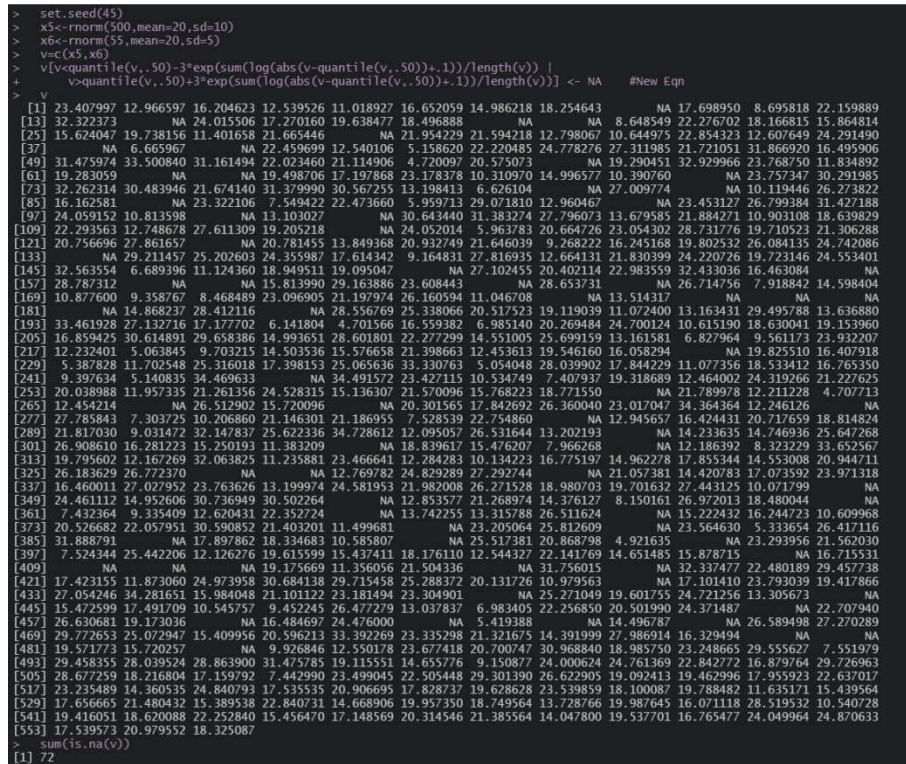


Figure 22. Sensitivity of the New Equation with variation in standard deviation.

References

[1] I. Ben-Gal, "Outlier detection," in Data mining and knowledge discovery handbook, Springer, 2005, pp. 131-146.

[2] P. L. Clark, "Number theory: A contemporary introduction", 2012.

[3] C. E. Shannon, A mathematical theory of communication, vol. 27, Bell System Technical Journal, 1948, pp. 379-423.

[4] D. Papadopolus, T. Palpanas, D. Gonupulos and V. Kalogeraki, "Distributed deviation detection in sensor networks", vol. 32, Acm sigmod record, 2003, pp. 77-82.

[5] J. Orsborne and A. Overbay, "The power of outliers (and why researchers should always check for them)", vol. 9, Practical Assessment, Research and Evaluation, 2004, p. 6.

[6] X. Li and J. Han, "Mining approximate top k subspace anomalies in multidimensional time series data", in VDLBD, 2007, pp. 447-458.

[7] V. Ganti, R. Ramakrishnan, J. Gehrke, A. Powell and J. French, "Clustering large datasets in arbitrary metric spaces", in Proceedings 15th International Conference on Data Engineering (Cat. No. 99CB36337), pp. 502-511.

[8] S. D. Bay and M. Schwabacher, "Mining distance-based outliers in near linear time with randomization and a sample pruning rule", in Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining, 2003, pp. 29-38.

[9] E. Eskin, A. Arnold, M. Prerav, L. Portnoy and S. Stolfo, "A geometric framework for unsupervised anomaly

detection", in Applications of data mining in computer security. Springer, 2002, pp. 77-101.

[10] J. Han and M. Kamber, "Data mining concepts and techniques", San Francisco: morgan kaufmann publishers.

[11] V. Barnett, "The ordering of multivariate data", vol. 139, Journal of royal statistical society Series A (General), 1976, pp. 318-344.

[12] P. J. Rousseeuw and M. Hubert, "Robust statistics for outlier detection", vol. 1, Wiley interdisciplinary reviews. Data mining and knowledge discovery, 2011, pp. 73-79.

[13] C. C. Aggarwals, J. Han, J. Wang and P. S. Yu, A framework for projected clustering of high dimensional data streams", vol. 30, in Proceedings of the Thirtieth international conference on very large databases, 2004, pp. 852-863.

[14] P. C. Wu, "The Central Limit Theorem and comparing means, trimmed means, one-step M-estimators and modified one-step M-estimators under non-normality", University of Southern California, 2002.

[15] A. Biswas and A. Bisaria, "A test of normality from allegorizing the bell curve or the gaussian probability distribution as memoryless and depthless like a black hole", vol. 14, Applied Mathematics Sciences, 2020, pp. 349-359.

[16] R. Lugannani and S. Rice, "Saddle point approximation for the distribution of the sum of independent random variables", vol. 12, Advances in applied probability, 1980, pp. 475-490.

[17] C. Leys, C. Ley, O. Klein, P. Bernard and L. Licata, Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median", vol. 49, Journal of experimental social psychology, 2013, pp. 764-766.

- [18] P. J. Rouseeuw and B. C. Van Zomeren, "Unmasking multivariate outliers and leverage points", vol. 85, *Journal of the American Statistics association*, 1990, pp. 633-639.
- [19] V. L. Sourd, "Performance measurement for traditional investment", vol. 58, *Financial Analysis Journal*, 2007, pp. 36-52.
- [20] P. V. Hippel, "Mean, median and skew: correcting a textbook rule", vol. 13, *Journal of statistics Education*, 2005.
- [21] E. M. Knorr and R. T. Ng, "Finding intensional knowledge of distance-based outliers", vol. 99, in *Vldb*, 1999, pp. 211-222.
- [22] W. Dixon, "Processing data for outliers", vol. 9, *Biometry*, 1953, pp. 74-89.
- [23] F. Angiulli and C. Pizzuti, "fast outlier detection in high dimensional spaces, principals of data mining and knowledge discovery", 2002.
- [24] B. Troon, "Estimating average variation about the population mean using geometric measure of variation", 2020.
- [25] T. Li, H. Fan, J. Garcia and J. M. Corchado, "Second order statistics analysis and comarison between arithmetic and geometric average fusion: Application to multi-sensor target tracking", vol. 51, *Information Fusion*, 2019, pp. 233-243.
- [26] V. Barnett and T. Lewis, "Outliers in statistical data, Wiley series in Probability and Mathematical statistics. Applied Probability and Statistics 1984.