



ILJS-19-002

Assessment of Robustness of some Measures of Variation with Normal and Non-Normal Data Sets

Adeleke^{1*}, M. O., Adeleke¹, B. L. and Jimoh², K.

¹Department of Statistics, University of Ilorin, Ilorin, Nigeria.

² Department of Physical Sciences, Al-Hikmah University, Ilorin Nigeria.

Abstract

In any investigation where numerical values are obtained, it is always desirable to have a typical value for all the observations, and the mean as a measure of central tendency is commonly used. Reliability of the value of the mean is strengthened when a corresponding measure of variation (also known as dispersion) for the data is obtained. This paper therefore, presents a study of robustness of some measures of dispersion namely, the variance, standard deviation, absolute mean deviation with divisor 'n' (AMD(n)), and absolute mean deviation with divisor 'n-1' (AMD(n-1)). The level of robustness of the measures of dispersion in this paper was facilitated by the adoption of simulation technique that utilized the following: small sample sizes; and large sample sizes, for both normal and non-normal data sets of different specifications. Overall, the results obtained showed that AMD (n-1) gave values that were closest in magnitude to standard deviation. The implication of the findings herein is that all the three measures of spread proved to be robust, however AMD (n-1) is a better substitute for the standard deviation.

Keyword: Mean; Standard Deviation; Absolute Mean Deviation; Normal and Non-Normal Data; Data Simulation.

1. Introduction

Several situations exist in life, that as a matter of necessity have to be investigated. Meanwhile in the course of conducting an investigation, data are collected and subsequently analysed, for the purpose of eliciting useful information on human existence for the purpose of charting a rewarding course of action to mitigate existing challenges or to sustain the existing level of comfort or satisfaction. Data can be obtained from several sources, few of which are Official Research, Business Activity Records, census statistical survey and planned experiment.

*Corresponding Author: Adeleke M. O.
Email: amariam47@yahoo.com

Therefore, for any set of data collected in respect of an investigation, it is paramount to obtain a numerical value that utilizes all the available data points in terms of the average or the centre of the data. As much as it is desirable to have a measure of the average of the entire data set, it is also paramount to have a measure of dispersion or spread of the observations that are collected on the basis of the effort of the investigator. This is expected to be a desirable indicator of establishing whether or not the resulting measure of average gives a value around which all values are closely clustered or otherwise.

One of the most commonly used measures of spread is the variance. This is due to the fact that the variance, unlike several other methods of spread utilizes all the data points. Meanwhile, standard deviation, which is also commonly used, is the positive radical of the variance (McDonald, 2014; Rodrigues, 2017). Two other measures of spread that possess the property of utilizing all the data points are considered in this study for the purpose of evaluating them in different standard settings that are expressed to satisfy the following: normality and non-normality; and small and large sample sizes. Several other measures of spread failed to possess this unique feature (Berry *et al.*, 2019).

Absolute mean deviation is a measure of dispersion that utilises all the observations for its computation, and commonly used equation uses the divisor 'n', see for example Berry *et al.* (2019), Johnson *et al.* (2007) and Keller *et al.* (2003). However, Adeleke (2006) gave the formula for calculating absolute mean deviation with the divisor 'n-1'. Meanwhile, Adegboye (2009) asserted that the formula for calculating absolute mean deviation should have the median as the reference statistic in contrast to the mean see also Hana *et al.* (2017) and Nahmias and Olsen (2015). The foregoing varied assertions by different authors may be unarguably due to the fact that absolute mean deviation utilises all observations, in contrast to several other measures of dispersion that are either of the absolute or relative type.

2. Materials and Methods

The work focuses on the assessment of three measures of variation that are commonly used in statistics to describe data. These measures are variance, standard deviation and mean absolute deviation, see for example El-Amir (2012), Sharma (2007) and Wasserman (2017). These three measures were evaluated via simulation studies, using R software, by varying the sample sizes with the imposition of the normality and non-normality conditions (Horton *et al.*, 2015; R Core Team, 2019). The three measures considered were computed using the formulae given below.

The variance is computed using the definition form, which is expressed as follows:

$$\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2. \quad (1)$$

The standard deviation is computed as:

$$\sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}. \quad (2)$$

The mean absolute deviation is:

$$\frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|. \quad (3)$$

Finally, the mean absolute deviation, with divisor 'n-1' is computed as follows:

$$\frac{1}{n-1} \sum_{i=1}^n |x_i - \bar{x}|. \quad (4)$$

Samples with respectively, small and large sample sizes were considered in the study, with the focus of investigating the performance as well as reliability of the different scaled forms of the deviation of each of the observations from the mean value, in preference to the standard deviation. The scaling factor of concern in this study, have two distinct forms, which are, 'n', and 'n-1'. In statistical theory, when the sample size is quite large, the difference between any two statistics that are with the divisors 'n' and 'n-1' approaches zero.

It is of note that the scaling factor n-1, when 'n-1' is used as the denominator term in equation (4) above, gives one form of the absolute mean deviation, while the second form of the absolute mean deviation has the divisor 'n' as in equation (3)

Further, the denominator term of equation (1) above, provides an estimate that is unbiased of the population variance in contrast to the use of 'n' as the denominator which results in a biased estimator.

Preference for denominator "n-1" for the variance instead of 'n' is premised on the statistical property of unbiasedness. In fact, further justification of the denominator "n-1" is also anchored on the theory associated with the determination of degrees of freedom. The theory in the foregoing concerns the fact that the sample mean is obtained and utilized in the computation of the measures of variation that are considered in this paper. The fact that the mean has to be obtained prior to the calculation of a statistic will attract a penalty of a single degree of freedom.

Definition: Unbiased Estimator

Suppose a parameter θ has an estimator T , then T is said to be unbiased if $E(T) = \theta$. On the other hand, if the estimator T for a parameter θ is biased, the expression $E(T) - \theta$ will be non-zero (Calonico *et al.*, 2018). Overall or be that as it may, there are several other properties that an estimator should satisfy to be preferred to the other.

Standard deviation is more amenable to rigorous or rather higher statistical treatment. Ease of computation has received less attention with the unhindered availability of standard statistical packages that can execute complex statistical computations faster than ever imagined. However, it might be a problematic estimator for skewed data where few observations are extreme (Leys *et al.*, 2013).

3. Result and Discussion

Data simulation was executed with the use of R-Statistical package. Samples of different sizes and under different configurations of normality and non-normality were specified and subsequently implemented 1000 times, in each of the cases.

3.1 Experimental Results obtained from Simulated Data

Tables 1 through 8 presented the summary, with the criteria of mean, minimum and maximum values for each of the 1000 trials of each measure of dispersion in this paper, with the different sample sizes. The mean gives an insight about how close the estimated value are to the true value of the measures. In addition, a measure of the precision of the four measure is examined using the range. The entries in the tables are indeed the summary of massive results of the simulation experiments. In particular, standard normal, that is Normal (0, 1), Normal (2, 4), Exponential (1) and Exponential (0.5) distributions were considered.

The true values of the variance for the non-normal distributions are 1 and 4 for the Exponential (1) and Exponential (0.5) respectively. These non-normal distributions are positively skewed and therefore, will have large extreme observations. As a result, measures, such as variance and standard error is often not useful as it is highly influenced by the extreme observations (Aslam *et al.*, 2019; Gorard, 2015; Mohini and Prajakt, 2012; Zhang, 2016).

Table 1: Small sample size using a standard normal distribution.

Measures		Sample size				
		5	10	15	20	25
Variance	Mean	0.9794	0.9902	0.991	0.9897	1.0054
	Min	0.0338	0.246	0.0807	0.3077	0.3514
	Max	5.3496	3.165	2.6558	2.3966	2.1221
Standard Deviation	Mean	0.9293	0.9701	0.9775	0.9823	0.9919
	Min	0.184	0.496	0.284	0.5547	0.5928
	Max	2.3129	1.7791	1.6297	1.5481	1.4567
Mean absolute deviation	Mean	0.7037	0.7566	0.7665	0.7757	0.7836
	Min	0.1371	0.3158	0.2355	0.4239	0.4361
	Max	1.8348	1.5233	1.2958	1.2315	1.1614
Mean absolute dev (n-1)	Mean	0.8796	0.8407	0.8213	0.8166	0.8162
	Min	0.1714	0.3509	0.2523	0.4462	0.4543
	Max	2.2935	1.6926	1.3884	1.2964	1.2098

Evaluation of measures of spread for standard normal distributed data sets with varying sample sizes of 5, 10, 15, 20, and 25 and keeping track of the summary values, using mean, minimum and maximum values.

Table 2: Large sample size using a standard normal distribution.

Measures		Sample size							
		30	40	50	60	70	80	90	100
Variance	Mean	1.0028	0.9922	0.9918	1.0028	1.0027	0.9971	0.9912	0.9921
	Min	0.3467	0.392	0.4599	0.5552	0.5525	0.6195	0.5757	0.6104
	Max	2.0877	1.7868	1.8217	1.6242	1.5601	1.4896	1.504	1.4493
Standard Deviation	Mean	0.993	0.9897	0.9904	0.9973	0.9978	0.9954	0.9928	0.9935
	Min	0.5888	0.6261	0.6782	0.7451	0.7433	0.7871	0.7588	0.7813
	Max	1.4449	1.3367	1.3497	1.2744	1.2491	1.2205	1.2264	1.2039
Mean absolute deviation	Mean	0.7866	0.7857	0.7868	0.7921	0.7949	0.7918	0.7895	0.7909
	Min	0.4778	0.4745	0.5004	0.5664	0.5714	0.6206	0.5975	0.6299
	Max	1.2409	1.0809	1.063	1.0627	1.0406	0.9896	0.9815	0.9939
Mean absolute dev (n-1)	Mean	0.8137	0.8059	0.8028	0.8055	0.8064	0.8018	0.7984	0.7989
	Min	0.4943	0.4867	0.5106	0.576	0.5797	0.6285	0.6042	0.6362
	Max	1.2836	1.1086	1.0847	1.0808	1.0557	1.0022	0.9925	1.004

Evaluation of measures of spread for standard normal distributed data sets with varying sample sizes of 30, 40, 50, 60, 70, 80, 90, and 100 and keeping track of the summary values, using mean, minimum and maximum values.

Table 3: Small sample size using an Exponential (1) distribution.

Measures		Sample size				
		5	10	15	20	25
Variance	Mean	0.9252	1.0182	0.9859	1.0509	0.9968
	Min	0.0048	0.0577	0.0826	0.1311	0.1787
	Max	16.2239	11.0199	5.4426	6.4412	3.9689
Standard Deviation	Mean	0.8372	0.9313	0.942	0.9835	0.964
	Min	0.0694	0.2403	0.2873	0.3621	0.4227
	Max	4.0279	3.3196	2.3329	2.5379	1.9922
Mean absolute deviation	Mean	0.6323	0.7036	0.7094	0.7322	0.7203
	Min	0.0582	0.1689	0.1991	0.2791	0.3168
	Max	3.5036	2.0494	1.7429	1.5637	1.445
Mean absolute dev (n-1)	Mean	0.7903	0.7818	0.7601	0.7707	0.7503
	Min	0.0728	0.1876	0.2133	0.2938	0.33
	Max	4.3796	2.2771	1.8674	1.646	1.5052

Evaluation of measures of spread for exponential (1) distributed data sets with varying sample sizes of 5, 10, 15, 20, and 25 and keeping track of the summary values, using mean, minimum and maximum values.

Table 4: Large sample size using an Exponential (1) distribution.

Measures		Sample size							
		30	40	50	60	70	80	90	100
Variance	Mean	0.9961	0.997	0.9836	1.006	1.0012	1.02	1.005	0.9949
	Min	0.2134	0.2187	0.317	0.3636	0.2535	0.3068	0.3507	0.4503
	Max	4.6852	4.1661	3.475	2.4265	3.0528	2.5295	2.7839	2.551
Standard Deviation	Mean	0.9695	0.977	0.9735	0.989	0.986	0.9978	0.9922	0.9881
	Min	0.462	0.4677	0.563	0.603	0.5035	0.5539	0.5922	0.671
	Max	2.1645	2.0411	1.8641	1.5577	1.7472	1.5904	1.6685	1.5972
Mean absolute deviation	Mean	0.7224	0.7271	0.7246	0.7325	0.7297	0.7368	0.7335	0.7313
	Min	0.369	0.3606	0.4081	0.4042	0.4077	0.4477	0.4746	0.5148
	Max	1.3092	1.4932	1.2324	1.1143	1.2814	1.0821	1.0626	1.1446
Mean absolute dev (n-1)	Mean	0.7473	0.7457	0.7394	0.7449	0.7403	0.7461	0.7418	0.7387
	Min	0.3817	0.3699	0.4165	0.411	0.4136	0.4533	0.4799	0.52
	Max	1.3543	1.5315	1.2576	1.1332	1.3	1.0958	1.0745	1.1561

Evaluation of measures of spread for exponential (1) distributed data sets with varying sample sizes of 30, 40, 50, 60, 70, 80, 90, and 100 and keeping track of the summary values, using mean, minimum and maximum values.

Table 5: Small sample size using a Normal (1, 4) distribution.

Measures		Sample size				
		5	10	15	20	25
Variance	Mean	4.0826	4.0189	3.9985	4.0447	4.028
	Min	0.019	0.249	0.9689	1.1985	1.1913
	Max	19.3846	12.7059	11.6355	9.4195	9.3713
Standard Deviation	Mean	1.9003	1.9503	1.9663	1.9851	1.9854
	Min	0.138	0.499	0.9843	1.0948	1.0915
	Max	4.4028	3.5645	3.4111	3.0691	3.0613
Mean absolute deviation	Mean	1.4402	1.5175	1.5428	1.5667	1.5691
	Min	0.1179	0.3931	0.8164	0.8014	0.7515
	Max	3.055	2.9475	2.4821	2.4916	2.5185
Mean absolute dev (n-1)	Mean	1.8003	1.6861	1.653	1.6492	1.6344
	Min	0.1473	0.4368	0.8747	0.8436	0.7828
	Max	3.8187	3.275	2.6594	2.6227	2.6234

Evaluation of measures of spread for normal (1, 4) distributed data sets with varying sample sizes of 5, 10, 15, 20, and 25 and keeping track of the summary values, using mean, minimum and maximum values.

Table 6: Large sample size using a Normal (1, 4) distribution.

Measures		Sample size							
		30	40	50	60	70	80	90	100
Variance	Mean	4.0362	3.9911	3.9867	4.017	3.9441	4.0024	3.998	3.9891
	Min	1.2938	1.6372	1.9494	2.0453	1.9714	2.2732	2.6117	2.2936
	Max	8.6679	7.208	7.5413	6.7388	6.8742	6.2846	6.7128	5.9889
Standard Deviation	Mean	1.9928	1.9858	1.9864	1.9957	1.9784	1.9941	1.9939	1.9922
	Min	1.1375	1.2795	1.3962	1.4302	1.4041	1.5077	1.6161	1.5145
	Max	2.9441	2.6848	2.7461	2.5959	2.6219	2.5069	2.5909	2.4472
Mean absolute deviation	Mean	1.5779	1.5767	1.579	1.5846	1.5718	1.586	1.5855	1.5843
	Min	0.906	1.0316	1.1009	1.0477	1.1019	1.172	1.2127	1.221
	Max	2.3729	2.1429	2.1981	2.1521	2.158	2.0324	2.1194	1.9932
Mean absolute dev (n-1)	Mean	1.6324	1.6171	1.6112	1.6115	1.5946	1.606	1.6033	1.6003
	Min	0.9373	1.058	1.1234	1.0655	1.1178	1.1869	1.2263	1.2334
	Max	2.4547	2.1979	2.243	2.1886	2.1893	2.0581	2.1433	2.0133

Evaluation of measures of spread for normal (1, 4) distributed data sets varying sample sizes of 30, 40, 50, 60, 70, 80, 90, and 100 and keeping track of the summary values, using mean, minimum and maximum values.

Table 7: Small sample size using an Exponential (0.5) distribution.

Measures		Sample size				
		5	10	15	20	25
Variance	Mean	4.2063	3.7769	4.0152	3.903	3.9957
	Min	0.0123	0.2009	0.3298	0.565	0.6466
	Max	53.6765	41.8969	23.1786	22.1632	32.5796
Standard deviation	Mean	1.7618	1.7955	1.8993	1.888	1.928
	Min	0.1111	0.4482	0.5743	0.7517	0.8041
	Max	7.3264	6.4728	4.8144	4.7078	5.7079
Mean absolute deviation	Mean	1.3338	1.3561	1.4222	1.4173	1.4515
	Min	0.0818	0.294	0.4662	0.5128	0.654
	Max	5.1178	4.717	3.4867	2.8583	2.9348
Mean absolute dev (n-1)	Mean	1.6672	1.5068	1.5238	1.4919	1.512
	Min	0.1022	0.3266	0.4995	0.5398	0.6812
	Max	6.3972	5.2411	3.7357	3.0088	3.0571

Evaluation of measures of spread for exponential (0.5) distributed data sets varying sample sizes of 5, 10, 15, 20, and 25 and keeping track of the summary values, using mean, minimum and maximum values.

Table 8: Large sample size using an Exponential (0.5) distribution.

Measures		Sample size							
		30	40	50	60	70	80	90	100
Variance	Mean	3.9261	4.1234	4.0488	4.0447	3.9621	3.9673	4.0117	3.958
	Min	0.4797	0.7889	0.9667	1.3309	1.5155	1.3625	1.6394	1.5873
	Max	13.5184	19.545	15.725	10.631	10.935	10.599	8.6996	11.136
Standard Deviation	Mean	1.9276	1.9784	1.978	1.9801	1.9645	1.9682	1.981	1.9706
	Min	0.6926	0.8882	0.9832	1.1536	1.2311	1.1672	1.2804	1.2599
	Max	3.6767	4.421	3.9655	3.2605	3.3068	3.2556	2.9495	3.337
Mean absolute deviation	Mean	1.4346	1.4679	1.4645	1.4654	1.4581	1.4554	1.4639	1.4609
	Min	0.5697	0.7027	0.8022	0.9049	0.9354	0.9108	0.9738	0.99
	Max	2.8423	2.7744	2.249	2.2119	2.2666	2.4418	2.1385	2.0318
Mean absolute dev (n-1)	Mean	1.4841	1.5055	1.4943	1.4903	1.4792	1.4738	1.4803	1.4757
	Min	0.5893	0.7207	0.8186	0.9203	0.949	0.9223	0.9847	1.0000
	Max	2.9403	2.8455	2.2949	2.2494	2.2995	2.4728	2.1625	2.0523

Evaluation of measures of spread for exponential (0.5) distributed data sets varying sample sizes of 30, 40, 50, 60, 70, 80, 90, and 100 and keeping track of the summary values, using mean, minimum and maximum values.

3.2 Discussion of Results

The results obtained in this paper covered the following settings: small sample size with standard normal, Exponential (1), Normal (1, 4), and Exponential (0.5); and large sample size with standard normal, Exponential (1), Normal (1, 4), and Exponential (0.5).

Results in column 1 of Table 1 indicate the following, for 1000 samples of size 5; the mean of the variances was 0.9794, while the minimum and maximum variances were 0.0338 and 5.3496; the mean of the standard deviation was 0.9293, while the minimum and maximum standard deviation were 0.184 and 2.3129; the mean of the mean deviation with the divisor 'n' was 0.7037, while the corresponding minimum and maximum values were 0.1371 and 1.8348; and the mean of the mean absolute deviation was 0.8796, while the corresponding minimum and maximum values were 0.1714 and 2.2935. All other entries in Table 1, as well as the entries in Tables 2 through 8, have similar interpretations.

In all the tables, it was observed that the mean estimate of variances and the corresponding standard deviations were close to the true value. Moreover, the values were closer to the true values and more precise as the sample size increases.

For the non-normal distributions, the variability of the estimates is higher compared to the normal distributions, as evident from the range, which is the difference between minimum and maximum values. This is in line with the knowledge that for skewed distribution, the standard deviation might not be a valid measure of dispersion. However, this did not alter the trend of the highest value, for the standard deviation followed by AMD(n-1), with the least value for AMD(n).

The implication of the foregoing interpretations is that the values obtained for the standard deviation, each of which was obtained as a positive radical of the variance, had the corresponding values obtained in respect of the absolute mean deviation with the divisor 'n-1' as the next lower values in each of the three categories of the mean, minimum and maximum values. Hence, the absolute mean deviation with the divisor 'n-1' is preferable to the absolute mean deviation with the divisor 'n' as a valid substitute of the standard deviation.

4. Conclusion

This paper concluded that, with all the configurations considered in respect of sample size, normality and non-normality of data sets, the performance of the three measures of dispersion

were found to be consistent, in the magnitude recorded by each of them. Absolute mean deviation with the divisor 'n-1', that is AMD (n-1), is therefore recommended as a valid substitute for the standard deviation in preference to the AMD (n) even for non-normal distribution where the standard deviation is highly affected by extreme observations.

References

- Adegboye, O. S. (2009): *Descriptive Statistics*. Olad Publishers, Ilorin, Nigeria.
- Adeleke, B. L. and Isah, A. (2006): *College Statistics and Probability*. Y- Books, A Division of Associated Book Maker Ltd., Ibadan, Nigeria.
- Aslam, M., Rao, S., Al-Marshadi, A. H., Ahmad L. and Jun C. (2019): Control charts for monitoring process capability index using median absolute deviation for some popular distributions. *Processes*. **7**(5), 287-301.
- Berry, K. J., Johnston, J. E. and Mielke, P. W. (2019): *Central Tendency and Variability*. In: *A Primer of Permutation Statistical Methods*. Springer, Switzerland.
- Calonico, S., Cattaneo, M. D., Farrell M. H. (2018): On the effect of bias estimation on coverage accuracy in nonparametric inference. *Journal of the American Statistical Association*. **113**(522), 767-779.
- El-Amir, E. A. H. (2012): On uses of mean absolute deviation: decomposition, skewness and correlation coefficients. *Metron*. **70**, 145–164.
- Gorard, S. (2015): Introducing the mean absolute deviation 'effect' size. *International Journal of Research and Method in Education*. **38**(2), 105-114.
- Hana, B., Wu, Y. and Song, Y. (2017): A novel active contour model based on median absolute deviation for remote sensing river image segmentation. *Computers and Electrical Engineering*. **62**, 209-223.
- Horton, N. J. and Kleinman, K. (2015): *Using R and RStudio for Data Management, Statistical Analysis, and Graphics*. Second Edition. CRC Press, Boca Raton, USA.
- Johnson, R. and Kubly, P. (2007): *Elementary Statistics*. Tenth Edition. Thomson Higher Education, Belmont, USA.
- Keller, G. and Warrack, B. (2003): *Management and Economics*. Seventh Edition. Thomson Higher Education, Belmont, USA.
- Leys, C., Ley, C., Klein, O., Bernard, P. and Licata, L. (2013): Detecting outliers: Do not use standard deviation around the mean, use absolute deviation around the median. *Journal of Experimental Social Psychology*. **49**(4), 764–766.
- McDonald, J. H. (2014): *Handbook of Biological Statistics*. Third Edition. Sparky House Publishing, Baltimore, USA.
- Mohini, P. B. and Prajakt, J. B. (2012): What to use to express the variability of data: Standard deviation or standard error of mean? *Perspectives in Clinical Research*. **3**(3), 113-116.
- Nahmias, S. and Olsen, T. L. (2015): *Production and Operations Analysis*. Seventh Edition. Waveland Press. Long Grove, USA.
- R Core Team (2019): *R- A language and Environment for Statistical Computing*. Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org>.

- Rodrigues, C. F., Lima F. J. and Barbosa F. T. (2017): Importance of using basic statistics adequately in clinical research. *Brazilian Journal of Anesthesiology (English Edition)*. **67**(6), 619-625.
- Sharma, J. K. (2007): *Business Statistics*. Pearson Education, New Delhi, India.
- Wasserman, N. H., Casey S., Champion, J. and Huey, M. (2017): Statistics as unbiased estimators: exploring the teaching of standard deviation. *Research in Mathematics Education*. **19**(3), 236-256.
- Zhang, Z. (2016): Univariate description and bivariate statistical inference: the first step delving into data. *Annals of Translational Medicine*. **4**(5), 91-96.