

RESEARCH PAPER

Normality for Nonnormal Distributions

Koh Khong Liang* and Nor Aishah Ahad

School of Quantitative Sciences, Universiti Utara Malaysia, 06010 Sintok, Kedah, Malaysia

*Corresponding author:leonkk199@gmail.com

DOI: <https://doi.org/10.37134/jsml.vol8.2.7.2020>

Received: 24 March 2020; Accepted: 1 May 2020; Published: 09 June 2020

Abstract

It has been usually assumed that a sample data is normally distributed when the sample size is at least 30. This is the general rule in using central limit theorem based on the sample size being greater or equal to 30. Many literary works also assumed normality when sample size is at least 30. This study aims to determine the least required sample size that satisfy normality assumption from three non-normal distributions, Poisson, Gamma and Exponential distributions. Computer simulations are carried out to study the least required sample size for the three distributions. Through the study, it is found that sample data from Poisson and Gamma distributions need sample size less than 30, while Exponential needs more than 30 to achieve normality.

Keywords: Normal; Sample Size; Poisson; Exponential; Gamma

Abstrak

Terdapat andaian bahawa data sampel terabur secara normal apabila saiz sampel adalah sekurang-kurangnya 30. Ini adalah peraturan umum dalam menggunakan Teorem Had Memusat berdasarkan saiz sampel sama atau lebih besar daripada 30. Banyak kajian lepas juga menganggap normal apabila saiz sampel adalah sekurang-kurangnya 30. Kajian ini bertujuan untuk menentukan saiz sampel paling kecil yang diperlukan untuk memenuhi andaian kenormalan daripada tiga taburan tidak normal, Taburan Poisson, Gamma dan Ekponensial. Simulasi komputer dijalankan untuk mengkaji saiz sampel paling kecil yang diperlukan daripada ketiga-tiga taburan tersebut. Melalui kajian ini, didapati data sampel dari taburan Poisson dan Gamma memerlukan saiz sampel kurang daripada 30, manakala taburan Ekspensial memerlukan lebih daripada 30 untuk mencapai kenormalan.

Kata kunci: Normal; Saiz Sampel; Poisson; Ekspensial; Gamma

INTRODUCTION

Normality has been one of the important assumptions in inferential statistics. Inferential statistics is a branch of statistics which gives inferences on the population based on the sample collected. The samples obtained are used in hypothesis testing to make decision or conclusion on the population which the sample data represents.



Hypothesis testing is where normality becomes important. Hypothesis testing is defined as an assertion or deduction of the distribution of one or more random variables. If a hypothesis test strictly specifies the distribution, it is classified as simple hypothesis. Otherwise, it is called a composite hypothesis (Miller & Miller, 2014).

Statistical hypothesis testing is classified into two types which are parametric test and non-parametric test. A parametric test has a specific set of conditions for the test to be used. The conditions are usually assumptions such as the observations made from the population have to be: (1) independent of each other, (2) drawn from normally distributed populations, (3) have the same variances, and (4) the variables involved in the test must use the scale of interval or ratio (Bian, 2016). The test compares two sample means to determine whether the samples are significantly different. Without the assumptions, the test is not valid (Smith & Wells, 2006). Meanwhile, non-parametric tests do not consider any specific conditions of the parameters of the population like parametric tests. Strong measurement scales such as interval and ratio are not needed for non-parametric tests. Therefore, sample data with nominal and ordinal scale are tested using non-parametric methods (Bian, 2016).

Sample data that follows non-normal distribution are far more common in the real world compared to data that follows the normal distribution (Das & Imon, 2016). This situation gives a crucial issue for researchers as they cannot perform those parametric statistical tests. However, the sample data can approximate the normal distribution by increasing the sample size because as stated by the Central Limit Theorem (CLT) when the sample size from a non-normal distribution is large enough, the data will follow normal distribution (Kwak & Kim, 2017). In practice, statisticians and researchers have accepted the criterion of the sample size $n \geq 30$ to assume the distribution of sample mean approximated to normal distribution (Chang, Huang, & Wu, 2006; Chang, Wu, Ho, & Chen, 2008).

The Central Limit Theorem is the most fundamental theory in modern statistics. Parametric tests with assumptions that sample data from a population with fixed parameters cannot be modelled with a normal distribution without this theorem (Kwak & Kim, 2017). With the increase in sample size, the distribution varies less from the population mean, thus smaller variance. From this, it can be deduced that when sample size approaches infinity, the distribution approximates the normal distribution (Kwak & Kim, 2017). However, the question is how large is the sample size considered large enough for the sample data to approximate the normal distribution?

Ahad, Yaacob, Othman, Ng & Teoh (2011) questioned when CLT can be applied, and what exact sample size can be proven large enough to apply CLT. In the study, Chi-Squared distribution was the subject of study, where the mean, $m = 3$, variance, $v^2 = 6$, and the least sample size was 20. The researchers showed that the requirement for sample size to be at least 30 to achieve normality is not always true for all distributions.

There are many other non-normal distributions which real world data may follow. This study focuses on the three non-normal distributions which are Poisson distribution, Exponential distribution and Gamma distribution. These three distributions are chosen because they are very useful in many different fields. For example, Poisson distribution which counts the number of events occurred in a region of space or a period of time (Anacleto, 2018) can be used to model the waiting time and service time in a service industry and the chances of a game ties or a team wins in the sports industry (Tse, 2014). Poisson distribution is also chosen to represent



distributions which are slightly skewed. Therefore, the parameter for Poisson distribution is set as $\lambda = 3$. The distribution is plotted in Figure 1.

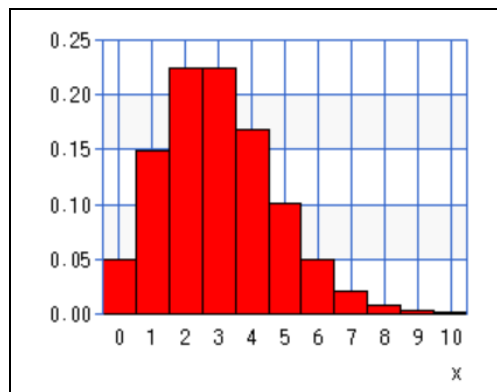


Figure 1. Poisson Distribution, $\lambda = 3$.

Exponential distribution is also used in the electronics field to model lifetime in life testing (Epstein, 1958). It is included in this study to represent distributions which are highly skewed. The parameter used is $\lambda = 1$. The distribution is highly skewed as shown in Figure 2.

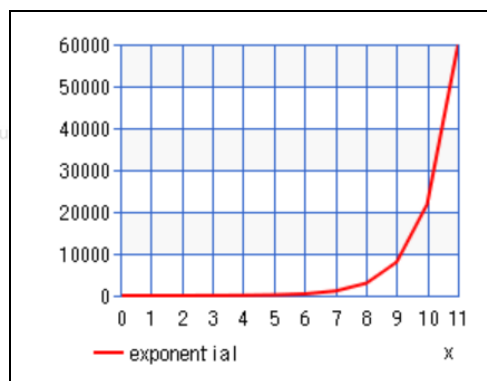


Figure 2. Exponential Distribution, $\lambda = 1$.

Gamma distribution is very useful survival analysis to conduct frailty modelling (Ngesa, & Orwa, 2019). In this study, it is used to represent moderately skewed distributions. The shape parameter was set as $\alpha = 2$ and scale parameter set as $\beta = 1$. The skewness of the graph is shown in Figure 3.

Although normality is needed to conduct parametric tests, most real world data do not achieve normality. While the data can achieve normality through the increment of sample size as stated in CLT when the sample size is large enough, the sample data approximates normal distribution, it is not known how much increment is enough for sample data following different distributions. Therefore, this study aims to identify the minimum sample sizes for Poisson distribution, Exponential distribution and Gamma distribution to achieve normality. These three distributions are chosen because these distributions are widely used in various industries as aforementioned.

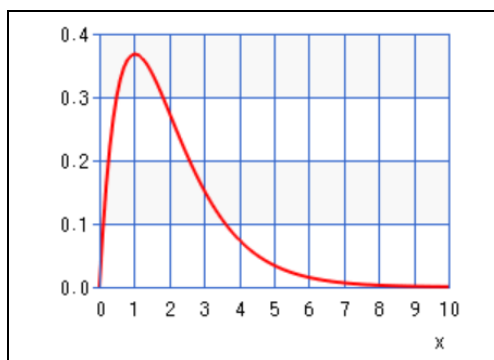


Figure 3. Gamma Distribution, $\alpha = 2$ and $\beta = 1$.

Normality tests verify the normality of the data set from a population before proceeding to statistical inference procedure (Laha, 2006). There are four most common normality tests namely, Shapiro-Wilk, Kolmogorov-Smirnov, Cramer-Von-Mises and Anderson-Darling tests. Furthermore, graphical methods are also used to test normality such as Quantile-Quantile plot (Q-Q plot).

Q-Q plot is an important diagnostic plot used to check the assumption of normality of a sample data set (Stine, 2016). The quantiles are values which divide a probability distribution into equal intervals, meaning that every interval of the total population is proportional. The aim of a Q-Q plot is to determine whether two data sets come from the same distribution. Since the focus of this study is on normality, the data set will be determined whether they come from a normally distributed with different sample sizes.

METHODOLOGY

The objective of this study is to determine the adequate sample sizes for Poisson, Exponential and Gamma distribution to achieve normality. To test the normality of the sample, normality tests that will be used are Shapiro-Wilk test, Kolmogorov-Smirnov test, Cramer-Von-Mises test, Anderson-Darling test and Q-Q plots.

Different sample sizes will be generated from each distribution, and normality tests will be applied on each sample. If normality is not achieved, the sample size will be increased and tested again until normality is reached for each distribution. To compare the results of the normality tests used, the smallest sample size for each distribution is compared.

The simulated data used in this study will be generated using SAS software. The data are the results of random number generated to follow the distribution set with fixed parameters such as mean and scale. The only change will be the number of sample data generated that will be used for normality tests.

The algorithm for the normality tests that will be carried out as follow:

- Step 1. Set the seed number for SAS to fix previous sample data generated and add new random data.
- Step 2. Set number of sample data that will be generated. The number of sample data will be set to 200.
- Step 3. Set the type of distribution, its parameters, and sample size. The first generated

- sample data is from Poisson distribution with the starting size of $n = 5$ because the minimum sample size to conduct Shapiro Wilk test is $n > 3$ (NCSS (n.d.)).
- Step 4. Calculate the sample mean for each sample.
 - Step 5. Test the 200 generated samples using the sample mean with the methods of Shapiro-Wilk, Kolmogorov-Smirnov, Cramer-Von-Mises, Anderson-Darling and Q-Q plot.
 - Step 6. If the generated sample data is found to be not normal, repeat the Step 4, using the increment of previous sample size.
 - Step 7. When the tests show that the generated sample data has achieved its normality, reduce the sample data by one until the sample data is detected to be not normal. Such step is to find the sufficient and exact sample size for Poisson distribution.
 - Step 8. Repeat Step 3 until Step 6 using the Exponential, $\lambda = 1$ and Gamma distributions, $\alpha = 2$, $\beta = 1$.
 - Step 9. Tabulate the results from the four numerical normality test methods in order to compare the results from four different test methods and Q-Q plots.

RESULTS AND DISCUSSION

Results for Poisson Distribution

The parameter value of Poisson distribution was set as $\lambda = 3$. The sample means of the 200 generated sample sizes were tested for normality using the four normality test methods as set by Statistical Analysis System (SAS) and Q-Q plots.

Table 1. Normality Tests for Poisson Distribution, $\alpha = 0.05$

| Poisson | <i>p</i> -value | | | | | |
|--------------------|-----------------|---------|---------|----------|----------|----------|
| | $n = 5$ | $n = 8$ | $n = 9$ | $n = 10$ | $n = 15$ | $n = 20$ |
| Shapiro-Wilk | 0.0084 | 0.0016 | 0.3316 | 0.5402 | 0.5814 | 0.2658 |
| Kolmogorov-Smirnov | <0.0100 | <0.0100 | 0.0507 | 0.0279 | >0.1500 | 0.03552 |
| Cramer-Von-Mises | <0.0050 | 0.0168 | 0.1763 | 0.1874 | >0.2500 | 0.1413 |
| Anderson-Darling | <0.0050 | 0.0094 | 0.2343 | >0.2500 | >0.2500 | 0.1535 |

Table 1 shows the results from the four normality test methods. It shows that when the sample size is 5, the sample data is found to be not normal. However, when the sample is increased by five to 10, the Shapiro-Wilk test, Cramer-Von-Mises test and Anderson Darling test show that the sample is normal at 5% significance level.

To obtain the exact sample size where the sample begins to achieve normality, the sample size is reduced by one until it is found to be not normal. From Table 1, it is known that when sample size is increased until 9 the sample data is still normal, but when reduced to 8, the sample data became not normal. Based on this result, it can be said that the minimum sample size for Poisson distribution (with $\lambda=3$) to achieve normality is 9 at 5% significance level.

From Table 1, three normality test methods which are Shapiro-Wilk, Cramer-Von-Mises and Anderson-Darling aligned with each other because the numerical results provide the same decision. However, Kolmogorov-Smirnov test does not show the same result as the other three methods.

In addition, the numerical values from Kolmogorov-Smirnov test method is not consistent and provide inconsistent decisions unlike the other three tests. Also, unlike other test methods, the Kolmogorov-Smirnov alternate between normal and not normal as the sample size increased.

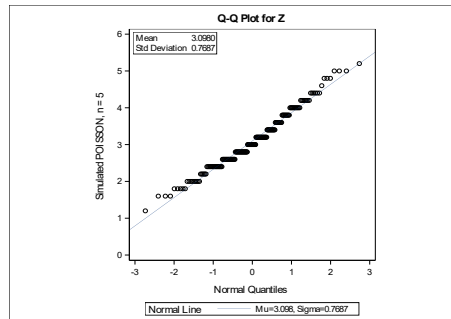


Figure 4. Q-Q Plot for Poisson, $n = 5$.

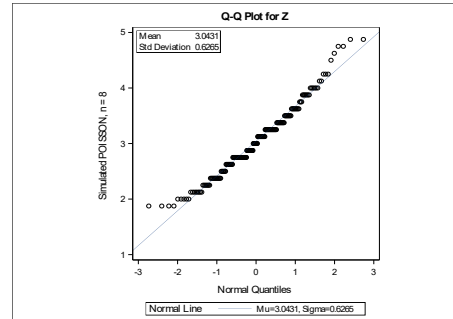


Figure 5. Q-Q Plot for Poisson, $n = 8$.

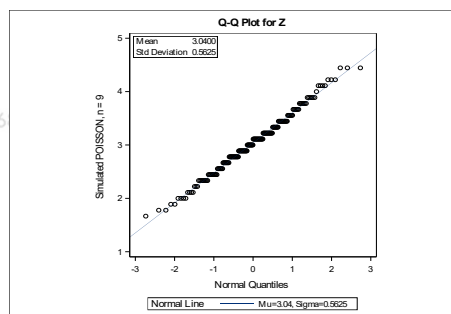


Figure 6. Q-Q Plot for Poisson, $n = 10$.

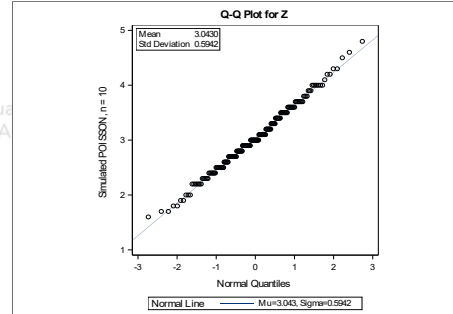


Figure 7. Q-Q Plot for Poisson, $n = 9$.

Figures 1 to 4 show the Q-Q plot results for sample data following Poisson distribution with sample sizes 5, 8, 9 and 10. Based on the result in the Q-Q plots, it can be observed that the samples are approaching normal as the sample sizes increased. For samples with sample sizes of 5 and 8, heavy tails are detected, and it can be said that both sample sizes are not normal. Whereas, the samples with sizes 9 and 10 are normal because the points are scattered closely around the diagonal. The Q-Q plots seem to align with the results of the normality tests except Kolmogorov-Smirnov test.

Results for Exponential Distribution

The parameter of Exponential distribution which is λ equal to 1 is set as default in SAS. The generated sample data are tested with four normality test methods as set in SAS and Q-Q plots.

Table 2. Normality Tests for Exponential Distribution, $\alpha = 0.05$

| Exponential | p -value | | | | |
|--------------------|------------|----------|----------|----------|----------|
| | $n = 5$ | $n = 10$ | $n = 15$ | $n = 20$ | $n = 25$ |
| Shapiro-Wilk | <0.0010 | <0.0001 | 0.0001 | 0.0065 | 0.0027 |
| Kolmogorov-Smirnov | 0.0133 | 0.0286 | <0.0100 | 0.0603 | 0.1304 |
| Cramer Von Mises | <0.0050 | 0.0199 | <0.0050 | 0.0327 | 0.0278 |
| Anderson-Darling | <0.0050 | 0.0073 | <0.0050 | 0.0181 | 0.0145 |

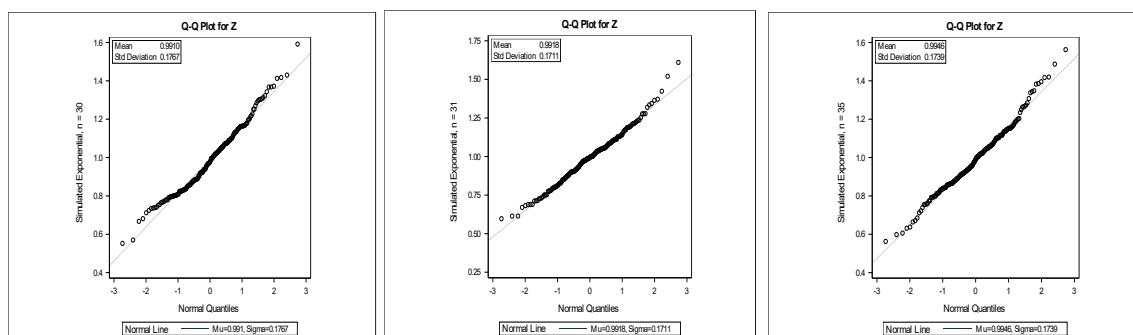
Table 3. Normality Tests for Exponential Distribution, $\alpha = 0.05$

| Exponential | p -value | | | | | |
|--------------------|------------|----------|----------|----------|----------|----------|
| | $n = 30$ | $n = 31$ | $n = 32$ | $n = 33$ | $n = 34$ | $n = 35$ |
| Shapiro-Wilk | 0.0119 | 0.1535 | 0.0516 | 0.0927 | 0.7408 | 0.0764 |
| Kolmogorov-Smirnov | 0.0284 | >0.1500 | 0.0914 | >0.1500 | >0.1500 | >0.1500 |
| Cramer-Von-Mises | 0.0293 | >0.2500 | 0.0928 | >0.2500 | >0.2500 | 0.1510 |
| Anderson-Darling | 0.0094 | >0.2500 | 0.0937 | >0.2500 | >0.2500 | 0.0857 |

Table 2 shows the results from four normality tests for sample data following Exponential distribution with sample sizes from 5 until 25, whereas Table 3 presents samples sizes from 30 until 35. Based on Table 2 and Table 3, it can be observed that as sample sizes increase, the p -values also increase for all four normality tests. This indicates that the sample data are approaching normality. Using the starter sample of size 5, which is then increased by 5 (i.e. $n = 5, 10, 15, \dots, i$), it has been found that the generated sample data from Exponential distribution achieved its normality at 5% significant level at sample of size 35 (i.e. $i = 35$).

The sample size is then reduced by one until $n = 31$. A line is drawn between sample size 30 and 31 because at sample size 31, the sample data achieves normality. This shows that the minimum sample size needed for Exponential distribution to achieve normality is 31 at 5% significance level.

Figures 5, 6, and 7 show the graphical test results for sample sizes 30, 31 and 35. Figures 5 to 7 show that only when the sample sizes are greater than 31, the sample data are approximately normal. For samples with sizes less than 31, the tails are heavy, thus it has denied the normality characteristics.

**Figure 8.** Exponential, $n = 30$. **Figure 9.** Exponential, $n = 31$ **Figure 10.** Exponential, $n = 35$.

Results for Gamma Distribution

The shape parameter of this distribution was set as $\alpha = 2$, $\beta = 1$ set in SAS. 200 samples with this parameter are generated and tested with the four normality test methods and Q-Q plot.

Table 4. Normality Tests for Gamma Distribution, $\alpha = 0.05$

| Gamma | p-value | | | | | | |
|--------------------|---------|----------|----------|----------|----------|----------|----------|
| | $n = 5$ | $n = 10$ | $n = 15$ | $n = 19$ | $n = 20$ | $n = 24$ | $n = 25$ |
| Shapiro-Wilk | 0.002 | 0.0006 | 0.0207 | 0.0121 | 0.0263 | 0.0103 | 0.0800 |
| Kolmogorov-Smirnov | 0.0232 | 0.0497 | 0.0271 | 0.0499 | >0.1500 | >0.1500 | >0.1500 |
| Cramer-Von-Mises | <0.0050 | 0.0194 | 0.0070 | 0.0121 | 0.2211 | >0.2500 | 0.2049 |
| Anderson-Darling | <0.0050 | 0.0071 | 0.0069 | 0.0090 | 0.1195 | >0.2500 | 0.1065 |

Table 4 shows the results from the four normality tests for the sample data following Gamma distribution with samples sizes from 5 until 25. Table 4 shows that when the sample size is 5, the sample data with shape parameter $\alpha = 2$ is denied of normality. When the sample size is increased to 20, three normality test methods which are Kolmogorov-Smirnov, Cramer-Von-Mises and Anderson-Darling show the normality at 5% significance level. Only Shapiro-Wilk test shows a different result. For Shapiro-Wilk test, the sample size has to increase to 25 for the sample to achieve normality at 5% significance level.

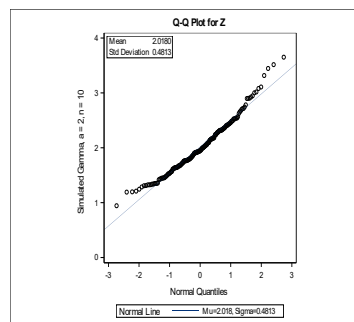


Figure 11. Gamma, $\alpha = 2, n = 10$

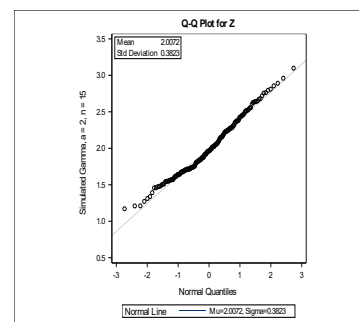


Figure 12. Gamma, $\alpha = 2, n = 15$

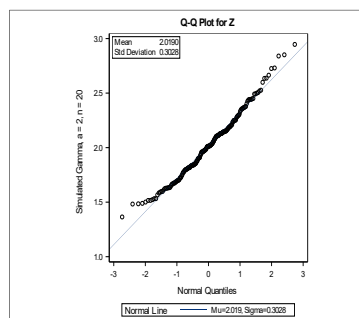


Figure 13. Gamma, $\alpha = 2, n = 20$

Figures 8 to 10 show the graphical test for normality for sample sizes 10, 15 and 20. When the sample size is increased to 15 as in Figure 9, the sample is almost normal. In addition, when the sample size further increases to 20 as in Figure 10, the skewness sample data is almost negligible. The results from the Q-Q plots do not align with the three normality tests because the three tests show normality at sample size 20. However, the results from the three normality tests should be given priority because graphical tests should be considered reference as it is subjective and cannot provide a formal conclusion (Yap & Sim, 2011).

Skewness and Sample Size

The degree of skewness of each distribution is as follows: (1) Poisson distribution -slightly skewed, (2) Gamma distribution – moderately skewed, and (3) Exponential distribution – highly skewed. From the findings above, the sample sizes needed for these three distributions to achieve normality also increase as its degree of skewness increases. This suggests that the more skewed the sample data is, the larger the sample size is needed to achieve normality.

CONCLUSIONS

This study aimed to determine how large the sample size is needed for three different non-normal distributions to achieve normality which are Poisson distribution, Exponential distribution and Gamma distribution. The sample size for each distribution was generated and increased until the sample data achieved normality. From this study, it is concluded that the number of sample size for a sample data to achieve normality is not necessarily 30. The minimum number of sample sizes for Poisson ($\lambda=3$) distribution, Exponential ($\lambda=1$) distribution and Gamma ($\alpha=2, \beta=1$) distribution are 9, 31 and 20 respectively. It also suggests that the sample size needed for a sample data to achieve normality is affected by its skewness. If the data is highly skewed, the larger the sample size needed.

The limitations of this study are the limited distribution shape and types of distributions. The shape the distribution affects the normality of its sample data. This study only focused on three shapes which are slightly skewed, moderately skewed and highly skewed which are subject to different interpretation by different researchers. Future research could study on the sample size needed to achieve normality as the shape of the distribution changes

ACKNOWLEDGMENT

The authors would like to express their gratitude to School of Quantitative Sciences, Universiti Utara Malaysia for all the facilities which contribute to the completion of this project.

REFERENCES

- Ahad, N. A., Yaacob, C. R., Othman, A. R., Ng, S. L., & Teoh, S. H. (2011). Central Limit Theorem in a Skewed Leptokurtic Distribution. *Jurnal Sains dan Matematik*, 26-33.
- Anacleto, O. (2018, February). Introduction to probability distributions. South Bridge, Edinburgh, The United Kingdom.
- Bian, H. (2016). *Non-parametric Tests*. Retrieved January 2, 2020, from <http://core.ecu.edu/ofe/statisticsresearch/Non-Parametric%20Tests.pdf>
- Chang, H. J., Huang, K. C., & Wu, C. H. (2006). Determination of Sample Size in Using Central Limit Theorem for Weibull Distribution. *Information and Management Sciences*, 17(3), 31-46.
- Chang, H. J., Wu, C. H., Ho, J. F., & Chen, P. Y. (2008). On Sample Size in Using Central Limit Theorem for Gamma Distribution. *Information and Management Sciences*, 19(1), 153-174.
- Das, K. R., & Imon, A. H. (2016). A Brief Review of Tests for Normality. *American Journal of Theoretical and Applied Statistics*, 5-12.
- Epstein, B. (1958). *The Exponential Distribution and Its Role in Life Testing*. Michigan: Department of Mathematics, Wayne State University.
- Kwak, S. G., & Kim, J. H. (2017). Central Limit Theorem: the cornerstone of modern statistics. *Statistical Round*, 144-156.
- Kiche, J., Ngesa, O. & Orwa, G. (2019). On Generalized Gamma Distribution and Its Application to Survival Data. *International Journal of Statistics and Probability*, 8(5), 85-102.
- Laha, A. K. (2006). *A Note on Tests of Normality*. Ahmedabad: Indian Institute of Management.
- Miller, I., & Miller, M. (2014). *John E. Freund's Mathematical Statistics with Applications*. Pearson Education Limited.
- NCSS.com. (n.d.). Normality Tests. Retrieved January 2, 2020, from https://ncss-wpengine.netdna-ssl.com/wp-content/themes/ncss/pdf/Procedures/NCSS/Normality_Tests.pdf
- Smith, Z. R., & Wells, C. S. (2006). *Central Limit Theorem and Sample Size*. New York: University of Massachusetts Amherst.
- Stine, R. A. (2016). Explaining Normal Quantile-Quantile Plots through. *The American Statistician*, 145-147.
- Tse, K.-K. (2014). Some Applications of the Poisson Process. *Applied Mathematics*, 3011-3017.
- Yap, B. W., & Sim, C. H. (2011). Comparisons of various types of normality tests. *Journal of Statistical Computation and Simulation*, 2141-2155.