



Preliminary Tests of Homogeneity - Type I Error Rates under Non-Normality



Tanweer Ul Islam*

Department of Economics, National University of Sciences and Technology (NUST), Pakistan

Submission: November 23, 2017; Published: May 22, 2018

*Corresponding author: Tanweer Ul Islam, Department of Economics, National University of Sciences and Technology (NUST), Pakistan;
Email: tanweer.ul.islam@gmail.com

Abstract

Many statistical procedures utilize preliminary tests to enhance the accuracy of the final inferences. Preliminary tests like Goldfeld-Quandt (GQ) and Levene-type tests are used to assess the assumption of equality of population variances with normality as the underlying distributional assumption. Such tests must be used with care as the final inferences are conditional on the performance of these tests at first stage. This study explores the size distortions of GQ and Levene-type tests under non-normality. The results do not warrant the use of GQ & Levene test under non-normality as the size distortions are as high as 88 & 48% for the respective statistics. However, the modified form of Levene test (BF-test) retains its size properties except for the multi-model alternatives with relatively big outliers.

Keywords: Size Distortions; Levene test; Equality of variances; Statistical procedures; Preliminary test; Chow-test; ANOVA; Homoscedasticity; Heteroskedastic model; Accuracy; Null hypothesis; Deficit data; Unequal sample sizes; Skewness; Kurtosis; Maximum size; Goldfeld-Quandt; Non-normality; Robust form; Multi-model distributions; Outliers

Abbreviations: GQ: Goldfeld-Quandt; DGP: Data Generating Process

Introduction

Many statistical procedures utilize preliminary test(s) to enhance the accuracy of the final inference. For example, in time series regression model, the Chow-test is widely used to test the presence of any structural change in the Data Generating Process (DGP), employs the Goldfeld-Quandt [1] test (GQ) as a preliminary test to assess the assumption of homogeneity of variances. The GQ- test is usually applied prior to the Chow test with normality as the underlying distributional assumption. Several other statistical procedures in the field of medical & social sciences, for example, One-way ANOVA makes use of the Levene's and the BF-tests as preliminary tests to assess equality of population variances. Such kind of preliminary tests are used in wide variety of applications, for example, public deficit data [2], regression analysis [3], audit pricing [4], capital structure [5], medicine [6], surgery [7], arthroplasty [8] & neuro imaging [9]. Furthermore, applications of the Levene-type tests have been surveyed in detail by Gastwirth et al. [10].

These preliminary tests must be used with care as the final inferences are conditional on the performance of these preliminary tests at first stage [11]. The GQ and the Levene's type tests assume the normality of the data while assessing the equal population variances. Although, some authors reassure robustness of modified Levene's type tests to normality but this study reemphasizes the use of diagnostic tests for normality for validating inferences made from regression models and from

other statistical procedures which utilize GQ & Levene's type preliminary tests. This study explores the impact of non-normality on the performance of the GQ & Levene's type tests. Since I plan to use numerical methods, the alternative (non-normal) space must be narrowed down to something sufficiently small to permit exploration by numerical methods. At the same time, the space should be large enough to provide a good approximation to the full space of alternatives – failing that, it should be large enough to approximate the distributions conventionally used in simulations studies to assess the performance of normality tests [12]; Pearson et al. [13]; Thadewald et al. [14], Zhang, et al. [15], Yazici, et al. [16], Romao et al. [17], Yap, et al. [18] and Bispo, et al. [19], Islam [20]. The distributions used as alternative space cover a wide range of real world applications in the field of Social Sciences, Genomics, Neuro Sciences and Bayesian Econometrics modelling. Type- I error rates for the GQ and Levene's type tests have been computed against the selected class of non-normal space to explore the impact of non-normality on their performance.

The preliminary tests

Some common statistical procedures like t-test, ANOVA & Chow test assume that variances of the populations from which k different samples are drawn are equal. The GQ & Levene's type tests assess this assumption. They test the null hypothesis that the population variances are homogeneous.

The goldfeld-quandt (GQ) test

For this test, it is assumed that the observations can be divided into two groups in such a way that under the hypothesis of homoscedasticity, the disturbance variances would be the same in the two groups, whereas under the alternative, the disturbance variances would differ systematically. The most favorable case for this would be the group-wise heteroskedastic model

$$y_i = x_i' \beta + \varepsilon_i,$$

Such that $\sigma_i^2 = \sigma^2 x_i^2$ for some variable x. To test explicitly, the suggested procedure is, by ranking the observations based on this x and dropping the central 'c' values, we can separate the observations into those with high and low variances. The test is applied by dividing the sample into two groups with and observations such that $n_1 + n_2 = n - c$. To obtain the statistically independent variances estimators, the regression is then estimated separately with the two sets of observations. The test statistic is

$$F[n_1 - K, n_2 - K] = \frac{e_1' e_1 / (n_1 - K)}{e_2' e_2 / (n_2 - K)}$$

Where, it is assumed that the disturbance variance is larger in the first sample. (If not, then reverse the subscripts.) Under the null hypothesis of homoscedasticity, this statistics has an F distribution with $n_1 - K$ & $n_2 - K$ degrees of freedom. A larger value than the standard F table value at the given level of significance leads to the rejection of the null hypothesis.

The levene-type tests

The Levene's type tests are used to assess the underlying assumption of homogeneity of variances. Statistical procedures

which typically assume equality of variances include analysis of variance (ANOVA) and t-tests. The Levene's test (1960) and the Brown-Forsythe (1974) test are often used as a preliminary test to validate the inferences drawn from the ANOVA and t-tests. The ANOVA is used to assess whether the k populations have a common mean For this, k samples $x_{i1}, x_{i2}, \dots, x_{in}$, of size n_i with respective means, μ_i and variances, $\sigma_i^2, i=1, \dots, k$ are drawn from each of k populations. To test the equality of means, the standard F-test assumes that the k populations has a common variance, σ^2 . To test the homogeneity of variances assumption, Levene proposed the following statistic.

$$F = \frac{(N - k) \sum_{i=1}^K n_i (Z_i - Z_{..})^2}{(k - 1) \sum_{i=1}^K \sum_{j=1}^{n_i} (Z_{ij} - Z_i)^2}$$

Where

$$Z_{ij} = \begin{cases} |Y_{ij} - \bar{Y}_i|, & \bar{Y}_i \text{ is a mean of } i^{\text{th}} \text{ group} \\ |Y_{ij} - \tilde{Y}_i|, & \tilde{Y}_i \text{ is a median of } i^{\text{th}} \text{ group} \end{cases}$$

$Z_{..}$ is the mean of all Z_{ij} and Z_i is the mean of the Z_{ij} for group i

The Levene's statistic is approximately F-distributed with $k - 1$ and $N - k$ degrees of freedom. The Brown-Forsythe test uses the median instead of mean. The Levene's type test based on median is recommended in the literature as these are robust statistics comparative to Levene's test against non-normality of data.

Simulation study & t- I error rates

Monte Carlo procedures are conducted to compute the type-I error rates for the GQ & Levene's type tests. These type- I error rates are obtained on the basis of 100,000 samples from the selected distributions (Table 1) for equal and unequal sizes of samples. Unequal sample sizes are chosen in 1:2, 1:3 & 1:4 ratios.

Table 1: Type-I Error Rates (equal samples& level of significance = 0.05).

Sr. No.	Distributions		Type-I Error Rates								
			n=25			n=50			n=75		
	D1	D2	GQ	BF	Levene	GQ	BF	Levene	GQ	BF	Levene
1	N (0,1)	N (0,1)	5.08%	4.01%	5.34%	5.07%	4.60%	5.22%	5.03%	4.71%	5.05%
2	Gamma (1,1)	Gamma (1,1)	27.27%	4.80%	14.50%	29.44%	4.80%	13.90%	30.41%	4.91%	13.82%
3	Beta (0.5,0.5)	Beta (0.5,0.5)	0.10%	2.28%	5.63%	0.04%	3.27%	5.29%	0.02%	3.62%	5.10%
4	Gamma (0.1250,1)	Gamma (0.1250,1)	64.72%	3.10%	30.62%	65.85%	3.99%	28.38%	66.39%	4.24%	27.22%
5	Gamma (0.25,1)	Gamma (0.25,1)	53.05%	4.32%	25.45%	54.36%	4.66%	23.92%	55.74%	4.80%	22.85%
6	Gamma (3,1)	Gamma (3,1)	13.86%	4.34%	8.82%	14.92%	4.74%	8.70%	15.33%	4.81%	8.72%
7	T (3)	T (3)	30.78%	3.78%	5.56%	36.85%	4.25%	5.16%	40.16%	4.44%	5.02%
8	T (5)	T (5)	17.00%	4.08%	5.63%	19.53%	4.47%	5.20%	20.53%	4.62%	5.06%
9	Chi2 (1.5)	Chi2 (1.5)	32.16%	4.85%	16.23%	34.43%	4.93%	15.86%	35.62%	4.86%	15.46%

10	Chi2 (2)	Chi2 (2)	27.51%	4.81%	14.29%	29.70%	4.86%	13.96%	30.59%	4.92%	13.67%
11	Chi2 (4)	Chi2 (4)	18.04%	4.45%	10.56%	19.28%	4.71%	10.27%	19.96%	4.86%	10.02%
12	Gamma (0.018,1)	Gamma (0.018,1)	88.20%	0.30%	47.68%	87.22%	0.99%	41.31%	86.86%	1.60%	39.05%
13	Gamma (0.0267,1)	Gamma (0.0267,1)	84.51%	0.62%	43.08%	83.65%	1.64%	38.50%	83.49%	2.22%	36.20%
14	LN (1,1)	LN (1,1)	50.70%	4.02%	19.17%	55.94%	4.21%	18.12%	58.87%	4.44%	18.26%
15	LN (1,2)	LN (1,2)	78.83%	2.12%	29.90%	79.84%	2.40%	28.89%	81.98%	2.59%	29.01%
16	Gamma (1,2)	Gamma (1,2)	27.69%	4.69%	14.19%	29.89%	4.84%	13.85%	30.43%	4.86%	13.77%
17	EV (1, 1)	EV (1,1)	14.98%	4.27%	8.30%	16.50%	4.75%	7.97%	17.01%	4.86%	8.01%
18	Logistic (0,1)	Logistic (0,1)	10.54%	4.21%	5.37%	11.22%	4.62%	5.22%	11.32%	4.70%	5.02%
19	Logistic (1,2)	Logistic (1,2)	10.43%	4.06%	5.42%	11.28%	4.48%	5.23%	11.45%	4.73%	5.23%
20	Laplace (0,1)	Laplace (0,1)	18.76%	4.35%	5.91%	19.71%	4.67%	5.53%	20.09%	4.83%	5.49%
21	NCX2 (1,3)	NCX2 (1,3)	19.32%	4.77%	11.93%	20.56%	4.96%	10.87%	21.13%	4.86%	11.48%
22	NCX2 (1,1)	NCX2 (1,1)	30.57%	4.98%	16.81%	31.92%	4.98%	15.93%	32.89%	5.02%	15.77%
23	Weibull (0.5,1)	Weibull (0.5,1)	27.59%	4.73%	14.41%	29.47%	4.84%	13.92%	30.41%	4.87%	13.86%
24	Weibull (1,2)	Weibull (1,2)	6.12%	4.18%	6.71%	6.21%	4.61%	6.39%	6.31%	4.72%	6.48%
25	Tukey (10)	Tukey (10)	20.71%	4.91%	8.09%	19.47%	5.03%	7.14%	18.98%	4.92%	6.80%

Performance of the GQ Test

In general, the GQ test performed poorly in terms of its size when evaluated over the entire range of selected alternative space for all sample sizes (Table 1 & 2). At 5% level of significance, the size of the GQ test goes up to 88% against highly skewed and heavy tailed alternatives both for the equal and unequal sample sizes. The size of the test is undervalued when the

alternative belongs to symmetric short tail class of distributions. The tenacious size distortions do not improve with the increase in sample size (Figure 1a & 1b). The size distortions are more than 10% and less 20% only for those alternatives where both skewness and kurtosis statistics are not far away from the normal distribution benchmark values; 0 & 3 respectively. Size distortions increase with the increase in value of either of the statistics- skewness and kurtosis.

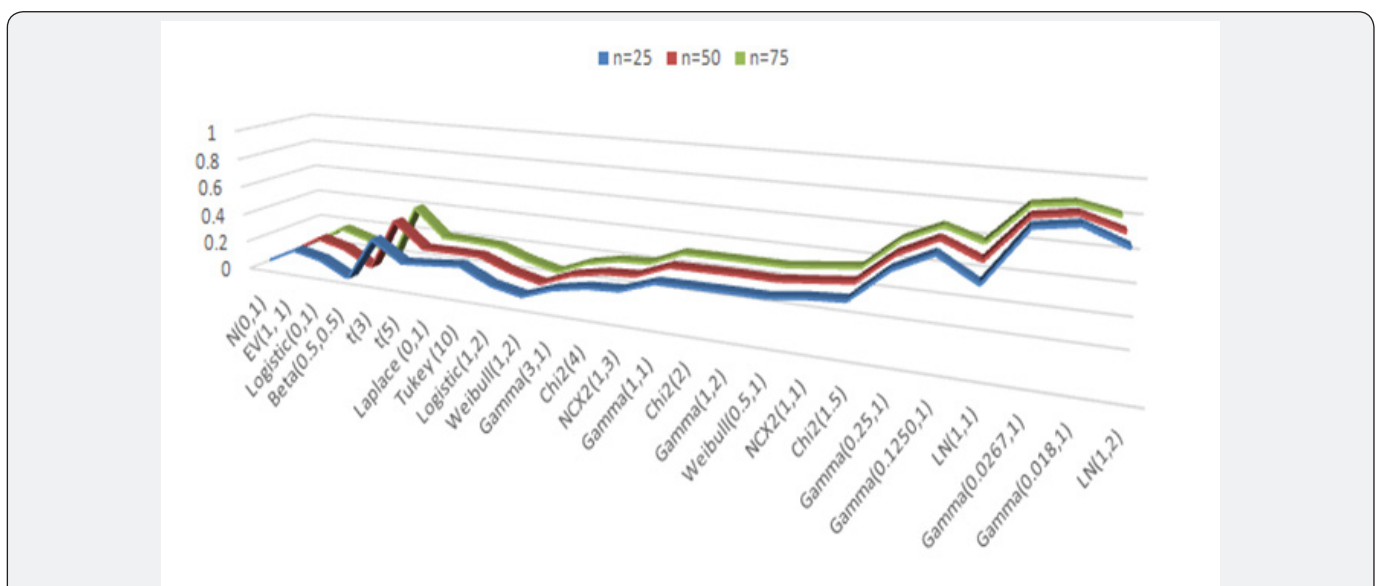


Figure 1(a): Size of GQ-test (equal samples).

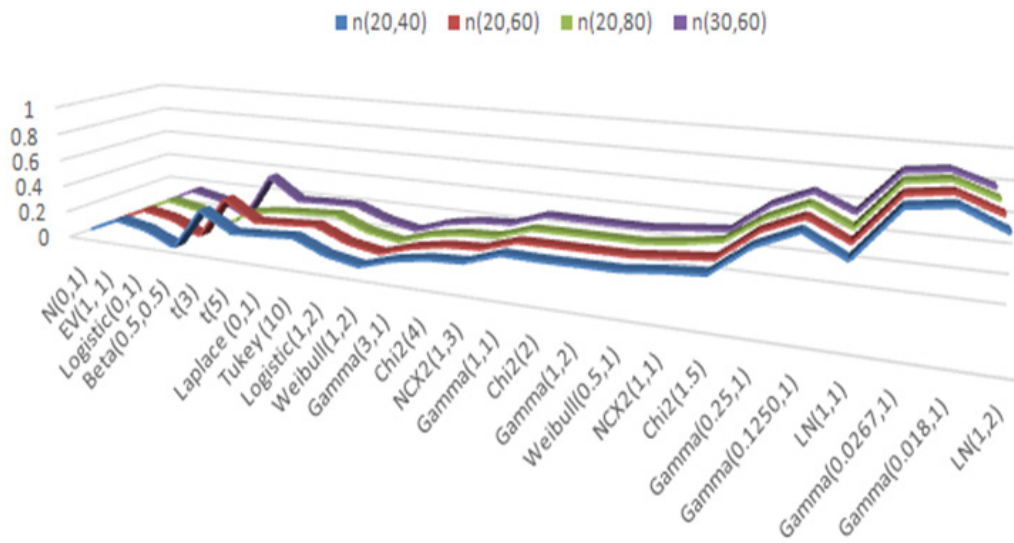


Figure 1(b): Size of GQ-test (unequal samples).

Performance of the levene-type tests

Robust form of Levene’s test proposed by Brown-Forsythe (1974), BF-test, performed exceptionally well in terms of size properties against all alternative distributions except for the multi-model distributions where the size of the test is underestimated (Table 1 & 2). The size of BF-test improves with the increase in sample size except for the cases where the alternative distribution contains few extreme outliers relative to rest of the sample data. The Levene’s test performance is not

satisfactory in comparison to its robust form (BF-test) which is based on median instead of arithmetic mean. The size of the test is more than 10% when the alternative space belongs to the group with skewness more than one and kurtosis more than five. Maximum size distortion reaches to as high as 48% for sample size of 25 (Figure 2a). There is a slight improvement in size distortions as the sample increases (Figure 2a & 2b). Mostly, the significant distortions are against the alternative distributions containing outliers with high values of skewness and kurtosis.

Table 2: Type-I Error Rates (unequal samples& level of significance = 0.05).

Sr. No.	Distributions		Type-I Error Rates											
			n1=20 & n2=40			n1=20 & n2=60			n1=20 & n2=80			n1=30 & n2=60		
	D1	D2	GQ	BF	Levene	GQ	BF	Levene	GQ	BF	Levene	GQ	BF	Levene
1	N (0,1)	N (0,1)	5.10%	4.30%	5.30%	5.10%	4.50%	5.20%	5.00%	4.60%	5.10%	5.00%	4.50%	5.10%
2	Gamma (1,1)	Gamma (1,1)	27.70%	4.70%	14.20%	27.40%	4.70%	13.60%	27.40%	4.60%	13.20%	29.10%	4.90%	14.00%
3	Beta (0.5,0.5)	Beta (0.5,0.5)	0.10%	3.60%	5.90%	0.10%	4.30%	6.30%	0.10%	5.20%	6.40%	0.00%	3.50%	5.60%
4	Gamma (0.1250,1)	Gamma (0.1250,1)	64.80%	3.70%	28.70%	65.00%	4.20%	27.20%	65.20%	4.40%	25.70%	65.60%	4.00%	28.00%
5	Gamma (0.25,1)	Gamma (0.25,1)	52.80%	4.30%	24.10%	53.30%	4.30%	23.00%	53.20%	4.30%	22.30%	54.20%	4.50%	23.50%
6	Gamma (3,1)	Gamma (3,1)	13.70%	4.60%	8.70%	13.60%	4.60%	8.40%	13.90%	4.60%	8.40%	14.50%	4.80%	8.60%
7	t(3)	t(3)	31.00%	4.00%	6.00%	32.00%	4.40%	6.10%	13.90%	4.50%	6.30%	35.10%	4.30%	5.70%
8	t(5)	t(5)	17.20%	4.30%	5.50%	17.30%	4.30%	5.40%	17.10%	4.50%	5.50%	18.50%	4.30%	5.30%
9	Chi2 (1.5)	Chi2 (1.5)	32.30%	4.80%	16.10%	32.50%	4.60%	15.50%	35.10%	4.50%	15.30%	33.80%	4.80%	15.60%

10	Chi2 (2)	Chi2 (2)	27.10%	4.70%	13.90%	27.60%	4.70%	13.70%	27.50%	4.60%	13.10%	28.70%	4.90%	13.90%
11	Chi2 (4)	Chi2 (4)	17.90%	4.60%	10.30%	18.00%	4.70%	10.10%	17.80%	4.60%	9.90%	18.80%	4.80%	10.20%
12	Gamma (0.018,1)	Gamma (0.018,1)	88.10%	2.20%	36.90%	87.80%	4.50%	28.40%	87.90%	6.00%	24.00%	87.40%	2.60%	36.60%
13	Gamma (0.0267,1)	Gamma (0.0267,1)	84.50%	2.60%	36.50%	84.30%	4.50%	29.90%	84.50%	5.50%	25.50%	84.10%	2.90%	35.50%
14	LN (1,1)	LN (1,1)	51.20%	4.00%	18.30%	52.30%	4.40%	17.30%	52.60%	4.40%	16.60%	54.30%	4.10%	18.20%
15	LN (1,2)	LN (1,2)	76.30%	3.00%	28.00%	77.60%	4.20%	24.40%	78.80%	5.40%	21.90%	78.70%	3.10%	27.70%
16	Gamma (1,2)	Gamma (1,2)	27.40%	4.90%	14.10%	27.50%	4.60%	13.60%	27.60%	4.60%	13.30%	28.60%	4.80%	14.00%
17	EV (1,1)	EV (1,1)	14.90%	4.50%	8.10%	14.80%	4.60%	8.00%	14.90%	4.60%	8.00%	15.90%	4.60%	7.90%
18	Logistic (0,1)	Logistic (0,1)	10.60%	4.30%	5.30%	10.70%	4.40%	5.30%	10.60%	4.70%	5.20%	11.10%	4.50%	5.30%
19	Logistic (1,2)	Logistic (1,2)	10.40%	4.30%	5.40%	10.40%	4.30%	5.20%	10.50%	4.50%	5.30%	11.10%	4.40%	5.10%
20	Laplace (0,1)	Laplace (0,1)	18.60%	4.50%	5.80%	18.70%	4.40%	5.40%	18.80%	4.50%	5.40%	19.40%	4.60%	5.50%
21	NCX2 (1,3)	NCX2 (1,3)	19.10%	4.80%	11.70%	19.30%	4.70%	11.40%	19.10%	4.70%	11.10%	20.10%	4.90%	11.80%
22	NCX2 (1,1)	NCX2 (1,1)	30.40%	4.90%	16.20%	30.30%	4.70%	16.20%	30.60%	4.60%	15.50%	31.50%	5.00%	15.90%
23	Weibull (0.5,1)	Weibull (0.5,1)	27.50%	4.70%	13.90%	27.70%	4.70%	13.50%	27.60%	4.70%	13.30%	28.80%	4.80%	13.70%
24	Weibull (1,2)	Weibull (1,2)	6.10%	4.40%	6.50%	6.20%	4.50%	6.50%	6.10%	4.60%	6.30%	6.20%	4.60%	6.40%
25	Tukey (10)	Tukey (10)	20.40%	4.90%	7.90%	20.40%	4.80%	7.40%	20.10%	4.80%	7.10%	20.00%	4.90%	7.40%

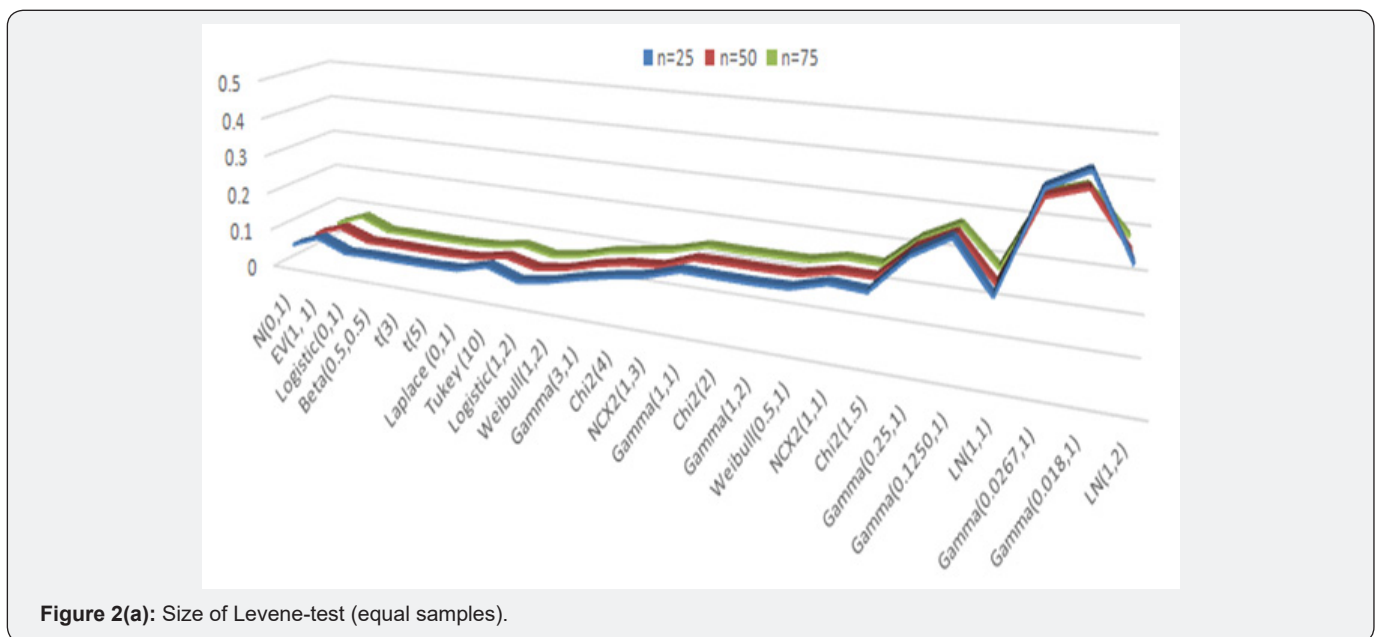


Figure 2(a): Size of Levene-test (equal samples).

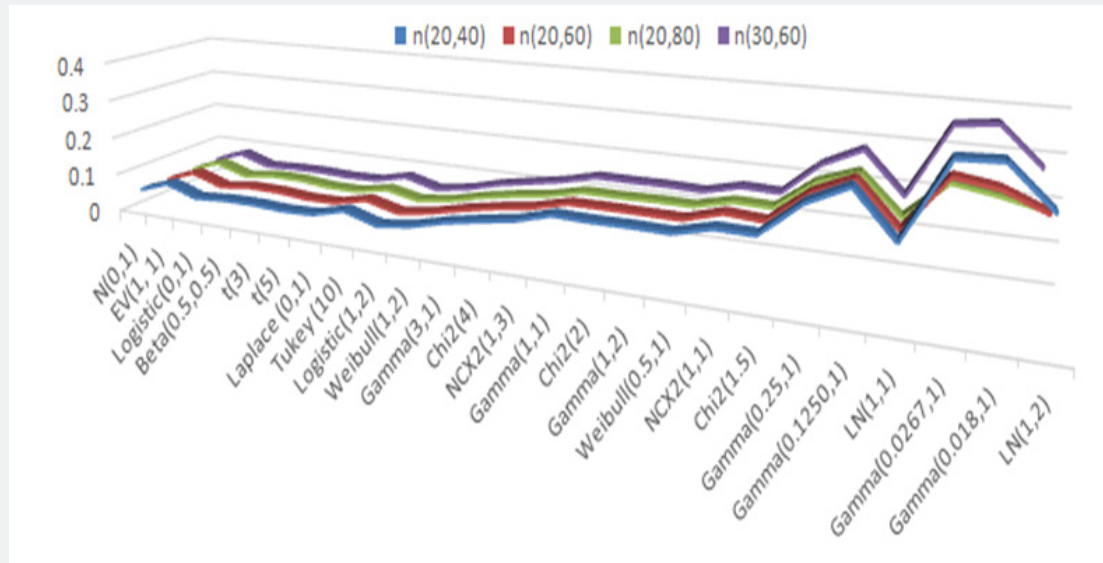


Figure 2(b): Size of Levene-test (equal samples).

Conclusion

Preliminary tests of homogeneity such as Goldfeld-Quandt (1965) and Levene-type tests are used to assess the assumption of homogeneity of variances which serves as the underlying assumption of many statistical procedures including Chow-test and one-way ANOVA. These preliminary tests assume the normality of data while assessing the equal population variances. Such kind of preliminary tests should be used with care as the final inferences are conditional on the performance of these tests at first stage. This study explores the impact of non-normality of the size distortions of these tests. At 5% level of significance, the size of the GQ test goes up to 88% against highly skewed and heavy tailed alternatives both for the equal and unequal sample sizes (Table 1 & 2). The size of the Levene test is more than 10% when the alternative space belongs to the group with skewness more than one and kurtosis more than five. Maximum size distortion reaches to as high as 48% for sample size of 25 (Figure 2a). Robust form of Levene’s test proposed by Brown-Forsythe, BF-test, performed exceptionally well in terms of size properties against all alternative distributions except for the multi-model distributions where the size of the test is underestimated (Table 1 & 2).

In general, both the statistics, GQ & Levene tests, suffer from severe size distortions when the alternatives belong to non-normal distributional space. However, the robust or modified form of Levene test (BF-test) perform well against the selected non-normal space except for few alternative distributions which are multi-model and contains big outliers. This study does not recommend the use of GQ & Levene test for assessing the assumption of equality of populations variances when the distribution is non-normal. Although, the modified form of Levene’s test (BF-test) retains its size properties however, the use is not recommended in case the distribution is multi-model and contains relatively big outliers.

References

1. Goldfeld SM, Quandt RE (1965) Some Tests for Homoscedasticity. *Journal of the American Statistical Association* 60(310): 539-547.
2. Correia Md, Neck R, Panagiotidis T, Richter C (2008) An empirical investigation of the sustainability of the public deficit in Portugal. *Springer-Verlag* 5(1): 209-223.
3. Zeileis A, Hothorn T (2002) Diagnostic Checking in Regression Relationships. *R News* 3(3): 7-10.
4. Francis JR, Simon DT (1987) A Test of Audit Pricing in the Small-Client Segment of the U. S. Audit Market. *The Accounting Review* 6(1): 145-157.
5. Tang CH, Jang S (2007) Revisit to the determinants of capital structure: A comparison between lodging firms and software firms. *Hospitality Management* 26: 175-187.
6. Banks ML, Roma PG, Folk JE (2011) Effects of the delta-opioid agonist SNC80 on the abuse liability of methadone in rhesus monkeys: a behavioural economic analysis. *Psychopharmacology* 16(3): 431-439.
7. Baiarda FU, Grobbelaar AO (2009) A comparison of one- versus two-stage surgery in an experimental model of functional muscle transfer with interposed nerve grafting. *J Plast Reconstr Aesthet Surg* 62(18):1042-1047.
8. Chawda M, Hucker P, Whitehouse SL, Crawford RW, English H, et al. (2009) Comparison of Cemented vs Uncemented Positioning Using an Imageless Navigation System. *J Arthroplasty* 24(8): 1170-113.
9. Grinband J, Wager TD, Ferrera VP, Hirsch J (2008) Detection of time-varying signals in event-related fMRI designs. *Neuroimage* 43(3): 509-520.
10. Gastwirth JL, Gel YR, Miao W (2009) The Impact of Levene’s Test of Equality of Variances on Statistical Theory and Practice. *Statistical Science* 24(3): 343-360.
11. Schucany WR, Ng HK (2006) Preliminary Goodness-of-fit Tests for Normality do not validate the One-Sample Student t. *Communication in Statistics- Theory and Methods* 35(12): 2275-2286.
12. Shapiro SS, Wilk MB, Chen HJ (1968) A Comparative Study of Various Tests for Normality. *Journal of the American Statistical Association* 63(324): 1343-1372.

13. Pearson ES, D Agostino RB, Bowman KO (1977) Tests for departure from normality: Comparison of power. *Biometrika* 64(02): 231-246.
14. Thadewald T, Büning H (2007) Jarque-Bera test and its competitors for testing normality- A power comparison. *Journal of Applied Statistics* 34(1):87-105.
15. Zhang J, Wu Y (2005) Likelihood-ratio tests for normality. *Computational Statistics & Data Analysis* 49: 709-721.
16. Yazici B, Yolacan S (2007) A comparison of various tests of normality. *Journal of Statistical Computation and Simulation* 77(02): 175-183.
17. Romao X, Delgado R, Costa A (2010). An empirical power comparison of univariate goodness-of-fit tests for normality. *Journal of Statistical Computation and Simulation* 80(5): 1-47.
18. Yap BW, Sim CH (2011) Comparisons of various types of normality tests. *Journal of Statistical Computation and Simulation* 81(12): 1-15.
19. Bispo R, Marques T A, Pestana D (2012) Statistical power of goodness-of-fit tests based on the empirical distribution function for type-I right-censored data. *Journal of Statistical Computation and Simulation* 82(2): 173-181.
20. Islam T U (2017) Stringency-based ranking of normality tests. *Communications in Statistics- Simulation and Computation* 46(1): 655-668.



This work is licensed under Creative Commons Attribution 4.0 License
DOI: [10.19080/BBOJ.2018.06.555699](https://doi.org/10.19080/BBOJ.2018.06.555699)

Your next submission with Juniper Publishers will reach you the below assets

- Quality Editorial service
- Swift Peer Review
- Reprints availability
- E-prints Service
- Manuscript Podcast for convenient understanding
- Global attainment for your research
- Manuscript accessibility in different formats
(Pdf, E-pub, Full Text, Audio)
- Unceasing customer service

Track the below URL for one-step submission

<https://juniperpublishers.com/online-submission.php>