

**A MODIFIED SHAPIRO – FRANCIA TEST FOR NORMALITY
OF DISTURBANCES IN LINEAR STATISTICAL MODEL**

Vijaya Kumar K*¹ and Balasiddamuni P²

¹Department of Statistics, S. G. S. Arts College, TTD, Tirupati, Andhra Pradesh, India.

²Department of Statistics, Sri Venkateswara University, Tirupati, Andhra Pradesh, India.

(Received on: 07-03-14; Revised & Accepted on: 24-03-14)

ABSTRACT

In the present study, a procedure has been discussed for testing normality of disturbances in linear statistical model. A modified version of Shapiro-Francia Test for normality of disturbances based on internally studentized residuals has been proposed as large sample test.

1. INTRODUCTION

In many empirical studies using the linear regression model, it is standard practice to assume that the random disturbances have zero means, and are normally, independently and identically distributed. Since, failure of these assumptions to hold can lead to biased and inconsistent estimators; and incorrect inferences, it is sensible to test these assumptions rather than to presume that they are correct.

Tests of normality are statistical inference procedures designed to test that the underlying distribution of a random disturbance variable is normally distributed. There is a long history of these tests, and there are a plethora of them available for use. The recent interest in goodness-of-fit tests for normality has arisen from the exciting work of Shapiro and Wilk (1965).

The problem of non-normal disturbances is one of the serious problems in the theory of Econometrics. The presence of non-normal disturbances disturbs the optimal properties of the Ordinary Least Squares (OLS) estimators of the parameters of a linear statistical model. Hence, there is a necessity of detecting the problem of non-normal disturbances so that the various techniques can be applied to avoid the presence of non-normal disturbances in the linear statistical model.

In the present study, an attempt is made to propose a modified version of Shapiro-Francia test statistic based on Internally studentized residuals, for testing normality of disturbances in linear statistical models.

2. REVIEW OF THE LITERATURE

The main contribution relating to the problem of testing normality of the disturbances in the linear statistical models, has been made by Geary (1935, 1947), Shapiro and Wilk (1965), Shapiro, Wilk and Chen (1968), Shapiro and Francia (1972), Huang and Bolch (1974), Filliben (1975), Weisberg and Bingham (1975), Locke and Spurrier (1976), Gastwirth and Owens (1977), Spiegelhather (1977), Pearson, D'A gostino and Bowman (1977), White and Mac Donald (1980), Pierce and Gray (1982), Royston (1982), La Riccia (1986), Jarque and Bera (1987), Bonett and Woodward (1990), Kalirajan and Jayasuriya (1991), Al-shiha and Yang (1996), Markatou and Manos (1996) and Rahman and Govindarajulu (1997).

Most of the tests for normality discussed in the literature have used either the Ordinary Least Squares (OLS) residual vector (e) or the best linear unbiased scalar (BLUS) residuals vector (\tilde{e}) or the recursive residual vector ($\tilde{\tilde{e}}$).

Corresponding author: Vijaya Kumar K*¹

¹Department of Statistics, S. G. S. Arts College, TTD, Tirupati, Andhra Pradesh, India.

E-mail: vijayneeraja73@gmail.com

SHAPIRO AND FRANCIA W¹ TEST FOR NORMALITY

Shapiro and Wilk (1965) suggested w test statistic for normality as

$$W = \frac{\left(\sum_{i=1}^n a_i y_i\right)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \tag{2.1}$$

Where, $a^1 = [a_1 a_2 \dots a_n] = \frac{m^1 v^{-1}}{[m^1 v^{-1} v^{-1} m^1]^{1/2}}$ (2.2)

$V = ((V_{ij}))$ is (n x n) covariance matrix;

$m^1 = [m_1 m_2 \dots m_n]$ = vector of expected values of standard normal order statistics.

$y^1 = [y_1 y_2 \dots y_n]$ = vector of ordered random observations

The elements of V are given only for sample sizes upto 20 and approximations were developed by Shapiro and Wilk (1965) for calculating coefficients {a_i} for using the w test statistic upto the sample size 50.

The overcome the main drawback of Shapiro and Wilk (1965) W test statistic for large samples, the observation {y_i} may be treated as if they were independent and the identify matrix I_n can be substituted for V⁻¹ in the estimation of the slope of the regression line

$$y_i = \mu + \sigma X_i, i = 1, 2, 3, \dots, n. \tag{2.3}$$

Consequently, Shapiro and Francia (1972) suggested an approximate W¹ test statistic for normality as

$$w^1 = \frac{\left(\sum_{i=1}^n b_i y_i\right)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \tag{2.4}$$

Where, $b^1 = [b_1 b_2 \dots b_n] = \frac{m^1}{\sqrt{m^1 m}}$ (2.5)

Values of ‘m’ are given for different sample sizes by Harter (1961). The null distribution of W¹ statistic was approximated by an empirical sampling study. Shapiro and Francia (1972) conducted an empirical sampling study to compare the sensitivities of the W and W¹ test statistics in detecting non-normality. The empirical percentage points of W¹ test statistic for different sample sizes were tabulated and concluded that in general, the W¹ test appears to be more sensitive than the W test when the alternative distribution was continuous and symmetric with high fourth moment (as compared to the normal distribution); when it was near normal and when it was discrete and skewed. The two tests W and W¹ appeared to be equivalent for alternative distributions which were continuous and asymmetric with high fourth moment and discrete and symmetric. The W test is superior to the W¹ test for other alternative distributions. Overall, the differentials in the power were small.

3. STUDENTIZATION OF RESIDUALS

By Studentization, the OLS residuals can be transformed to have a selected covariance structure and the null distributions of the transformed residuals are independent of the scale parameters.

Since, $Var(e) = \sigma^2(I - V)$ or $Var(e_i) = \sigma^2(I - v_{ii}) \forall i = 1, 2, 3, \dots, n.$, the OLS residuals have the distribution which is scale dependent. Margolin (1977) defined the studentized residuals as the division of a scale dependent

statistic say U, by a scale estimate T, so that the ratio $S = \frac{U}{T}$ has a distribution that is free of the nuisance scale parameters.

David (1981) distinguished two type of studentization of residuals namely

- (i) Internal studentization, where U and T are dependent and
- (ii) External studentization, where U and T are independent.

(A) INTERNAL STUDENTIZATION OF RESIDUALS

Under internal studentization, the ‘Internally studentized Residuals’ are

$$e_i^+ = \frac{e_i}{\hat{\sigma} \sqrt{1-v_{ii}}}, i = 1, 2, \dots, n. \tag{3.1}$$

Where $\hat{\sigma} = \sqrt{\frac{\sum e_i^2}{n-k}}$ (3.2)

Beckman and Trussell (1974) have proved that, $e_i^{+2}/n-k$ follows a Beta distribution with parameters $\frac{1}{2}$ and $\frac{n-k-1}{2}$. The probability density function of e_i^+ is given by

$$f(e_i^+) = \frac{\Gamma\left(q + \frac{1}{2}\right)}{\Gamma(q)\Gamma\left(\frac{1}{2}\right)\sqrt{n-k}} \left(1 - \frac{e_i^{+2}}{n-k}\right)^{q-1}; \tag{3.3}$$

and $|e_i^+| \leq \sqrt{n-k}$

Here, $q = \frac{n-k-1}{2}$

It follows that $E[e_i^+] = 0, \forall i = 1, 2, 3, \dots, n$

$Var(e_i^+) = 1$; and

$$cov(e_i^+, e_j^+) = \frac{-v_{ij}}{\sqrt{(1-v_{ii})(1-v_{jj})}}, \forall i \neq j \tag{3.4}$$

(B) EXTERNAL STUDENTIZATION

By external studentization, the ‘Externally Studentized Residuals’ are defined by

$$e_i^{++} = \frac{e_i}{\hat{\sigma}_{(i)} \sqrt{1-v_{ii}}} \tag{3.5}$$

Where, $\hat{\sigma}_{(i)}^2 = \frac{(n-k)\hat{\sigma}^2 - \frac{e_i^2}{1-v_{ii}}}{n-k-1}$

or $\hat{\sigma}_{(i)}^2 = \hat{\sigma}^2 \left[\frac{n-k - \frac{e_i^{+2}}{n-k-1}}{n-k-1} \right]$ (3.6)

Here, $\hat{\sigma}_{(i)}^2$ is the residual mean square computer without ith observation on both Y and X's. Under normality, $\hat{\sigma}_{(i)}^2$ and e_i are independent.

The probability density function of e_i^{++} is the Student's t-distribution with (n-k-1) degree of freedom.

A relationship between e_i^{++} and e_i^+ is given by

$$e_i^{++} = e_i^+ \left[\frac{n-k-1}{n-k-e_i^{+2}} \right]^{1/2} \tag{3.7}$$

This shows that e_i^{++2} is a monotonic transformation of e_i^{+2} .

Studentized residuals provide way to examine the information in the residuals, both because they have equal variances and because they are easily related to the t-distribution in many situations.

4. A MODIFIED SHAPIRO – FRANCIA TEST FOR NORMALITY OF DISTURBANCES

Consider the classical linear regression model,

$$Y_{n \times 1} = X_{n \times k} \beta_{k \times 1} + \epsilon_{n \times 1} \tag{4.1}$$

Where, X is a known nonstochastic (nxk) matrix of rank k; β is a (kx1) vector of unknown parameters; Y is the (nx1) vector of random observation; and ϵ is an (nx1) random vector whose elements are assumed to be normally distributed with mean zero and constant variance σ^2 .

Since, the disturbance vector ϵ is unobservable, the various tests such as tests for specification errors in the model, tests for normality etc., generally rely on sample residuals such as OLS or BLUS or Recursive or studentized or predicted residuals. However, OLS residuals are considered to be inferior to other residuals by many econometricians, since the OLS residuals vector has a singular normal distribution and the elements of the vector are not independently and identically distributed.

The OLS residual vector e associated with the true disturbance vector ϵ is given by

$$e = M \epsilon \tag{4.2}$$

where, $M = [I - X(X'X)^{-1}X']$ is an (nxn) symmetric idempotent matrix with rank (n-k). Thus e has a singular normal distribution with zero mean vector and covariance matrix is given by $\sigma^2 M$.

The application of the tests to the disturbances offers a set of benchmark results for gauging the performance of the test as applied to identically and independently distributed random variables.

The recent interest in goodness of fit tests for normality has arisen from the exciting work of Shapiro and Wilk (1965),

The (n-k) x 1 matrix of BLUS residuals is given by

$$\tilde{e} = A' \epsilon \tag{4.3}$$

Where A' is (n-k) xn such that $A'A = I_{n-k}$ and

$$E[\tilde{e}\tilde{e}'] = \sigma^2 I_{n-k}. \text{ Let } A = ((a_{ij}))$$

Thus, the BLUS residuals are linear combinations of the disturbance and a set of coefficients such that

$$\left. \begin{aligned} \sum_{r=1}^n a_{ri} a_{rj} &= 0 \\ \sum_{r=1}^n (a_{ri} a_{rj})^2 &\neq 0 \end{aligned} \right\} \forall i \neq j \tag{4.4}$$

By applying the following theorem given by Kingman and Grabill (1970):

Let Y_1, Y_2, \dots, Y_n be i.i.d. random variables with cdf denoted by F and let n be a positive integer. Suppose that L_1 and L_2 are defined by

$$L_1 = \sum_{i=1}^n a_i Y_i ; \quad L_2 = \sum_{i=1}^n b_i Y_i \quad (4.5)$$

Where a_i and b_i are constants such that

$$\sum_{i=1}^n a_i b_i = 0 \quad \text{and} \quad \sum_{i=1}^n (a_i b_i)^2 \neq 0$$

Then F is the cdf of a normally distributed random variable if and only if L_1 and L_2 are independent”.

Using this theorem, we find that if the disturbances are independently and identically distributed, the BLUS residuals can be independent if and only if the disturbances are normally distributed. BLUS residual vector $\tilde{\mathbf{e}}$ has a nonsingular normal distribution with a covariance matrix given by $\sigma^2 \mathbf{I}_{n-k}$.

Suppose $\tilde{e}_{(1)} \leq \tilde{e}_{(2)} \leq \dots \leq \tilde{e}_{(n-k)}$ denote an ordered random sample of $(n-k)$ BLUS residuals derived from n OLS residuals.

Let $\mathbf{m}^l = (m_1, m_2, \dots, m_{n-k})$ be the vector of expected values of ordered BLUS residuals and

$$\mathbf{b}^l = \frac{\mathbf{m}^l}{\sqrt{\mathbf{m}^l \mathbf{m}^l}}, \text{ where } \mathbf{b}^l = (b_1, b_2, \dots, b_{n-k}) \quad (4.6)$$

By using Shapiro and Francia (1972) approximate Analysis of variance test of normality, a modified test statistic for normality is given by

$$W^* = \frac{\left[\sum_{i=1}^{n-k} b_i \tilde{e}_{(i)} \right]^2}{\sum_{i=1}^{n-k} [\tilde{e}_{(i)} - \bar{\tilde{e}}]^2} \quad (4.7)$$

$$\text{where, } \bar{\tilde{e}} = \frac{\sum_{i=1}^{n-k} \tilde{e}_i}{n-k} \quad (4.8)$$

This test can be performed by using the critical values tabulated by Shapiro and Francia (1972)

Two more modification of test statistic, using the approximations suggested by Weisburg and Binham (1975) and Blom (1956) are given by:

$$(a) \quad \mathbf{b}^{*l} = \frac{\mathbf{m}^{*l}}{\sqrt{\mathbf{m}^{*l} \mathbf{m}^{*l}}} \quad (4.9)$$

$$\text{where } m^* = \phi^{-1} \left(\frac{i - 3/8}{n + 1/4} \right), \quad i=1, 2, \dots, (n-k) \quad (4.10)$$

and ϕ^{-1} denotes the inverse of the standard normal cumulative distribution function. Then the test statistic is given by

$$W^{**} = \frac{\left[\sum_{i=1}^{n-k} b_i^* \tilde{e}_{(i)} \right]^2}{\sum_{i=1}^{n-k} [\tilde{e}_{(i)} - \bar{\tilde{e}}]^2} \quad (4.11)$$

$$(b) \mathbf{b}^{**i} = \frac{\mathbf{m}^{**i}}{\sqrt{\mathbf{m}^{**i} \mathbf{m}^{**}}}, \quad (4.12)$$

where, $\mathbf{m}^{**} = \phi^{-1} (P_i), \quad i = 1, 2, \dots, (n - k)$ (4.13)

Here, $P_i = \frac{i}{n + 1}, \quad i = 1, 2, 3, \dots, (n - k)$

$$W^{***} = \frac{\left[\sum_{i=1}^{n-k} \mathbf{b}_i^{**} \tilde{\mathbf{e}}_{(i)} \right]^2}{\sum_{i=1}^{n-k} \left[\tilde{\mathbf{e}}_{(i)} - \bar{\tilde{\mathbf{e}}} \right]^2} \quad (4.14)$$

The main advantages of these test statistics that they remove the need to have tables of weights for computation of the test statistics.

The proposed modified test statistic will provide Goodness of fit test for normality for all sample sizes without the assumption of zero-correlation among the ordered residuals.

It should be noted that, although the above tests assume that the residuals are independent, both Haung and Bolch (1974) and Ramsey (1974) reported from their Monte Carlo studies, where the correlated OLS residual vector led to a more powerful test than that obtained using the BLUS residual vector.

Now, we propose a modified Shapiro-Francia large sample test based on Internally studentized residual vector \mathbf{e}^+ as follows.

$$W^+ = \frac{\left[\sum_{i=1}^n \mathbf{b}_i^+ \mathbf{e}_{(i)}^+ \right]^2}{\sum_{i=1}^n \mathbf{e}_{(i)}^{+2}} \quad (4.15)$$

where, $\mathbf{b}^+ = \frac{\mathbf{m}^+ \Sigma^{-1}}{\sqrt{\mathbf{m}^+ \Sigma^{-1} \mathbf{m}^+}} ; \quad (4.16)$

$\mathbf{m}^+ = (\mathbf{m}_1^+, \mathbf{m}_2^+, \dots, \mathbf{m}_n^+) =$ Vector of expected values of ordered Internally studentized residuals;

$\mathbf{e}_{(1)}^+ \leq \mathbf{e}_{(2)}^+ \leq \dots \leq \mathbf{e}_{(n)}^+$ are the ordered studentized residuals; $\sum ((\sigma_{ij}))$ is the nxn covariance matrix of ordered Internally studentized residuals; and

$$\begin{aligned} \sigma_{ij} &= 1 \quad \forall i = j = 1, 2, \dots, n \\ &= \frac{-v_{ij}}{\sqrt{(1 - v_{ii})(1 - v_{jj})}} \quad \forall i \neq j = 1, 2, \dots, n \end{aligned} \quad (4.17)$$

Here, $\mathbf{V} = ((v_{ij})) = \mathbf{X}(\mathbf{X}^t \mathbf{X})^{-1} \mathbf{X}^t =$ Hat matrix (4.18)

Under the following assumptions given by white and Mac Donald (1980) for the large sample behavior of the modified tests for nonnormality is the linear regression model, it can be shown that the proposed modified Shapiro-Francia Large sample test statistic W^+ is consistent estimator of the true statistic W .

Assumption (i) The linear regression model is

$$Y_i = X_i^T \beta + \epsilon_i, \quad i = 1, 2, \dots, n \quad (4.19)$$

Where $|X_i|$ is a sequence of uniformly bounded fixed (kx1) vector such that

$$\frac{X^T X}{n} \rightarrow Q_{xx}, \text{ a positive definite matrix} \quad (4.20)$$

$|\epsilon_i|$ is a sequence of i.i.d. random variables

$$\text{with } \left. \begin{aligned} E[\epsilon_i] &= 0, \quad \forall i=1, 2, \dots, n \\ E[\epsilon_i^2] &= \sigma^2, \quad 0 < \sigma^2 < \infty \end{aligned} \right\} \quad (4.21)$$

β is a (kx1) vector of parameters

This assumption is sufficient to ensure

$$\hat{\beta} \xrightarrow{a.s} \beta \quad (4.22)$$

Assumption (ii) (a) The density of ϵ_i , say $f = dF$ is uniformly continuous, positive on the interval of support and bounded;

(b) Define

$$\left. \begin{aligned} X_{ip}^+ &= \text{Max}\{X_{ip}, 0\} \quad \text{and} \\ X_{ip}^- &= X_{ip}^+ - X_{ip}; \quad i = 1, 2, \dots, n \\ & \quad p=1, 2, \dots, k \end{aligned} \right\} \quad (4.23)$$

$$\text{Max}\{X_{ip}^+ : 1 \leq i \leq n\} = 0 \left\{ \left[\sum_{i=1}^n (X_{ip}^+)^2 \right]^{1/2} \right\} \text{ and}$$

$$\text{Max}\{X_{ip}^- : 1 \leq i \leq n\} = 0 \left\{ \left[\sum_{i=1}^n (X_{ip}^-)^2 \right]^{1/2} \right\}, \quad p = 1, 2, \dots, k$$

These two assumptions imply the necessary and sufficient conditions for the asymptotic normality of $\sqrt{n}(\hat{\beta} - \beta)$.

Under the above two conditions, it can be shown that

$$\left| W^+ - W \right| \xrightarrow{p} 0 \quad (4.24)$$

Where W is Shapiro-Francia test statistic in terms of disturbances ϵ_i 's.

Small values of W^+ are significant, that indicates nonnormality of disturbances. A more precise significance level may be associated with an observed W^+ value by using the critical values given by Shapiro and Francia (1972).

5. CONCLUSIONS

In the present study a modified version of Shapiro-Francia test for normality of disturbances in a linear statistical model has been discussed by using Internally Studentized Residuals. This proposed test statistic is a consistent estimator of the true test statistic under the certain assumptions given by White and Mac Donald (1980). This kind of study can be further extended by deriving power functions of the tests, for the comparative study of the various tests of normality of disturbances.

REFERENCES

1. Al-shiha, A.A. and Yang, S.S. (1996), "An Improved Approximation for the Shapiro-Wilk Test Statistic for Normality", Pakistan Journal of Statistics, 12,215-230.
2. Beckman, R. and Trussell, H. (1974), "The Distribution of an Arbitrary Studentized Residual and the Effects of updating in Multiple Regression", JASA, 69, 199-201.
3. Blom, G. (1956), "Statistical Estimates and Transformed Beta Variables", New York: John Wiley.
4. Bonett, D.G. and Woodward, J.A. (1990), "Testing Residual Normality in the ANOVA Model", Journal of Applied Statistics, 17, 383-387
5. David, H.A. (1981), Order Statistics, 2nd Edition Wiley, New York.
6. Filliben, J.J. (1975), "The Probability plot Correlation Co-efficient Test for Normality", Technometrics, 17, 111-117.
7. Gastwirth, J.L. and Owens, M.E.B.(1977), "On Classical Tests of Normality", Biometrika, 64, 135-139
8. Geary, R.C. (1935), "Note on the Correlation Between p_z and W^1 ", Biometrika, 27, 353-355.
9. Geary, R.C. (1947), "Testing for Normality", Biometrika, 34, 209-242.
10. Harter, H.L. (1961), "Expected Values of Normal Order Statistics", Biometrika, 48, 151-165.
11. Huang, C.J. and Bolch, B.W. (1974), "On the Testing of Regression Disturbances for Normality", JASA, 69, 330-335.
12. Jarque, C.M. and Bera, A.K. (1987), "A Test for Normality of Observations and Regression Residuals", International Statistical Review, 55, 163-172.
13. Kalirajan, K.P. and JayaSuriya, S.K.W. (1991), "Simultaneous Testing of Regression Disturbances for Heteroscedasticity and Non-Normality", Journal of Applied Statistics, 18, 307-312.
14. Kingman, A. and Gray Bill, E.A. (1970), "A Nonlinear Characterization of the Normal Distributon", The Annals of Mathematical Statistics, 41, 1889-95.
15. La Riccia, V.N. (1986), "Optimal Goodness-of-Fit Tests for Normality Against Skewness and Kurtosis Alternatives", Journal of Statistical Planning and Inference, 13, 67-69.
16. Locke, C. and Spurrier, J.D. (1976), "The Use of U-Statistics for Testing Normality Against Non-symmetric Alternatives", Biometrika, 63, 143-147.
17. Margolin, B.H. (1977), "The Distribution of Internally Studentized Statistics Via Laplace Transform Inversion", Biometrika, 64, 573-582.
18. Markatou, M. and Manos, G. (1996), "Robust Tests in Nonlinear Regression Models", Journal of Statistical Planning and Inference, 55, 205-217.
19. Pearson, E.S., D'Agostino, R.B. and Bowman, K.O. (1977), "Tests for Departure for Normality: Comparision of Powers", Biometrika, 64, 231-246.
20. Pierce, D.A. and Gray, R.J. (1982), "Testing Normality of Errors in Regression Models", Biometrika, 69, 233-236
21. Rahman, M.M. and Govindarajulu, Z. (1997), "A Modification of the Test of Shapiro and Wilk for Normality", Journal of Applied Statistics, 24, 219-235.
22. Ramsey, J.B. (1974), "Classical Model Selection Through Specification Error Tests", in P. Zarembka, ed., Frontiers in Econometrics, New York: Academic.
23. Royston, J.P. (1982), "An Extension of Shapiro and Wilk's W Test for Normality to Large Samples", Applied Statistics, 2, 115-124.
24. Shapiro, S.S. and Francia, R.S. (1972), "An Approximate Analysis of Variance Test for Normality", JASA, 67, 215-216.
25. Shapiro, S.S. and Wilk, M.B. (1965), "An Analysis of Variance Test for Normality (Complete samples)", Biometrika, 52, 591-611.
26. Shapiro, S.S., Wilk, M.B. & Chen, H.J. (1968), "A Comparative Study of Various Tests for Normality", JASA, 63, 1343-1372.
27. Spiegelhalter, D. J. (1977), "A Test for Normality against symmetric Alternatives", Biometrika, 64, 415-418.
28. Vijaya Kumar, K. (2013), "Nonnormal Disturbances; Testing Normality in Linear Statistical Models", Ph.D., Thesis, S.V. University.
29. Weisberg, S. and Bingham, C. (1975), "An Approximate Analysis of Variance Test for Non-Normality suitable for Machine Calculation", Technometrics, 17, 133-134
30. White, H. and Mac Donald, G.M. (1980), "Some Large-Sample Tests for Nonnormality in the Linear Regression Model", JASA, 75, 16-28.

Source of support: Nil, Conflict of interest: None Declared