

## THE ASYMPTOTIC EQUIVALENCE OF SOME MODIFIED SHAPIRO-WILK STATISTICS—COMPLETE AND CENSORED SAMPLE CASES

BY STEVE VERRILL<sup>1</sup> AND RICHARD A. JOHNSON<sup>2</sup>

*Lawrence Livermore National Laboratory and  
University of Wisconsin-Madison*

The Shapiro-Wilk statistic and its modifications are widely applied in tests for normality. We establish the asymptotic equivalence of a class of statistics based on different choices of normal scores. In particular, we conclude that the Shapiro-Francia, Filliben, Weisberg-Bingham and de Wet-Venter versions of the statistic are asymptotically equivalent. Our results also apply to the Type I and Type II censored data cases.

**1. Introduction.** Shapiro and Wilk (1965) proposed a test of normality based on the statistic

$$(1.1) \quad W \equiv \frac{(\mathbf{m}'\mathbf{V}^{-1}\mathbf{Y})^2}{(\mathbf{m}'\mathbf{V}^{-1}\mathbf{V}^{-1}\mathbf{m})\sum_{i=1}^n(Y_{in} - \bar{Y})^2},$$

where  $Y_{in}$  is the  $i$ th order statistic from some sample, and  $\mathbf{m}$ ,  $\mathbf{V}$  are the expectation vector and covariance matrix of the order statistics for a standard normal sample. Shapiro and Francia (1972) introduced the modified  $W$ -statistic

$$(1.2) \quad W' \equiv \frac{(\mathbf{m}'\mathbf{Y})^2}{(\mathbf{m}'\mathbf{m})\sum_{i=1}^n(Y_{in} - \bar{Y})^2},$$

which is easier to calculate for large samples. In order to further reduce computation costs, Weisberg and Bingham (1975) proposed modifying  $W'$  by replacing  $m_{in}$  by  $H((i - \frac{3}{8})/(n + \frac{1}{4}))$ . Here,  $H$  is the inverse of the  $N(0,1)$  distribution function.

Filliben (1975) and Ryan and Joiner (1973) noted that the Shapiro-Francia statistic could be written as

$$(1.3) \quad W' = \left[ \sum_{i=1}^n (Y_{in} - \bar{Y})(m_{in} - \bar{m}) \right]^2 / \left[ \sum_{i=1}^n (Y_{in} - \bar{Y})^2 \sum_{i=1}^n (m_{in} - \bar{m})^2 \right],$$

which is the square of the correlation coefficient associated with a normal

---

Received March 1983; revised May 1986.

<sup>1</sup>Work performed under the auspices of the U.S. Department of Energy at the Lawrence Livermore National Laboratory under Contract W-7405-Eng-48.

<sup>2</sup>Research supported by Office of Naval Research Contract No. N00014-78-C-0722.

AMS 1980 subject classifications. Primary 62F99, 62E20; secondary 62G99.

Key words and phrases. Correlation tests of normality, modified Shapiro-Wilk statistics, Shapiro-Francia statistic, asymptotic equivalence, Type I censoring, Type II censoring.

probability plot. This interpretation of  $W'$  led Filliben to propose

$$(1.4) \quad r_F \equiv \frac{\sum_{i=1}^n (Y_{in} - \bar{Y})(M_{in} - \bar{M})}{\sqrt{\sum_{i=1}^n (Y_{in} - \bar{Y})^2 \sum_{i=1}^n (M_{in} - \bar{M})^2}}$$

as an alternate version of the Shapiro–Wilk statistic. Here,  $M_{in}$  is the median of the  $i$ th order statistic from a standard normal sample.

de Wet and Venter (1972, 1973) first obtained the asymptotic distribution of the modified Shapiro–Wilk statistic

$$(1.5) \quad r \equiv \frac{\sum_{i=1}^n (Y_{in} - \bar{Y})(H_{in} - \bar{H})}{\sqrt{\sum_{i=1}^n (Y_{in} - \bar{Y})^2 \sum_{i=1}^n (H_{in} - \bar{H})^2}},$$

where  $H_{in} \equiv H(i/(n+1))$ .

Type II censored data versions of these correlation statistics were introduced in the Monte Carlo studies of Smith and Bain (1976) and Gerlach (1980).

Leslie, Stephens and Fotopoulos (1986) established the asymptotic equivalence of Shapiro and Wilk's  $W$  and de Wet and Venter's  $r$ .

In this paper we prove that the members of a class of statistics that includes  $\sqrt{W'}$  and  $r_F$  have the same asymptotic distribution as  $r$ . Moreover, our results show that this asymptotic equivalence also holds for the Type I and Type II censored data cases.

**2. The main result.** The asymptotic theory has been established for correlation statistics defined in terms of the scores  $H(i/(n+1))$ . de Wet and Venter (1972, 1973) showed that

$$(2.1) \quad 2n(1-r) - a_n \rightarrow_L \sum_{m=3}^{\infty} (W_m^2 - E(W_m^2))/m,$$

where the  $W_m$  are independent standard normals and  $\{a_n\}$  is a known sequence of constants.

The censored data versions of (2.1) were established by Verrill and Johnson (1983) [see also Verrill (1981)]. The Type I and Type II censored data cases produce the same limit. In essence, it differs from the complete sample limit in just two ways. When censoring is present the  $W_m$ 's are correlated, and there are six additional terms in the series. With Type II censoring at  $Y_{[n\delta]}$ ,  $0 < \delta < 1$ , the statistic is

$$(2.2) \quad r_{n,\delta} \equiv \frac{\sum_{i=1}^{[n\delta]} (Y_{in} - \bar{Y}_{n,\delta})(H_{in} - \bar{H}_{n,\delta})}{\sqrt{\sum_{i=1}^{[n\delta]} (Y_{in} - \bar{Y}_{n,\delta})^2 \sum_{i=1}^{[n\delta]} (H_{in} - \bar{H}_{n,\delta})^2}},$$

where  $\bar{Y}_{n,\delta} \equiv \sum_{i=1}^{[n\delta]} Y_{in}/[n\delta]$  and  $\bar{H}_{n,\delta} \equiv \sum_{i=1}^{[n\delta]} H_{in}/[n\delta]$ . For Type I censoring, the upper limit of summation is replaced by the number of uncensored observa-

tions. The censored data statistic (2.2) also has modified versions which correspond to the scores  $m_{in}$ ,  $M_{in}$  and  $H((i - \frac{3}{8})/(n + \frac{1}{4}))$ .

**THEOREM 2.1.** *Let the observations be independent and identically distributed normal random variables. Then the correlation statistics based on  $H(i/(n + 1))$ ,  $m_{in}$ ,  $M_{in}$  and  $H((i - \frac{3}{8})/(n + \frac{1}{4}))$  are asymptotically equivalent. This equivalence also holds for Type I and Type II censoring.*

**PROOF.** In the next section we will establish a theorem that identifies a condition on the scores that ensures asymptotic equivalence. In Lemma 3.2, we present a simpler sufficient condition which is then verified to hold, in Lemma 3.3, for the special scores mentioned above.

**3. Conditions for equivalence.** The asymptotic distribution results are known for the choice of scores  $H(i/(n + 1))$ . The proof of the equivalence for correlation statistics based on other scores will be presented in terms of Type II censoring.

Let  $b_{1n} \leq \dots \leq b_{nn}$  be a sequence of potential scores. Define

$$(3.1) \quad r_{n,\delta}(b) \equiv \frac{\sum_{i=1}^{[n\delta]} (Y_{in} - \bar{Y}_{n,\delta})(b_{in} - \bar{b}_{n,\delta})}{\sqrt{\sum_{i=1}^{[n\delta]} (Y_{in} - \bar{Y}_{n,\delta})^2 \sum_{i=1}^{[n\delta]} (b_{in} - \bar{b}_{n,\delta})^2}}.$$

In this notation,  $r_{n,\delta} = r_{n,\delta}(H)$ .

**THEOREM 3.1.** *Let the observations be i.i.d. normal random variables and suppose that*

$$(3.2) \quad \lim_{n \rightarrow \infty} \log \log(n) \sum_{i=1}^n (b_{in} - H_{in})^2 = 0.$$

Then

$$n(r_{n,\delta}(b) - r_{n,\delta}(H)) \rightarrow_P 0.$$

**PROOF.** The proof is given for Type II censoring. Type I censoring can be treated in the same manner. We write  $m = [n\delta]$  and introduce vectors  $\mathbf{Y} = (Y_{1n} - \bar{Y}_{n,\delta}, \dots, Y_{mn} - \bar{Y}_{n,\delta})^T$ ,  $\mathbf{H} = (H_{1n} - \bar{H}_{n,\delta}, \dots, H_{mn} - \bar{H}_{n,\delta})^T$  and  $\mathbf{b} = (b_{1n} - \bar{b}_{n,\delta}, \dots, b_{mn} - \bar{b}_{n,\delta})^T$ .

For any vector  $\mathbf{a}$  we write  $\|\mathbf{a}\| = (\mathbf{a}^T \mathbf{a})^{1/2}$ .

Since  $\|\mathbf{Y}\|$  is obviously of exact order  $n^{1/2}$  in probability, we have to show that

$$(3.3) \quad n^{1/2}(\alpha - \beta)^T \mathbf{Y} \rightarrow_P 0,$$

where  $\alpha = \mathbf{H}/\|\mathbf{H}\|$  and  $\beta = \mathbf{b}/\|\mathbf{b}\|$ . If  $\eta\mathbf{H}$  is the projection of  $\mathbf{b}$  on  $\mathbf{H}$ , we have  $\mathbf{b} = \eta\mathbf{H} + \mathbf{c}$  with  $\mathbf{c}^T \mathbf{H} = 0$  and in view of (3.2)

$$(3.4) \quad \|\mathbf{b} - \mathbf{H}\|^2 = (\eta - 1)^2 \|\mathbf{H}\|^2 + \|\mathbf{c}\|^2 = o((\log \log n)^{-1}) = o(1)$$

as  $n \rightarrow \infty$ . Since  $\|\mathbf{H}\|^2$  is of exact order  $n$ , this implies that

$$\begin{aligned} |\eta - 1| &= o(n^{-1/2}), & \|\mathbf{c}\|^2 &= o((\log \log n)^{-1}), \\ \|\mathbf{b}\|^2 &= \eta^2 \|\mathbf{H}\|^2 + \|\mathbf{c}\|^2 = \|\mathbf{H}\|^2(1 + 2(\eta - 1) + o(n^{-1})), \\ \boldsymbol{\beta} &= \frac{\mathbf{b}}{\|\mathbf{H}\|}(1 - (\eta - 1) + o(n^{-1})), \\ (\boldsymbol{\alpha} - \boldsymbol{\beta}) &= \frac{\mathbf{H}}{\|\mathbf{H}\|} [(1 - \eta)^2 + o(n^{-1})] - \frac{\mathbf{c}}{\|\mathbf{H}\|} [2 - \eta + o(n^{-1})] \\ &= \frac{\mathbf{H}}{\|\mathbf{H}\|} o(n^{-1}) - \frac{\mathbf{c}}{\|\mathbf{H}\|} (1 + o(n^{-1/2})), \end{aligned}$$

where the order symbols refer to scalar quantities. Since  $\mathbf{c}^T \mathbf{H} = 0$ , it follows that

$$(3.5) \quad n^{1/2}(\boldsymbol{\alpha} - \boldsymbol{\beta})^T \mathbf{H} = \|\mathbf{H}\| o(n^{-1/2}) = o(1),$$

$$\begin{aligned} (3.6) \quad \left[ n^{1/2}(\boldsymbol{\alpha} - \boldsymbol{\beta})^T (\mathbf{Y} - \mathbf{H}) \right]^2 &\leq n \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|^2 \|\mathbf{Y} - \mathbf{H}\|^2 \\ &= \|\mathbf{Y} - \mathbf{H}\|^2 \left[ o(n^{-1}) + O\left(\frac{n\|\mathbf{c}\|^2}{\|\mathbf{H}\|^2}\right) \right] \\ &= \|\mathbf{Y} - \mathbf{H}\|^2 o((\log \log n)^{-1}). \end{aligned}$$

Together, (3.5) and (3.6) are sufficient to prove (3.3) and the theorem if we can show that

$$(3.7) \quad E\|\mathbf{Y} - \mathbf{H}\|^2 \leq E \sum_{i=1}^m (Y_{in} - H_{in})^2 \leq \sum_{i=1}^n E(Y_{in} - H_{in})^2 = O(\log \log n).$$

We have

$$\sum_{i=1}^n E(Y_{in} - H_{in})^2 = \sum_{i=1}^n E(Y_{in} - m_{in})^2 + \sum_{i=1}^n (m_{in} - H_{in})^2.$$

By Lemmas 3.2 and 3.3, the second term is  $o(\log \log n)$ .

The quantity  $\sum_{i=1}^n E(Y_{in} - m_{in})^2$  was well studied in connection with asymptotic expansions. Combining (A2.22) in Albers, Bickel and van Zwet (1976) with (5.56) in Bickel and van Zwet (1978) and the fact that the Albers et al. condition  $R_2$  holds since

$$(3.8) \quad \lim_{t \rightarrow 0,1} t(1-t)|H(t)|H'(t) = \lim_{x \rightarrow \pm\infty} \frac{|x|\Phi(x)(1-\Phi(x))}{\phi(x)} = 1,$$

we find

$$(3.9) \quad \sum_{i=1}^n E(Y_{in} - m_{in})^2 = \log \log(n) + O(1),$$

which proves the theorem.  $\square$

LEMMA 3.2. Set  $c_{in} = \Phi(b_{in})$ . If for all  $i \in \{1, \dots, n\}$ ,  $n \geq N$ , we have

$$(3.10) \quad \frac{d}{n+1} < c_{in} < 1 - \frac{d}{n+1}, \quad \left| c_{in} - \frac{i}{n+1} \right| \leq \frac{K}{n+1},$$

for fixed positive  $N$ ,  $d$  and  $K$ , then

$$(3.11) \quad \sum_{i=1}^n (b_{in} - H_{in})^2 = O((\log n)^{-1}).$$

PROOF. Using (3.8) and taking logs, we see that  $H^2(t) = O(-\log[t(1-t)])$  for  $t \rightarrow 0$  or  $1$ . Hence  $H'(t) = O(\{t(1-t)\}^{-1}\{-\log[t(1-t)]\}^{-1/2})$ . Now the assumptions of the lemma imply that for points  $\tilde{c}_{in}$  between  $c_{in}$  and  $i/(n+1)$ ,

$$\begin{aligned} \sum_{i=1}^n (b_{in} - H_{in})^2 &= \sum_{i=1}^n \left( H(c_{in}) - H\left(\frac{i}{n+1}\right) \right)^2 \\ &= \sum_{i=1}^n [H'(\tilde{c}_{in})]^2 \left( c_{in} - \frac{i}{n+1} \right)^2 \\ &\leq \frac{K^2}{(n+1)^2} \sum_{i=1}^n [H'(\tilde{c}_{in})]^2 \\ &= O\left( n^{-1} \int_{d/n}^{1-d/n} [H'(t)]^2 dt \right) \\ &= O\left( n^{-1} \int_{d/n}^{1-d/n} \frac{dt}{t^2(1-t)^2 \log[t(1-t)]} \right) \\ &= O((\log n)^{-1}). \quad \square \end{aligned}$$

We now show that the Shapiro-Francia, Weisberg-Bingham and Filliben scores satisfy condition (3.10).

LEMMA 3.3. (i) (Shapiro and Francia). Let  $Z_{in}$  be the  $i$ th order statistic from a sample of  $n$   $N(0, 1)$  random variables. Then  $b_{in}^{(1)} \equiv E(Z_{in})$  satisfies condition (3.10).

(ii) Let  $b_{in}^{(2)} \equiv H((i+a)/(n+b))$  where  $a > -1$  and  $a < b$ . Then  $b_{in}^{(2)}$  satisfies condition (3.10). [This class of scores contains the  $H((i-\alpha)/(n-2\alpha+1))$ ,  $\alpha < 1$  scores. The Weisberg-Bingham score is a special case of these.]

(iii) (Filliben).

$$b_{in}^{(3)} \equiv \begin{cases} H\left(1 - \left(\frac{1}{2}\right)^{1/n}\right) & \text{for } i = 1, \\ H\left((i - 0.3175)/(n + 0.365)\right) & \text{for } i \in \{2, \dots, n - 1\}, \\ H\left(\left(\frac{1}{2}\right)^{1/n}\right) & \text{for } i = n \end{cases}$$

satisfies condition (3.10).

(iv) (Filliben). Let  $M_{i_n}$  denote the median of the  $i$ th order statistic from a sample of  $n$   $U(0, 1)$  random variables. Then  $b_{i_n}^{(4)} \equiv H(M_{i_n})$  satisfies (3.10).

PROOF. (i) It is known [see David (1981), page 77] that for  $i \leq (n + 1)/2$ ,

$$(i - \frac{1}{2})/n \leq \Phi(E(Z_{i_n})) \leq i/(n + 1)$$

and for  $i \geq (n + 1)/2$ ,

$$(i - \frac{1}{2})/n \geq \Phi(E(Z_{i_n})) \geq i/(n + 1).$$

Thus (3.10) holds.

(ii) Clear.

(iii) For  $i = \{2, \dots, n - 1\}$  the result is clear. For  $i = 1$ , the first condition in (3.10) holds if

$$1 - \left(\frac{1}{2}\right)^{1/n} > d/(n + 1) \quad \text{or} \quad \left(1 - \frac{d}{n + 1}\right)^n > \frac{1}{2}.$$

Take  $d = \frac{1}{4}$  and  $n$  sufficiently large. The second condition follows from  $1 - (\frac{1}{2})^{1/n} < 1/(n + 1)$  for  $n$  sufficiently large. By symmetry, (3.10) holds for  $i = n$ .

(iv) Since  $U_{i_n}$  is Beta( $i, n + 1 - i$ ),  $M_{1_n} = 1 - (\frac{1}{2})^{1/n}$  and  $M_{n_n} = (\frac{1}{2})^{1/n}$ . We handled these in (iii). Also, for  $i < (n + 1)/2$ ,

$$\begin{aligned} P[U_{i_n} < \text{Mode}] &= P\left[\text{Bin}\left(n, \frac{i - 1}{n - 1}\right) > i - 1\right] \\ &< \frac{1}{2} \\ &< P\left[\text{Bin}\left(n, \frac{i}{n + 1}\right) > i - 1\right] \\ &= P[U_{i_n} < \text{Mean}], \end{aligned}$$

where the inequalities are given in Johnson and Kotz (1969), page 53. Consequently,  $M_{i_n}$  lies between the mean and the mode of  $U_{i_n}$  which are  $i/(n + 1)$  and  $(i - 1)/(n - 1)$ . Thus it is easy to show that  $H(M_{i_n})$  satisfies (3.10). Essentially the same proof holds for  $i \geq (n + 1)/2$ .  $\square$

REMARK. We have established the asymptotic equivalence of various test statistics under the null hypothesis of normality. To study asymptotic power, we could consider contiguous alternatives. By the definition of contiguity, the equivalence of the statistics will continue to hold. These are the main facts that one would want to know from an asymptotic analysis.

REMARK. Under the conditions of Lemma 3.2, we can use the decomposition of  $\mathbf{b}$  [above (3.4)] and the expression for  $\|\mathbf{b}\|^2$  [below (3.4)] to show that the scores  $b_{i_n}$  and  $H_{i_n}$  are highly correlated. In particular,  $n(\text{corr}(b_{i_n}, H_{i_n}) - 1) \rightarrow 0$ .

**Acknowledgment.** The authors wish to thank Willem van Zwet for suggestions that made the proof both more concise and more insightful.

## REFERENCES

- ALBERS, W., BICKEL, P. J. and VAN ZWET, W. R. (1976). Asymptotic expansions for the power of distribution free tests in the one-sample problem. *Ann. Statist.* 4 108–156.
- BICKEL, P. J. and VAN ZWET, W. R. (1978). Asymptotic expansions for the power of distributionfree tests in the two-sample problem. *Ann. Statist.* 6 937–1004.
- DAVID, H. A. (1981). *Order Statistics*. Wiley, New York.
- DE WET, T. and VENTER, J. H. (1972). Asymptotic distributions of certain test criteria of normality. *South African Statist. J.* 6 135–149.
- DE WET, T. and VENTER, J. H. (1973). Asymptotic distributions for quadratic forms with applications to tests of fit. *Ann. Statist.* 1 380–387.
- FILLIBEN, J. J. (1975). The probability plot correlation coefficient test for normality. *Technometrics* 17 111–117.
- GERLACH, B. (1980). A correlation-type goodness-of-fit test for normality with censored sampling. *Math. Operationsforsch. Statist. Ser. Statist.* 11 207–218.
- JOHNSON, N. L. and KOTZ, S. (1969). *Discrete Distributions*. Houghton-Mifflin, Boston.
- LESLIE, J. R., STEPHENS, M. A. and FOTOPOULOS, S. (1986). Asymptotic distribution of the Shapiro–Wilk  $W$  for testing for normality. *Ann. Statist.* 14 1497–1506.
- RYAN, T. and JOINER, B. (1973). Normal probability plots and tests for normality. Technical Report, Pennsylvania State Univ.
- SHAPIRO, S. S. and FRANCA, R. S. (1972). An approximate analysis of variance test for normality. *J. Amer. Statist. Assoc.* 67 215–216.
- SHAPIRO, S. S. and WILK, M. B. (1965). Analysis of variance test for normality (complete samples). *Biometrika* 52 591–611.
- SMITH, R. M. and BAIN, L. J. (1976). Correlation type goodness-of-fit statistics with censored sampling. *Comm. Statist. A—Theory Methods* 5 119–132.
- VERRILL, S. P. (1981). Some asymptotic results concerning censored data versions of the Shapiro–Wilk goodness of fit test. Ph.D. thesis, Univ. of Wisconsin, Madison.
- VERRILL, S. P. and JOHNSON, R. A. (1983). The asymptotic distributions of censored data versions of the Shapiro–Wilk test of normality statistic. Technical Report 702, Dept. of Statistics, Univ. of Wisconsin, Madison.
- WEISBERG, S. and BINGHAM, C. (1975). An approximate analysis of variance test for nonnormality suitable for machine calculation. *Technometrics* 17 133–134.

LAWRENCE LIVERMORE NATIONAL LABORATORY  
P.O. Box 808  
LIVERMORE, CALIFORNIA 94550

DEPARTMENT OF STATISTICS  
UNIVERSITY OF WISCONSIN  
1210 WEST DAYTON STREET  
MADISON, WISCONSIN 53706