

ON DATA DRIVEN NEYMAN'S TESTS

BY

WILBERT C. M. KALLENBERG (ENSCHEDÉ) AND TERESA LEDWINA* (WROCLAW)

Abstract. The smooth tests for testing uniformity were introduced by Neyman [15]. The data driven method of selecting the number of components in a smooth test for uniformity is discussed, including the first-order asymptotic null distribution, consistency, empirical critical values and Monte Carlo powers. The first-order asymptotic null distribution is not sufficiently precise for approximation tools. A substantial improvement is made in this paper by deriving a second-order approximation of the null distribution, which turns out to be very accurate in numerical examples. The approximations are based on the second-order behaviour of Schwarz's selection rule S under uniformity. The new results on S are of independent interest.

1. Introduction. Let X_1, X_2, \dots, X_n be i.i.d. r.v.'s. Consider the goodness-of-fit problem of testing the simple null hypothesis H_0 that the X_i 's have distribution function F_0 , where F_0 is a given continuous distribution function. Without loss of generality we assume that under H_0 the distribution of X_i is the uniform distribution on $[0, 1]$.

Neyman's [15] paper of 1937, called '*Smooth test for goodness of fit*', is seen as the starting point of a subbranch of goodness-of-fit tests, the smooth tests of goodness-of-fit. In fact, Neyman's test is a particular example of the whole class of smooth goodness-of-fit tests, taking the Legendre polynomials as the orthonormal system. In [16], pp. 19-22, it is memorized that Neyman's test is intended as an "optimal" competitor to Pearson's chi-squared test. *Optimal* means here that the test is locally uniformly most powerful symmetric, unbiased for testing uniformity against alternatives of the form

$$(1.1) \quad c(\theta) \exp \left\{ \sum_{j=1}^k \theta_j \pi_j(x) \right\}, \quad 0 \leq x \leq 1,$$

where $\{\pi_j\}$ are the orthonormal Legendre polynomials on $[0, 1]$ and $c(\theta)$ is a normalizing constant. Alternatives of the form (1.1) were vaguely described

* Research supported by Grant MEN 341114.

as *smooth*. The test statistic for testing $H_0: \theta_j = 0, j = 1, \dots, k$, corresponding to uniformity, is given by

$$(1.2) \quad N_k = \sum_{j=1}^k \left\{ n^{-1/2} \sum_{i=1}^n \pi_j(X_i) \right\}^2.$$

Replacing π_j in (1.2) by

$$(1.3) \quad p_j(x) = \{1_{(d_{j-1}, d_j]}(x) - (d_j - d_{j-1})\} (d_j - d_{j-1})^{-1/2},$$

where $1_A(x)$ denotes the indicator function of the set A , we get Pearson's chi-square test

$$(1.4) \quad P_k = \sum_{j=1}^k \left\{ n^{-1/2} \sum_{i=1}^n p_j(X_i) \right\}^2 = \sum_{j=1}^k (O_j - E_j)^2 / E_j$$

with partition $0 = d_0 < d_1 < \dots < d_k = 1$, O_j and E_j denoting the observed and expected (under uniformity) number of observations in the j -th interval, respectively.

Hamdan [6], [8] implemented a suggestion of Lancaster to use different orthonormal systems, leading to what nowadays is called the *class of smooth tests for goodness-of-fit*, given by the test statistics

$$(1.5) \quad T_k = \sum_{j=1}^k \left\{ n^{-1/2} \sum_{i=1}^n \phi_j(X_i) \right\}^2,$$

where ϕ_0, ϕ_1, \dots is an orthonormal system in $L_2([0, 1])$ with $\phi_0(x) \equiv 1$. For more extensive and general information on the class of smooth tests we refer to Rayner and Best [18]. After the papers of Barton [2]–[4], Watson [21], and Hamdan [6]–[8] had appeared, there was no much progress in the field. A possible reason is that the required computations are quite heavy by hand. This is no longer a barrier with modern computers. Although smooth tests arose little interest for several years, nowadays Neyman's paper [15] is considered to be ingenious (cf. Le Cam and Lehmann [12], p. ix). Recently, there was a renewed interest in smooth tests. For an overview we refer to Rayner and Best [19], who conclude in reviewing several tests of fit: "don't use those other methods — use a smooth test!". The same conclusion is derived in Milbrodt and Strasser [14], p. 14. One of the important issues in applying T_k is the number k of components. We discuss three different approaches.

In the first one, k (and the orthonormal system) are selected such that the alternatives of particular interest are represented using only the first k components of the orthonormal system with k as small as possible. It requires that the user not automatically applies the test, but thinks why the goodness-of-fit test is executed and what types of alternatives are of particular interest. This imperative task is useful on itself, and therefore the method is more advantageous than disadvantageous. While especially focused on the selected (type of)

alternatives, the test keeps its omnibus character. A criterion of a simple structure to implement the idea is presented and extensively motivated in Inglot et al. [9]. This approach is closest to the classical application of smooth tests, since the test statistic is not modified and a formal test is executed.

The second approach is a more data-analytic approach. Obviously, it is not valid to apply a lot of significance tests, by varying k , to a data set and focus only on the most critical of them. Neyman was very clear about it as can be seen from the citation on p. 47 of [18]. Nevertheless, if the null hypothesis is rejected, the components may be used informally to suggest the nature of the departure from the null hypothesis. This second approach is described in detail in [18].

Returning to formal testing theory, the third approach concerns a data driven selection of k . This idea was introduced by Ledwina [13]. Roughly speaking, it works as follows. First, Schwarz's [20] selection rule is applied to find a suitable dimension S , say, of an exponential family model for the data. Then Neyman's test is applied within the fitted model, resulting in the test statistic N_S . So Schwarz's rule serves as a kind of first selection, followed by a more precise instrument, being Neyman's test in the "right" dimension. As with the original test of Neyman, the data driven version can also be applied using other orthonormal systems than the Legendre polynomials.

Observe that the "data driven approach" has also the extra advantage lying behind the idea of the second approach. Namely, when rejecting the null hypothesis, automatically a well-defined alternative model is provided for the data at hand.

The rest of this paper is mainly concerned with this third approach. In Section 2 the method is described and properties of the tests as obtained by Ledwina [13] and Kallenberg and Ledwina [11] are discussed. In Section 3 new results on Schwarz's selection rule are presented. The results may be of independent interest. They are applied in Section 4, where an accurate new approximation of the null distribution of the data driven version of smooth goodness-of-fit tests is derived.

2. A data driven version of smooth goodness-of-fit tests. Consider the class of test statistics given by (1.5). Assume that the functions ϕ_1, ϕ_2, \dots are bounded, but not necessarily uniformly bounded in $j = 1, 2, \dots$. Define for $k = 1, 2, \dots$ exponential families by their densities $p_\theta(x)$ with respect to Lebesgue measure on $[0, 1]$ of the form

$$(2.1) \quad p_\theta(x) = \exp\{\theta \circ \phi(x) - \psi_k(\theta)\},$$

where

$$(2.2) \quad \theta = (\theta_1, \dots, \theta_k), \quad \phi = (\phi_1, \dots, \phi_k), \quad \psi_k(\theta) = \log \int_0^1 \exp\{\theta \circ \phi(x)\} dx,$$

and \circ stands for the inner product in R^k . When there is no confusion, the dimension k is sometimes suppressed in the notation.

It is well known that

$$(2.3) \quad \lambda(\theta) = \psi'(\theta) = E_{\theta} \phi(X),$$

where ' denotes derivative. Moreover, by orthonormality, we have

$$(2.4) \quad \lambda(0) = 0, \quad \psi''(0) = I \text{ (the identity matrix).}$$

Writing

$$(2.5) \quad Y_n = (\bar{\phi}_1, \dots, \bar{\phi}_k), \quad \bar{\phi}_j = n^{-1} \sum_{i=1}^n \phi_j(x_i)$$

we obtain the expression for the density of X_1, \dots, X_n , each X_i having density (2.1), in the form

$$(2.6) \quad \exp\{n[Y_n \circ \theta - \psi_k(\theta)]\}.$$

Schwarz's [20] Bayesian information criterion (BIC) for choosing sub-models corresponding to successive dimensions yields

$$(2.7) \quad S = \min\{k: 1 \leq k \leq d(n), n \sup_{\vartheta \in R^k} \{Y_n \circ \vartheta - \psi_k(\vartheta)\} - \frac{1}{2}k \log n \\ \geq n \sup_{t \in R^j} \{Y_n \circ t - \psi_j(t)\} - \frac{1}{2}j \log n, j = 1, \dots, d(n)\}.$$

Although it is not mentioned in the notation, S depends of course on the upper bound $d(n)$ of the exponential families under consideration. Schwarz's BIC is further discussed in Section 3.

The data driven smooth test statistic is defined by

$$(2.8) \quad T_S = \sum_{j=1}^S \{n^{-1/2} \sum_{i=1}^n \phi_j(X_i)\}^2$$

with S given by (2.7). The null hypothesis of uniformity is rejected for large values of T_S .

The idea behind (2.8) is that Schwarz's rule gives a first indication about the true density of the observations by fitting an exponential family model to the data and that the finishing touch comes from the smooth test in the selected exponential family.

Numerical results in, e.g., [9] show that Neyman's test performs very well provided a good choice of k has been made. On the other hand, examples in Table 2 of [9] and in Tables 3-5 in [11] show that a considerable loss of power may occur when a wrong choice of k is made. This illustrates that a good procedure for choosing k based on the data is very welcome. Both by theoretical results and by the Monte Carlo simulation it will be argued that T_S is such a procedure. To do this we start by considering N_S , the data driven

Table 1 on pseudorandom sequences different from those used to prepare Table 1.

It is clearly seen from Table 1 that the difference between the asymptotic 0.05 (0.01) critical value equal to 3.841 (6.635) and the simulated ones for $d(n) \geq 2$ is substantial, even for $n = 120$. Therefore, we may conclude that the first-order limiting theorem (Theorem 2.1) is not sufficiently precise for approximation tools, and hence there is a need for a second-order limiting theorem to improve the accuracy of the approximation of N_S under H_0 . This is done, more generally, for T_S in Section 4 on the basis of the results of Section 3.

To show that N_S has good power properties under a wide class of alternatives we consider its consistency under the alternative distribution P on $[0, 1]$. Suppose that

$$(2.11) \quad E_P \pi_1(X) = \dots = E_P \pi_{K-1}(X) = 0, \quad E_P \pi_K(X) \neq 0$$

for some $K = K(P)$, X having distribution P . Essentially any alternative of interest satisfies (2.11). (Note that the orthonormal system $\{\pi_j\}$ is complete.) The consistency of N_S is given by Corollary 4.5 of [11] and reads as follows.

THEOREM 2.2. *If (2.11) holds, and if*

$$(2.12) \quad \liminf_{n \rightarrow \infty} d(n) \geq K(P), \quad \lim_{n \rightarrow \infty} d^3(n)n^{-1} \log n = 0,$$

then the test based on N_S is consistent at P .

Since $K(P)$ is fixed, (2.9) together with $\lim_{n \rightarrow \infty} d(n) = \infty$ implies that the test based on N_S is consistent against any alternative of the form (2.11). Apart from the theoretical justification in Theorem 2.2, the good power properties of N_S for a wide range of alternatives are shown by simulation results in Tables 3 and 4 of [13] and Tables 3–5 of [11]. For illustration here we present in one table simulated powers of N_S , the widely recommended tests of Watson (W), Anderson and Darling (AD), Neuhaus (N) and the recently introduced procedures of Bickel and Ritov (BR), and Eubank and LaRiccia (ELR). For definitions of the several test statistics we refer to [13] and [11]. The alternatives under consideration are given by the following densities:

$$g_1(x) = \frac{1}{4}\{x^{-1/2} + (1-x)^{-1/2}\}, \quad g_2(x) = 1.5(1-x)^{1/2},$$

$$g_3(x) = (3\sqrt{2})\{x^{1/2}1_{[0,1/2]}(x) + (1-x)^{1/2}1_{(1/2,1]}(x)\}, \quad g_4(x) = 1 + \frac{1}{2}\cos(2\pi x),$$

$$g_5(x) = \exp\{-0.5\pi_2(x) - 0.2\pi_4(x) - \psi(0, -0.5, -0.2)\} \text{ (cf. (2.1)).}$$

Table 3 and other simulation results show that N_S has a stable and relatively high power for a broad range of alternatives.

TABLE 3. Estimated powers (in %) based on 10000 samples in each case; $\alpha = 0.05$, $n = 50$, $d(n) = 10$

Alter-natives	N_S	W	AD	N	BR	ELR
g_1	81	65	69	64	63	63
g_2	58	39	71	62	60	24
g_3	61	59	16	50	43	24
g_4	57	61	33	56	48	27
g_5	63	55	15	46	40	25

The theoretical results on the asymptotic null distribution of N_S and its consistency may be extended to other orthonormal systems as well. To do this we make some assumptions. Define

$$(2.13) \quad V_k = \max_{1 \leq j \leq k} \sup_{x \in [0,1]} |\phi_j(x)|.$$

ASSUMPTION 1. *The following condition holds:*

$$(2.14) \quad \lim_{n \rightarrow \infty} d(n) V_{d(n)} (n^{-1} \log n)^{1/2} = 0.$$

For the orthonormal Legendre polynomials on $[0, 1]$ we get $V_k = (2k+1)^{1/2}$. The asymptotic null distribution and the consistency of T_S are given by Theorems 3.4 and 4.4 in [11] and read as follows:

THEOREM 2.3. *If Assumption 1 holds, then under H_0*

$$(2.15) \quad T_S \xrightarrow{d} \chi_1^2,$$

and the test based on T_S is consistent against any alternative P of the form (2.11), provided that $\liminf_{n \rightarrow \infty} d(n) \geq K(P)$.

3. Schwarz's selection rule. It is seen in Section 2 that the first-order result $T_S \xrightarrow{d} \chi_1^2$ is not sharp enough to imply accurate approximations of the null distribution. To improve the approximations we need some second-order results. Both for technical reasons and for more insight in what is going on, we have to derive precise results on the behaviour of S under H_0 . The notation P_0 refers to the null hypothesis, where X_i has the uniform distribution on $[0, 1]$. While at the first order only the event $\{S = 1\}$ plays a role (cf. Theorems 3.2 and 3.4 in [11]), for the second-order results both $\{S = 1\}$ and $\{S = 2\}$ come in

the picture. The event $\{S \geq 3\}$ may still be neglected under H_0 as is seen in the following theorem:

THEOREM 3.1. *If Assumption 1 holds, then*

$$(3.1) \quad P_0(S \geq 3) = O(n^{-1}(\log n)^{1/2}) \quad \text{as } n \rightarrow \infty.$$

Proof. Let K be fixed with $3 \leq K \leq d(n)$. Write

$$(3.2) \quad P_0(S \geq 3) = \sum_{k=3}^K P_0(S = k) + \sum_{k=K+1}^{d(n)} P_0(S = k).$$

By the definition of S and the application of Lemma 3.2 in [10], there exists for each $k \in \{3, \dots, K\}$ a constant c_k^* and there exists a constant \tilde{c}_K such that

$$(3.3) \quad P_0(S = k) \leq P_0(n \sup_{\vartheta \in \mathbb{R}^k} \{Y_n \circ \vartheta - \psi_k(\vartheta)\} \geq \frac{1}{2}(k-1)\log n) \\ \leq c_k^*(\log n)^{(1/2)(k-2)} \exp\{-\frac{1}{2}(k-1)\log n\} \leq \tilde{c}_K(\log n)^{1/2} n^{-1}.$$

For $k = K+1, \dots, d(n)$ it follows from Lemma 3.1 in [11] that for

$$\varepsilon \in (0, \min\{1, \frac{2}{3}kV_k^2\}),$$

uniformly for $k \in \{1, \dots, d(n)\}$,

$$P_0(S = k) \leq P_0(nY_n \circ Y_n \geq (1-\varepsilon)(k-1)\log n),$$

provided that n is sufficiently large. Note that here we use Assumption 1. Next we apply formula (2) of [17] with

$$\varrho = \{(1-\varepsilon)(k-1)\log n\}^{1/2}, \quad m = k, \quad \lambda = 1,$$

$$a = k^{1/2} V_k n^{-1/2} \{(1-\varepsilon)(k-1)\log n\}^{1/2},$$

which yields

$$P_0(nY_n \circ Y_n \geq (1-\varepsilon)(k-1)\log n) \leq c_1 \left(\frac{\varrho^2}{2}\right)^{(1/2)(k-1)} \frac{1}{\Gamma(k/2)} \exp\{-\frac{1}{2}\varrho^2(1-\eta(a))\},$$

where c_1 is an absolute constant, $\eta(a) \rightarrow 0$ as $a \rightarrow 0$, provided that $\varrho^2/2 \geq k$ and $a \leq 1$ ($\eta(a)$ is explicitly given on p. 188 of [17] and satisfies $\eta(a) \leq a$). In view of Assumption 1, $a \rightarrow 0$ as $n \rightarrow \infty$ (uniformly for $k \in \{1, \dots, d(n)\}$). Moreover, $\varrho^2/2 \geq k$ for sufficiently large n (independent of k). By taking ε small enough we therefore get for any $\zeta > 0$ and $k \in \{1, \dots, d(n)\}$

$$(3.4) \quad P_0(S = k) \leq n^{-(1/2)(1-\zeta)(k-1)}$$

for $n = n(\zeta)$ large enough, independent of $k \in \{1, \dots, d(n)\}$. Taking, for instance, $K = 3$ and $\zeta = 0.2$ we get for n sufficiently large

$$P_0(S \geq 3) \leq \tilde{c}_3 (\log n)^{1/2} n^{-1} + \sum_{k=4}^{\infty} n^{-0.4(k-1)} = \tilde{c}_3 (\log n)^{1/2} n^{-1} + \frac{n^{-1.2}}{1 - n^{-0.4}},$$

and hence (3.1) holds. ■

Remark 3.1. It is clear from the proof that the main contribution comes from $S = 3$, as expected. In other words, by virtually the same proof (cf. (3.3) and (3.4)) it follows that for each fixed K

$$(3.5) \quad P_0(S \geq K) = O(n^{-(1/2)(K-1)} (\log n)^{(1/2)(K-2)}) \quad \text{as } n \rightarrow \infty. \quad \blacksquare$$

In view of Theorem 3.1 we may concentrate on the events $\{S = 1\}$ and $\{S = 2\}$. The event $\{S = 1\}$ means that dimension 1 is more important than dimension 2, dimension 3, dimension 4, etc. Under H_0 it turns out that the comparison of dimension 1 with dimension 2 is up to the second order the only comparison of interest. Similarly, on the event $\{S = 2\}$ the only important contribution under H_0 comes from the comparison of dimension 2 with dimensions 1 and 3. This is shown in the next two theorems.

THEOREM 3.2. *There exists a positive constant c_2 such that for each event A*

$$(3.6) \quad P_0(A, S = 1) \leq P_0(A, n\bar{\phi}_2^2 \leq \log n + c_2 n^{-1/2} (\log n)^{3/2}) + O(n^{-1} (\log n)^{-1/2})$$

and

$$(3.7) \quad P_0(A, S = 1) \geq P_0(A, n\bar{\phi}_2^2 \leq \log n - c_2 n^{-1/2} (\log n)^{3/2}) + O(n^{-1} (\log n)^{-1/2}) - P_0(S \geq 3)$$

as $n \rightarrow \infty$, uniformly in A .

Proof. By the definition of S we get

$$P_0(A, S = 1) \leq P_0(A, n \sup_{\vartheta \in \mathbb{R}} \{Y_n \vartheta - \psi_1(\vartheta)\} - \frac{1}{2} \log n \geq n \sup_{t \in \mathbb{R}^2} \{Y_n \circ t - \psi_2(t)\} - \log n).$$

For each fixed k we have

$$\sup_{\theta \in \mathbb{R}^k} \{y \circ \theta - \psi_k(\theta)\} = \frac{1}{2} \|y\|^2 + O(\|y\|^3) \quad \text{as } y \rightarrow 0,$$

where $\|\cdot\|$ denotes the Euclidean norm in \mathbb{R}^k . Therefore, on the set

$$\{y: |y_j| \leq n^{-1/2} (2 \log n)^{1/2}, j = 1, \dots, k\}$$

we get

$$\left| \sup_{\theta \in \mathbb{R}^k} \{y \circ \theta - \psi_k(\theta)\} - \frac{1}{2} \|y\|^2 \right| \leq \bar{c}_k n^{-3/2} (\log n)^{3/2}$$

for some constant \bar{c}_k . Hence, uniformly in A ,

$$\begin{aligned} P_0(A, S = 1) &\leq P_0(A, n\{\frac{1}{2}\bar{\phi}_1^2 + \bar{c}_1 n^{-3/2}(\log n)^{3/2}\} + \frac{1}{2}\log n \\ &\geq n\{\frac{1}{2}(\bar{\phi}_1^2 + \bar{\phi}_2^2) - \bar{c}_2 n^{-3/2}(\log n)^{3/2}\}) \\ &\quad + P_0(n^{1/2}|\bar{\phi}_1| > (2\log n)^{1/2}) + P_0(n^{1/2}|\bar{\phi}_2| > (2\log n)^{1/2}) \\ &= P_0(A, n\bar{\phi}_2^2 \leq \log n + c_2 n^{-1/2}(\log n)^{3/2}) + O(n^{-1}(\log n)^{-1/2}) \end{aligned}$$

as $n \rightarrow \infty$, with $c_2 = 2(\bar{c}_1 + \bar{c}_2)$, where for $j = 1, 2$ we have used

$$P_0(n^{1/2}|\bar{\phi}_j| > (2\log n)^{1/2}) = O(n^{-1}(\log n)^{-1/2}),$$

which follows by standard large deviation theory. This completes the proof of (3.6).

To prove (3.7) note that

$$\begin{aligned} P_0(A, S = 1) &\geq P_0(A, n \sup_{\vartheta \in \mathbb{R}} \{Y_n \vartheta - \psi_1(\vartheta)\} - \frac{1}{2}\log n \\ &> n \sup_{t \in \mathbb{R}^2} \{Y_n \circ t - \psi_2(t)\} - \log n - P_0(S \geq 3). \end{aligned}$$

The rest of the proof is similar to the proof of (3.6). ■

THEOREM 3.3. *Let $d(n) \geq 2$. There exists a positive constant c_3 such that for each event A*

$$(3.8) \quad \begin{aligned} P_0(A, S = 2) &\leq P_0(A, n\bar{\phi}_2^2 \geq \log n - c_3 n^{-1/2}(\log n)^{3/2}, \\ &\quad n\bar{\phi}_3^2 \leq \log n + c_3 n^{-1/2}(\log n)^{3/2}) + O(n^{-1}(\log n)^{-1/2}) \end{aligned}$$

and

$$(3.9) \quad \begin{aligned} P_0(A, S = 2) &\geq P_0(A, n\bar{\phi}_2^2 \geq \log n + c_3 n^{-1/2}(\log n)^{3/2}, \\ &\quad n\bar{\phi}_3^2 \leq \log n - c_3 n^{-1/2}(\log n)^{3/2}) + O(n^{-1}(\log n)^{-1/2}) - P_0(S \geq 3) \end{aligned}$$

as $n \rightarrow \infty$, uniformly in A .

Proof. By the definition of S we get (as in the proof of Theorem 3.2)

$$\begin{aligned} P_0(A, S = 2) &\leq P_0(A, n \sup_{\vartheta \in \mathbb{R}^2} \{Y_n \circ \vartheta - \psi_2(\vartheta)\} - \log n \geq n \sup_{t \in \mathbb{R}} \{Y_n t - \psi_1(t)\} - \frac{1}{2}\log n, \\ &\quad n \sup_{\vartheta \in \mathbb{R}^2} \{Y_n \circ \vartheta - \psi_2(\vartheta)\} - \log n \geq n \sup_{\theta \in \mathbb{R}^3} \{Y_n \circ \theta - \psi_3(\theta)\} - \frac{3}{2}\log n) \\ &\leq P_0(A, n\{\frac{1}{2}(\bar{\phi}_1^2 + \bar{\phi}_2^2) + \bar{c}_2 n^{-3/2}(\log n)^{3/2}\} \\ &\quad - \log n \geq n\{\frac{1}{2}\bar{\phi}_1^2 - \bar{c}_1 n^{-3/2}(\log n)^{3/2}\} - \frac{1}{2}\log n, n\{\frac{1}{2}(\bar{\phi}_1^2 + \bar{\phi}_2^2) + \bar{c}_2 n^{-3/2}(\log n)^{3/2}\} \\ &\quad - \log n \geq n\{\frac{1}{2}(\bar{\phi}_1^2 + \bar{\phi}_2^2 + \bar{\phi}_3^2) - \bar{c}_3 n^{-3/2}(\log n)^{3/2}\} - \frac{3}{2}\log n) \end{aligned}$$

$$\begin{aligned}
& + \sum_{j=1}^3 P_0(n^{1/2}|\bar{\phi}_j| > (2\log n)^{1/2}) \\
& = P_0(A, n\bar{\phi}_2^2 \geq \log n - 2(\bar{c}_1 + \bar{c}_2)n^{-1/2}(\log n)^{3/2}, \\
& \quad n\bar{\phi}_3^2 \leq \log n + 2(\bar{c}_2 + \bar{c}_3)n^{-1/2}(\log n)^{3/2}) + O(n^{-1}(\log n)^{-1/2}).
\end{aligned}$$

By taking $c_3 = \max\{2(\bar{c}_1 + \bar{c}_2), 2(\bar{c}_2 + \bar{c}_3)\}$ the proof of (3.8) is completed. Inequality (3.9) is proved in a similar way. ■

Remark 3.2. The $O(n^{-1}(\log n)^{1/2})$ -terms in (3.6)–(3.9) may be replaced by $O(n^{-r}(\log n)^{-1/2})$ for arbitrary r . This is seen by replacing the set $\{y: |y_j| \leq n^{-1/2}(2\log n)^{1/2}, j = 1, \dots, k\}$ by the set $\{y: |y_j| \leq n^{-1/2}(2r\log n)^{1/2}, j = 1, \dots, k\}$ and noting that

$$P_0(n^{1/2}|\bar{\phi}_j| > (2r\log n)^{1/2}) = O(n^{-r}(\log n)^{-1/2}) \quad \text{as } n \rightarrow \infty. \quad \blacksquare$$

4. Second-order approximation to the null distribution of T_S . As was seen from Theorem 2.1 and Table 1, the first-order χ_1^2 -approximation of the null distribution of T_S is not very precise. Therefore, we consider the second-order approximation, based on the second-order limiting theorems for S under H_0 in the previous section. Denoting by U_1, U_2 two independent r.v.'s, each with a standard normal distribution, we shall prove the following result:

THEOREM 4.1. *Let $d(n) \geq 2$. If Assumption 1 holds, and Cramér's condition is satisfied, i.e.,*

$$(4.1) \quad \limsup_{|t| \rightarrow \infty} |E_{P_0} \exp\{it\phi_1(X_1)\}| < 1,$$

then

$$(4.2) \quad P_0(T_S \leq x) = \Pr(U_1^2 \leq x)\Pr(U_2^2 \leq \log n) \\ + \Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n) + O(n^{-1}\log n)$$

as $n \rightarrow \infty$, uniformly in x .

It is easily seen that the approximation on the right-hand side of (4.2) may also be written as

$$(4.3) \quad \Pr(U_1^2 \leq x) - \Pr(U_1^2 \leq x, U_1^2 + U_2^2 \geq x, U_2^2 \geq \log n).$$

In this formulation the correction term in comparison to the first-order approximation $\Pr(U_1^2 \leq x)$ is clearly seen. Corresponding to the simulation results in Table 1, indeed, $P_0(T_S \leq x)$ is overestimated by its first-order approximation. The two terms on the right-hand side of (4.2) are the leading terms of $P_0(T_S \leq x, S = 1)$ and $P_0(T_S \leq x, S = 2)$, respectively. For $x \leq \log n$

the second term disappears and the second-order approximation in that case is simply obtained by multiplying the first-order approximation by $\Pr(U_2^2 \leq \log n)$.

Condition (4.1) is satisfied if $\phi_1(X_j)$ has an absolutely continuous component, and therefore it is fulfilled in many cases, e.g., for Neyman's test with the orthonormal Legendre polynomials. To prove Theorem 4.1 we state and prove first several lemmas.

LEMMA 4.1. *For each fixed k and for each constant $c_4 \in \mathbb{R}$ we have*

$$(4.4) \quad P_0(n\bar{\phi}_k^2 \leq \log n + c_4 n^{-1/2}(\log n)^{3/2}) - P_0(n\bar{\phi}_k^2 \leq \log n) = O(n^{-1} \log n).$$

Proof. By standard large deviation theory (cf., e.g., [5], p. 553) we have

$$(4.5) \quad \begin{aligned} P_0(n\bar{\phi}_k^2 \leq \log n + c_4 n^{-1/2}(\log n)^{3/2}) - P_0(n\bar{\phi}_k^2 \leq \log n) \\ &= P_0(n\bar{\phi}_k^2 > \log n) - P_0(n\bar{\phi}_k^2 > \log n + c_4 n^{-1/2}(\log n)^{3/2}) \\ &= 2\Pr(U_1 > (\log n)^{1/2}) \{1 + O(n^{-1/2}(\log n)^{3/2})\} \\ &\quad - 2\Pr(U_1 > \{\log n + c_4 n^{-1/2}(\log n)^{3/2}\}^{1/2}) \{1 + O(n^{-1/2}(\log n)^{3/2})\} \\ &= O(n^{-1} \log n) \quad \text{as } n \rightarrow \infty. \quad \blacksquare \end{aligned}$$

LEMMA 4.2. *If Assumption 1 holds, then*

$$(4.6) \quad P_0(T_S \leq x, S = 1) = P_0(n\bar{\phi}_1^2 \leq x, n\bar{\phi}_2^2 \leq \log n) + O(n^{-1} \log n)$$

as $n \rightarrow \infty$, uniformly in x .

Proof. Taking $A = \{T_1 \leq x\}$ we see by Theorem 3.2 that there exists a positive constant c_2 such that

$$\begin{aligned} P_0(T_S \leq x, S = 1) &= P_0(T_1 \leq x, S = 1) \\ &\leq P_0(T_1 \leq x, n\bar{\phi}_2^2 \leq \log n + c_2 n^{-1/2}(\log n)^{3/2}) + O(n^{-1}(\log n)^{-1/2}) \\ &\leq P_0(T_1 \leq x, n\bar{\phi}_2^2 \leq \log n) + P_0(n\bar{\phi}_2^2 \leq \log n + c_2 n^{-1/2}(\log n)^{3/2}) \\ &\quad - P_0(n\bar{\phi}_2^2 \leq \log n) + O(n^{-1}(\log n)^{-1/2}) \end{aligned}$$

as $n \rightarrow \infty$, uniformly in x . The application of Lemma 4.1 and noting that $T_1 = n\bar{\phi}_1^2$ yield

$$P_0(T_S \leq x, S = 1) \leq P_0(n\bar{\phi}_1^2 \leq x, n\bar{\phi}_2^2 \leq \log n) + O(n^{-1} \log n)$$

as $n \rightarrow \infty$, uniformly in x . Similarly, the converse inequality is obtained by using (3.7) in combination with Theorem 3.1. \blacksquare

LEMMA 4.3. *For each constant $c_5 \in \mathbb{R}$ we have*

$$(4.7) \quad P_0(n\bar{\phi}_2^2 \geq \log n, n\bar{\phi}_3^2 > \log n - c_5 n^{-1/2}(\log n)^{3/2}) = O(n^{-1}(\log n)^{-1})$$

as $n \rightarrow \infty$.

Proof. Write

$$D_n = \{(y_1, y_2): y_1 \geq (\log n)^{1/2}, y_2 \geq \{\log n - c_5 n^{-1/2} (\log n)^{3/2}\}^{1/2}\}.$$

The set D_n is a convex Borel set and the point

$$a = ((\log n)^{1/2}, \{\log n - c_5 n^{-1/2} (\log n)^{3/2}\}^{1/2})$$

is the nearest point to the origin. The application of Assertion 2 in Remark 1 on p. 67 of [1] yields

$$\begin{aligned} P_0((n^{1/2} \bar{\phi}_2, n^{1/2} \bar{\phi}_3) \in D_n) &= \Pr((U_1, U_2) \in D_n) \{1 + O(n^{-1/2} (\log n)^{3/2})\} \\ &= O(n^{-1} (\log n)^{-1}) \quad \text{as } n \rightarrow \infty. \end{aligned}$$

By the same argument for the other three regions of similar form, which together with D_n constitute the event in (4.7), the proof is completed. ■

LEMMA 4.4. *If Assumption 1 holds, then*

$$(4.8) \quad P_0(T_S \leq x, S = 2) = P_0(n(\bar{\phi}_1^2 + \bar{\phi}_2^2) \leq x, n\bar{\phi}_2^2 \geq \log n) + O(n^{-1} \log n)$$

as $n \rightarrow \infty$, uniformly in x .

Proof. Taking $A = \{T_2 \leq x\}$ we see by Theorem 3.3 that there exists a positive constant c_3 such that

$$\begin{aligned} P_0(T_S \leq x, S = 2) &\leq P_0(T_2 \leq x, n\bar{\phi}_2^2 \geq \log n - c_3 n^{-1/2} (\log n)^{3/2}) + O(n^{-1} (\log n)^{-1/2}) \\ &\leq P_0(T_2 \leq x, n\bar{\phi}_2^2 \geq \log n) + P_0(\log n - c_3 n^{-1/2} (\log n)^{3/2} \leq n\bar{\phi}_2^2 < \log n) \\ &\quad + O(n^{-1} (\log n)^{-1/2}) \end{aligned}$$

as $n \rightarrow \infty$, uniformly in x . The application of Lemma 4.1 and noting that $T_2 = n(\bar{\phi}_1^2 + \bar{\phi}_2^2)$ yield

$$P_0(T_S \leq x, S = 2) \leq P_0(n(\bar{\phi}_1^2 + \bar{\phi}_2^2) \leq x, n\bar{\phi}_2^2 \geq \log n) + O(n^{-1} \log n)$$

as $n \rightarrow \infty$, uniformly in x .

For the converse inequality we consider (3.9) in combination with Theorem 3.1 and we obtain

$$\begin{aligned} P_0(T_S \leq x, S = 2) &\geq P_0(T_2 \leq x, n\bar{\phi}_2^2 \geq \log n) - P_0(\log n \leq n\bar{\phi}_2^2 < \log n + c_3 n^{-1/2} (\log n)^{3/2}) \\ &\quad - P_0(n\bar{\phi}_2^2 \geq \log n, n\bar{\phi}_3^2 > \log n - c_3 n^{-1/2} (\log n)^{3/2}) + O(n^{-1} (\log n)^{1/2}) \\ &\geq P_0(T_2 \leq x, n\bar{\phi}_2^2 \geq \log n) + O(n^{-1} \log n) \end{aligned}$$

by Lemmas 4.1 and 4.3. This completes the proof of the lemma. ■

LEMMA 4.5. *If (4.1) holds, then, uniformly for x , as $n \rightarrow \infty$*

$$(4.9) \quad P_0(n\bar{\phi}_1^2 \leq x) = \Pr(U_1^2 \leq x) + O(n^{-1}).$$

PROOF. By the standard Edgeworth expansion (cf., e.g., [5], Theorem 3 on p. 541) it follows that, uniformly in x ,

$$\begin{aligned} P_0(n\bar{\phi}_1^2 \leq x) &= P_0(n^{1/2}\bar{\phi}_1 \leq x^{1/2}) - P_0(n^{1/2}\bar{\phi}_1 < -x^{1/2}) \\ &= \Pr(U_1 \leq x^{1/2}) - \left\{ \frac{E\phi_1^3}{6n^{1/2}} \right\} \varphi(x^{1/2})(x-1) - \Pr(U_1 \leq -x^{1/2}) \\ &\quad + \left\{ \frac{E\phi_1^3}{6n^{1/2}} \right\} \varphi(-x^{1/2})(x-1) + O(n^{-1}) \\ &= \Pr(U_1^2 \leq x) + O(n^{-1}), \end{aligned}$$

where φ denotes the standard normal density. ■

LEMMA 4.6. *Uniformly for x as $n \rightarrow \infty$ we have*

$$(4.10) \quad P_0(n\bar{\phi}_1^2 \leq x, n\bar{\phi}_2^2 > \log n) = \Pr(U_1^2 \leq x)\Pr(U_2^2 > \log n) + O(n^{-1}\log n)$$

and

$$(4.11) \quad \begin{aligned} P_0(n(\bar{\phi}_1^2 + \bar{\phi}_2^2) \leq x, n\bar{\phi}_2^2 \geq \log n) \\ = \Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n) + O(n^{-1}\log n). \end{aligned}$$

PROOF. We give the proof of (4.11); the relation (4.10) can be shown in a quite similar way. Without loss of generality assume that $x \geq \log n$. Write

$$\begin{aligned} D_{n1} &= \{(y_1, y_2): y_2 \geq (\log n)^{1/2}, y_1^2 + y_2^2 \leq x\}, \\ D_{n2} &= \{(y_1, y_2): y_2 \leq -(\log n)^{1/2}, y_1^2 + y_2^2 \leq x\}. \end{aligned}$$

The sets D_{n1} and D_{n2} are convex. It follows by Assertion 2 in Remark 1 on p. 67 of [1] that for $i = 1, 2$, with $a = (0, (\log n)^{1/2})$ if $i = 1$ and $a = (0, -(\log n)^{1/2})$ if $i = 2$, we have

$$\begin{aligned} P_0((n^{1/2}\bar{\phi}_1, n^{1/2}\bar{\phi}_2) \in D_{ni}) &= \Pr((U_1, U_2) \in D_{ni}) \{1 + O(n^{-1/2}(\log n)^{3/2})\} \\ &= \Pr((U_1, U_2) \in D_{ni}) + O(n^{-1}\log n), \end{aligned}$$

which implies (4.11). ■

PROOF OF THEOREM 4.1. By (4.6), (4.8)–(4.11) and (3.1) we have, uniformly in x , as $n \rightarrow \infty$

$$\begin{aligned} P_0(T_S \leq x) &= P_0(T_S \leq x, S = 1) + P_0(T_S \leq x, S = 2) + P_0(T_S \leq x, S \geq 3) \\ &= P_0(n\bar{\phi}_1^2 \leq x) - P_0(n\bar{\phi}_1^2 \leq x, n\bar{\phi}_2^2 > \log n) \\ &\quad + P_0(n(\bar{\phi}_1^2 + \bar{\phi}_2^2) \leq x, n\bar{\phi}_2^2 \geq \log n) + O(n^{-1}\log n) \end{aligned}$$

$$= \Pr(U_1^2 \leq x) - \Pr(U_1^2 \leq x, U_2^2 > \log n) + \Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n) + O(n^{-1} \log n),$$

and hence (4.2) holds. ■

The approximation on the right-hand side of (4.2) is surprisingly simple and very accurate, especially when compared to the first-order approximation. The structure of the approximation clarifies that under H_0 up to the second order only $\{S = 1\}$ and $\{S = 2\}$ are important. Moreover, $\{S = 1\}$ and $\{S = 2\}$ may be further restricted to dimension 1 beats dimension 2, corresponding to $U_2^2 \leq \log n$, and dimension 2 beats dimension 1, corresponding to $U_2^2 > \log n$, respectively. (As a side-note, remark that here also dimension 2 beats dimension 3 is cancelled out.)

Both to investigate the accuracy of the approximation and to facilitate its evaluation we derive an upper and lower bound of $\Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n)$. Let for $x \geq \log n$

(4.12)

$$D^* = \{(y_1, y_2): y_1^2 + y_2^2 \geq x, 0 \leq y_1 \leq (x - \log n)^{1/2}, (\log n)^{1/2} \leq y_2 \leq x^{1/2}\}.$$

The area $a(x)$ of this set equals

$$(4.13) \quad a(x) = \{x^{1/2} - (\log n)^{1/2}\}(x - \log n)^{1/2} - \frac{1}{2}x \arccos\left(\frac{\log n}{x}\right)^{1/2} + \frac{1}{2}(\log n)^{1/2}(x - \log n)^{1/2},$$

while for $(y_1, y_2) \in D^*$ we have

$$(4.14) \quad n^{1/2} \exp(-x) \leq \exp\{-\frac{1}{2}(y_1^2 + y_2^2)\} \leq \exp(-\frac{1}{2}x).$$

Using symmetry we therefore get

$$(4.15) \quad \Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n) = 4\Pr(U_1^2 + U_2^2 \leq x, 0 \leq U_1 \leq (x - \log n)^{1/2}, (\log n)^{1/2} \leq U_2 \leq x^{1/2}) = 4\{\Pr(0 \leq U_1 \leq (x - \log n)^{1/2})\Pr((\log n)^{1/2} \leq U_2 \leq x^{1/2}) - (2\pi)^{-1}a(x)b(x)\}$$

with

$$(4.16) \quad n^{1/2} \exp(-x) \leq b(x) \leq \exp(-\frac{1}{2}x).$$

Replacing $\Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n)$ on the right-hand side of (4.2) by its lower and upper bounds obtained from (4.15) and (4.16) and inserting for x the empirical critical values of Table 1 we get the results compiled in Table 4. It is seen that for $d(n) = 2$ the approximation is excellent, even for small values of n . For $d(n) \geq 3$ we have a slight overestimation, but the approximation is

TABLE 4. Lower and upper bounds of the second-order approximation of $P_0(T_S \leq x)$ for empirical critical values x of Table 1 (probabilities $\times 10^5$)

$\alpha = 0.05$				$\alpha = 0.01$				
n	x	Lower bounds	Upper bounds	n	x	Lower bounds	Upper bounds	
20	5.367	94055	94912	20	8.451	98376	99095	
	6.537	96377	97341		10.216	99226	99662	
	7.213	97273	98182		11.332	99518	99817	
	7.617	97698	98552		12.183	99664	99885	
	7.854	97915	98733		12.997	99763	99926	
50	5.350	94758	95050	50	8.592	98654	99152	
	5.865	95784	96204		10.447	99378	99698	
	6.015	96043	96493		11.191	99545	99800	
	6.094	96173	96637		11.827	99652	99859	
	6.096	96176	96641		12.128	99693	99880	
	6.117	96210	96678		12.217	99704	99886	
	6.121	96217	96685		12.465	99734	99900	
80	5.260	94928	95028	80	8.511	98697	99109	
	5.508	95426	95579		9.853	99253	99578	
	5.581	95563	95732		10.239	99364	99659	
	5.592	95584	95755		10.386	99401	99686	
100	5.269	95122	95174	100	10.388	99402	99686	
	5.499	95562	95652		100	8.449	98704	99076
	5.557	95667	95768			9.941	99301	99597
	5.571	95692	95796			10.227	99380	99656
	5.581	95710	95815			10.352	99411	99679
	5.586	95719	95825			120	8.967	98981
120	5.032	94813	94819	10.272			99406	99664
	5.240	95214	95235	10.522	99464		99707	
	5.314	95353	95382	10.603	99482		99720	
	5.321	95366	95396	10.615	99485		99722	

also in that case rather accurate, and far more accurate than the first-order approximation.

As was noted before, if $x \leq \log n$, the second term in the approximation (4.2) disappears. On the other hand, if $x \geq 2\log n$, then

$$(4.17) \quad 0 \leq \Pr(U_2^2 \geq \log n) - \Pr(U_1^2 + U_2^2 \leq x, U_2^2 \geq \log n)$$

$$= \Pr(U_1^2 + U_2^2 > x, U_2^2 \geq \log n) \leq \Pr(U_1^2 + U_2^2 > x) = \exp(-\frac{1}{2}x) \leq n^{-1},$$

which implies for $x \geq 2\log n$

$$P_0(T_S \leq x) = \Pr(U_1^2 \leq x)\Pr(U_2^2 \leq \log n) + \Pr(U_2^2 \geq \log n) + O(n^{-1}\log n).$$

Therefore we may take the following very simple approximation of $P_0(T_S \leq x)$:

$$(4.18) \quad \begin{aligned} & \Pr(U_1^2 \leq x) \Pr(U_2^2 \leq \log n) && \text{if } x \leq \log n, \\ & \Pr(U_1^2 \leq x) \Pr(U_2^2 \leq \log n) + \Pr(U_2^2 \geq \log n) && \text{if } x \geq 2\log n, \\ & \text{linearize} && \text{if } \log n \leq x \leq 2\log n. \end{aligned}$$

To illustrate the approximation (4.18) consider $n = 50$ and the empirical 5% critical values x of Table 1. We obtain the following:

x	5.350	5.865	6.015	6.094	6.096	6.117	6.121
[approximation (4.18)] $\times 10^5$	93901	95068	95408	95587	95592	95639	95649

It is seen that the approximation works quite well in this case. As is seen from (4.17) the approximation (4.18) is larger than the right-hand side of (4.2) if $x \geq 2\log n$. It was seen in Table 4 that for the 1% critical values the right-hand side of (4.2) gives already a (slight) overestimation. Since as a rule in Table 4 the 1% critical values are larger than $2\log n$, the approximation (4.18) gives a (higher) overestimation of $P_0(T_S \leq x)$. A possible remedy is to replace the interval $[\log n, 2\log n]$ in (4.18) by, for instance, $[\log n, 3\log n]$ if small levels are concerned.

REFERENCES

- [1] A. K. Aleshkyavichene, *Multidimensional integral limit theorems for large deviation probabilities*, Theory Probab. Appl. 28 (1983), pp. 65–88.
- [2] D. E. Barton, *On Neyman's smooth test of goodness of fit and its power with respect to a particular system of alternatives*, Skand. Aktuarietidskr. 36 (1953), pp. 24–63.
- [3] – *A form of Neyman's ψ^2 test of goodness of fit applicable to grouped and discrete data*, ibidem 38 (1955), pp. 1–16.
- [4] – *Neyman's ψ^2 test of goodness of fit when the null hypothesis is composite*, ibidem 39 (1956), pp. 216–245.
- [5] W. Feller, *An Introduction to Probability Theory and Its Applications*, Vol. II, 2nd edition, Wiley, New York 1971.
- [6] M. A. Hamdan, *The powers of certain smooth tests of goodness of fit*, Austral. J. Statist. 4 (1962), pp. 25–40.
- [7] – *The number and width of classes in the chi-square test*, J. Amer. Statist. Assoc. 58 (1963), pp. 678–689.
- [8] – *A smooth test of goodness of fit based on the Walsh functions*, Austral. J. Statist. 6 (1964), pp. 130–136.
- [9] T. Inglot, W. C. M. Kallenberg and T. Ledwina, *Power approximations to and power comparison of smooth goodness-of-fit tests*, Scand. J. Statist. 21 (1994), pp. 131–145.
- [10] W. C. M. Kallenberg, *Bahadur deficiency of likelihood ratio tests in exponential families*, J. Multivariate Anal. 11 (1981), pp. 506–531.

- [11] — and T. Ledwina, *Consistency and Monte Carlo simulation of a data driven version of smooth goodness-of-fit tests*, Ann. Statist. (to appear).
- [12] L. LeCam and E. L. Lehmann, *J. Neyman — On the occasion of his 80th birthday*, ibidem 2 (1974), pp. vii–xiii.
- [13] T. Ledwina, *Data driven version of Neyman's smooth test of fit*, J. Amer. Statist. Assoc. 89 (1994), pp. 1000–1005.
- [14] H. Milbrodt and H. Strasser, *On the asymptotic power of the two-sided Kolmogorov-Smirnov test*, J. Statist. Plann. Inference 26 (1990), pp. 1–23.
- [15] J. Neyman, *'Smooth test' for goodness of fit*, Skand. Aktuarietidskr. 20 (1937), pp. 149–199.
- [16] — *Some memorable incidents in probabilistic/statistical studies*, in: *Asymptotic Theory of Statistical Tests and Estimation*, I. M. Chakravarti (Ed.), Academic Press, New York 1980, pp. 1–32.
- [17] A. V. Prohorov, *On sums of random vectors*, Theory Probab. Appl. 18 (1973), pp. 186–188.
- [18] J. C. W. Rayner and D. J. Best, *Smooth Tests of Goodness of Fit*, Oxford University Press, New York 1989.
- [19] — *Smooth tests of goodness of fit: An overview*, Internat. Statist. Rev. 58 (1990), pp. 9–17.
- [20] G. Schwarz, *Estimating the dimension of a model*, Ann. Statist. 6 (1978), pp. 461–464.
- [21] G. S. Watson, *Some recent results in chi-square goodness-of-fit tests*, Biometrics 15 (1959), pp. 440–468.

Faculty of Applied Mathematics
University of Twente
P.O. Box 217
7500 AE Enschede, The Netherlands

Institute of Mathematics
Technical University of Wrocław
Wybrzeże Wyspiańskiego 27
50-370 Wrocław, Poland

Received on 14.12.1993