

A SIGNIFICANCE TEST FOR COMPONENT ANALYSIS

BY PAUL G. HOEL

1. Introduction

During the last few years several papers and books have been written on various aspects of what has been termed component or factor analysis. This analysis has arisen from the psychological problem of describing the results on a series of tests in terms of a few distinct abilities or components. In much of such work it is claimed that there does not exist more than a certain number of components, the material discarded in order to substantiate such a claim being considered as due to random errors of sampling or errors of measurement. However, mere inspection of results or the calculation of standard errors of residual correlations is hardly sufficient to justify such conclusions, and therefore a significance test of some kind is necessary. Hotelling¹ considered such a test but based it upon an uncertain analogy with the analysis of variance and upon the legitimacy of using standard errors. The purpose of this paper is to derive a test which is more general in scope and in which all assumptions are explicitly stated.

If each test score is thought of as being made up of two parts, a true score and an error element, the assumption that there exists fewer components than the number of tests implies that the scatter diagram of the true scores will lie in a space of correspondingly smaller dimensionality. Consequently, an ideal test for the number of components would be one which would test the rank of the true moment matrix. In the case of normally distributed variables, this line of approach leads one to the sampling distribution of the generalized variance. Unfortunately, this distribution appears in unintegrated form; however, by considering its moments it is possible to find a good approximation to this exact distribution for samples which are not too small.

The paper proceeds by first finding two approximation distributions for the generalized variance, one for samples which are not too small and one for large samples. It then considers the type of population from which it will be assumed the sample was drawn, and finally applies the test to two numerical examples from recent literature along such lines.

2. Approximation Distributions

Suppose that N individuals have been drawn at random from an n variate normal population whose distribution is expressed by

$$(1) \quad P(x_1, x_2, \dots, x_n) = Ke^{-\sum_1^n \lambda_{ij} x_i x_j}$$

¹ Harold Hotelling, Analysis of a Complex of Statistical Variables into Principal Components, *The Journal of Educational Psychology*, September and October, 1933, pp. 21-25.

where $x_i = X_i - m_i$, $A_{ij} = \frac{\Delta_{ij}}{2\sigma_i\sigma_j\Delta}$, Δ is the determinant $|\rho_{ij}|$ and Δ_{ij} is the cofactor of ρ_{ij} in Δ , and $K = |A_{ij}|^{\frac{1}{2}}/(2\pi)^{n/2}$. If the observed values of the variables of the α th individual are denoted by $X_{i\alpha}$ ($i = 1, 2, \dots, n$), then the generalized sample variance is defined as $z = |a_{ij}|$, where $a_{ij} = \frac{1}{N} \sum_{\alpha=1}^N (X_{i\alpha} - \bar{X}_i)(X_{j\alpha} - \bar{X}_j)$. Wilks² has shown that in sampling from the population (1), the k th moment of the sampling distribution of z is given by

$$M_k = A^{-k} \frac{\Gamma\left(\frac{N+2k-1}{2}\right)\Gamma\left(\frac{N+2k-2}{2}\right)\dots\Gamma\left(\frac{N+2k-n}{2}\right)}{\Gamma\left(\frac{N-1}{2}\right)\Gamma\left(\frac{N-2}{2}\right)\dots\Gamma\left(\frac{N-n}{2}\right)}$$

where $A = N^n |A_{ij}|$. An inspection of the integrated form of the distribution of z in the case of $n = 1$ and $n = 2$ suggests that there likely exists a function of similar form for higher values of n whose k th moment can be made to differ from M_k only in higher powers of terms which contain N^{-1} as a factor. An investigation along such lines leads to the function

$$(2) \quad g(z) = Cz^m e^{-n\sqrt{az}}$$

where $C = \frac{a^{\frac{N-n}{2}} n^{\frac{N-n}{2}-1}}{\Gamma\left(n\frac{N-n}{2}\right)}$, $m = \frac{N-n-2}{2}$, $a = Aq$ and $q = 1 - \frac{(n-1)(n-2)}{2N}$.

It will be shown that the k th moment M'_k of $g(z)$ differs from M_k only in terms of magnitude less than the second and higher powers of k^2n/N or kn^2/N .

Multiplying $g(z)$ by z^k and integrating over the entire range of z will yield M'_k , which turns out to be

$$M'_k = \frac{\Gamma\left(n\frac{N-n+2k}{2}\right)}{a^k n^{nk} \Gamma\left(n\frac{N-n}{2}\right)}$$

Upon reducing the upper gamma function and performing successive steps of simple algebra

$$\begin{aligned} M'_k &= a^{-k} n^{-nk} \left(n\frac{N-n+2k}{2} - 1\right) \left(n\frac{N-n+2k}{2} - 2\right) \dots \left(n\frac{N-n}{2}\right) \\ &= N^{nk} a^{-k} 2^{-nk} \left(1 + \frac{2k-n-2/n}{N}\right) \left(1 + \frac{2k-n-4/n}{N}\right) \dots \\ &\quad \left(1 + \frac{2k-n-2kn/n}{N}\right). \end{aligned}$$

² S. S. Wilks, Certain Generalizations in the Analysis of Variance, *Biometrika*, Vol. XXIV, 1923, p. 477.

The terms in parentheses may be treated as the factored form of a polynomial of the nk th degree in unity. Thus the quantities $\frac{2k - n - 2/n}{N}$, etc., may be treated as the zeros with signs changed of the corresponding polynomial in x (say). As a result, the successive terms after the first in the non-factored form of this polynomial in unity are the sums of the products of these quantities taken one at a time, two at a time, etc. Upon performing this multiplication and letting $\phi = N^n/2^n A$, M'_k assumes the form

$$M'_k = \phi^k q^{-k} \left[1 - \frac{k(n^2 - nk + 1)}{N} + \dots \right]$$

where the neglected terms are in magnitude less than the second and higher powers of $k^2 n/N$ or kn^2/N . If M_k is handled in exactly the same manner, it will be found that

$$\begin{aligned} M_k &= A^{-k} \left(\frac{N + 2k - 1}{2} - 1 \right) \dots \left(\frac{N + 2k - 1}{2} - k \right) \dots \\ &\quad \left(\frac{N + 2k - n}{2} - 1 \right) \dots \left(\frac{N + 2k - n}{2} - k \right) \\ &= N^{nk} A^{-k} 2^{-nk} \left(1 + \frac{2k - 3}{N} \right) \dots \left(1 - \frac{1}{N} \right) \dots \\ &\quad \left(1 + \frac{2k - n - 2}{N} \right) \dots \left(1 - \frac{n}{N} \right) \\ &= \phi^k \left[1 - \frac{nk(n - 2k + 3)}{2N} + \dots \right] \end{aligned}$$

where the neglected terms are of the same order of magnitude as those neglected in the approximation to M'_k . Before a comparison of M_k and M'_k is possible, the factor q^{-k} of M'_k must be expanded and multiplied into the quantity in brackets. This operation yields the result

$$M'_k = \phi^k \left[1 - \frac{nk(n - 2k + 3)}{2N} + \dots \right].$$

Thus M_k and M'_k agree to within neglected terms. As a matter of fact, if the values of the neglected terms are considered more carefully, it will be found that the actual difference between M_k and M'_k is considerably less than the given upper bound for the magnitude of neglected terms would indicate. For example, when $n = 5$ the first term in the difference is $6k(k - .9)N^{-2}$, while $625k^2N^{-2}$ or $25k^4N^{-2}$ is the upper bound for this term when only general results are used. The general formula for the first term in this difference has been obtained, but since the remaining terms have not been investigated and since the type of problems to which the distribution $g(z)$ is to be applied does not

justify this refinement, it will not be considered here. Consequently, if one considers this distribution function as sufficiently determined by its low order moments and if one applies $g(z)$ only to problems in which N is fairly large compared with n^2 , then the function $g(z)$ will give a good approximation to the exact sampling distribution of z . Obviously, $g(z)$ is identical with the exact distribution for the known cases of $n = 1$ and $n = 2$. It is not possible under the above expansions to vary the constants in the form of $g(z)$ in such a manner as to obtain an approximation whose k th moment will agree with M_k to within still higher powers of comparable terms.

In order to test whether or not a sample value $z = Z$ can be reasonably assumed to have been obtained in random sampling from a population of type (i) with fixed A , it is necessary to calculate the probability P of obtaining in repeated samples a value of z greater than Z . Thus it is necessary to evaluate

$$P = 1 - \int_0^z g(z) dz.$$

Upon making the substitution $x = n\sqrt{az}$, and letting $p = n\frac{N-n}{2} - 1$ and $u = n\sqrt{aZ}\left(n\frac{N-n}{2}\right)^{-1} = nN\sqrt{\frac{Z}{\phi}\left[1 - \frac{(n-1)(n-2)}{2N}\right]} [2n(N-n)]^{-1}$, this integral can be reduced to the standard form of the incomplete gamma function. Hence P assumes the form

$$(3) \quad P = 1 - I(u, p)$$

where

$$I(u, P) = \frac{1}{\Gamma(p+1)} \int_0^{u\sqrt{p+1}} e^{-x} x^p dx.$$

In many applications of this distribution it will be found that the values of u and p lie beyond the tabled³ values of these constants. Consequently, it will often be sufficient to use the normal distribution to which the gamma distribution tends as N becomes large. This normal distribution will be considered next.

Rather than obtain a normal approximation to $g(z)$ or the gamma function to which $g(z)$ reduces after the above transformation, it is more illuminating to find the basic descriptive parameters of the exact distribution of z and from them obtain a normal approximation. Such a procedure will show how rapidly the distribution of z approaches normality with increasing N . By using the recurrence formula connecting M_{k+1} and M_k , which can be found directly from the ratio of these two moments, and expressing the necessary moments in

³ K. Pearson, Tables of the Incomplete Gamma Function, Biometric Laboratory (1922), Univ. of London.

terms of M_1 , it can be shown that these basic descriptive parameters are expressible in expanded form as follows:

$$\begin{aligned}
 m &= \phi \left[1 - \frac{n(n+1)}{2N} + \frac{n(n+1)(n-1)(3n+2)}{24N^2} + \dots \right] \\
 \sigma^2 &= \phi^2 \left[\frac{2n}{N} - \frac{n(2n^2 - n + 1)}{N^2} + \dots \right] \\
 \beta_1 &= \frac{2(3n-1)^2}{nN} \left[1 - \frac{(n+1)(5n-3)}{2(3n-1)N} + \dots \right] \\
 \beta_2 &= 3 \left[1 + \frac{4(3n-1)(4n-1)}{3nN} + \dots \right].
 \end{aligned}$$

These values suggest that

$$(4) \quad w = \sqrt{\frac{N}{2n}} \left[\frac{z}{\phi} - 1 \right]$$

will likely be distributed approximately normally with zero mean and unit variance. As a matter of fact, by using the second limit theorem of probability,⁴ it can be shown that the distribution of w approaches normality as N increases indefinitely. Hence, for samples in which N is large compared with n^2 , it will be sufficient to compare the value of w arising from a sample $z = Z$ with its variance of unity if a test of significance is desired. A better general approximation could have been obtained by centering the curve at $\phi \left[1 - \frac{n(n+1)}{2N} \right]$ rather than at ϕ ; however, since there is positive skewness and the true mean lies between these two values, there might arise some exaggeration in a significance test in doing so because the accuracy of such a test depends upon the accuracy of the approximation in the right hand tail of the curve.

Inspection of (3) and (4) shows that the only population parameter upon which these approximation distributions depend is ϕ . There are no assumptions necessary about the population means, or variances, or covariances, except in so far as they may be related when the value of ϕ is postulated. This means that either (3) or (4) enables one to test whether or not it is reasonable to assume that the sample variance $z = Z$ arose in random sampling from some normal population with ϕ equal to the postulated value.

3. Population Assumptions

Consider the set of variables u_1, u_2, \dots, u_n distributed according to the normal law

$$(5) \quad P(u_1, u_2, \dots, u_n) = K_1 e^{-\sum_1^n b_{ij} u_i u_j}$$

⁴ See, for example, Frechet and Shohat, A Proof of the Generalized Second Limit Theorem in the Theory of Probability, Transactions of the American Mathematical Society, Vol. 33, (1931), p. 533.

and the set of variables v_1, v_2, \dots, v_n distributed according to the normal law

$$(6) \quad P(v_1, v_2, \dots, v_n) = K_2 e^{-\sum_1^n c_i v_i^2}$$

where the v 's are uncorrelated with the u 's and with each other. The joint distribution of the u 's and v 's is expressed by

$$(7) \quad P(u_1, \dots, v_n) = K_3 e^{-\sum_1^n b_{ij} u_i u_j - \sum_1^n c_i v_i^2}$$

Upon writing down the determinant of the coefficients of these $2n$ variables, it will become evident that any one of its principal minors of any order can be expressed as the product of a principal minor of $|b_{ij}|$ with a principal minor of $|c_i|$. Since the distributions (5) and (6) are normal, the determinants $|b_{ij}|$ and $|c_i|$ are positive definite; consequently the determinant of the coefficients in (7) must also be positive definite.

Now consider the orthogonal transformation

$$y_i = \frac{u_i + v_i}{\sqrt{2}}, \quad i = 1, 2, \dots, n$$

$$y_i = \frac{u_i - v_i}{\sqrt{2}}, \quad i = n + 1, \dots, 2n.$$

Since the determinant of the coefficients in (7) is invariant under an orthogonal transformation, the resulting distribution of the y 's may be expressed by

$$(8) \quad P(y_1, y_2, \dots, y_{2n}) = K_4 e^{-\sum_1^{2n} d_{ij} y_i y_j}$$

where $|d_{ij}|$ is positive definite.

In order to obtain the distribution of the variables y_1, y_2, \dots, y_n , it is necessary to integrate (8) with respect to the variables y_{n+1}, \dots, y_{2n} over their range of values. If this integration is performed after the quadratic form in the exponent of (8) has been expressed as a sum of squares⁵ with coefficients which are the ratios of principal minors of $|d_{ij}|$, it will be clear that the integration leaves a quadratic form in the exponent which is also positive definite. Hence after the transformation $x_i = \sqrt{2}y_i (i = 1, 2, \dots, n)$ the distribution function of the variables $x_i = u_i + v_i (i = 1, 2, \dots, n)$ must be normal and may be expressed by (1). Thus it has been shown that if the true parts u_i of the variables x_i are normally distributed without error and if the error parts v_i are normally distributed but are uncorrelated with the u_i and with each other, then the variables x_i possess a normal distribution. The advantage of

⁵ See, for example, Risser and Traynard, *Les Principes de la Statistique Mathématique*, 1933, p. 225.

this formulation will become evident when the parameter ϕ is expressed in terms of the parameters of (5) and (6).

Since the v 's are uncorrelated with the u 's and with each other, the variance σ_i^2 of x_i is the sum of the variances of u_i and v_i , while the correlation ρ_{ij} between x_i and x_j may be expressed in terms of the correlation ρ'_{ij} between u_i and u_j and the variances $u_i^2, u_j^2, v_i^2, v_j^2$ of u_i, u_j, v_i, v_j respectively. These relationships are

$$(9) \quad \sigma_i^2 = \mu_i^2 + \nu_i^2, \quad \text{and} \quad \rho_{ij} = \frac{\rho'_{ij}}{\sqrt{(1 + \nu_i^2/\mu_i^2)(1 + \nu_j^2/\mu_j^2)}} \quad (i \neq j).$$

For simplicity of notation let $\lambda_i = \nu_i^2/\mu_i^2$. Now it is well known⁶ that ϕ can be expressed in the form

$$\phi = \sigma_1^2 \sigma_2^2 \cdots \sigma_n^2 | \rho_{ij} |.$$

If the values from (9) are inserted in $| \rho_{ij} |$ and if the resulting denominators of elements are factored out, ϕ will assume the form

$$\phi = \frac{\sigma_1^2 \sigma_2^2 \cdots \sigma_n^2 B}{(1 + \lambda_1) \cdots (1 + \lambda_n)}$$

where

$$B = \begin{vmatrix} 1 + \lambda_1 & \rho'_{12} & \cdots & \rho'_{1n} \\ \rho'_{12} & & & \vdots \\ \vdots & & & \vdots \\ \rho'_{1n} & \cdots & \cdots & 1 + \lambda_n \end{vmatrix}.$$

Following the methods of confluence analysis,⁷ B can be expressed as follows:

$$B = R + \sum_{\alpha=1}^n \lambda_\alpha R_{\alpha\alpha} + \sum_{\alpha < \beta} \lambda_\alpha \lambda_\beta R_{\alpha\beta\alpha} + \cdots + \lambda_1 \lambda_2 \cdots \lambda_n$$

where $R = | \rho'_{ij} |$, $R_{\alpha\alpha}$ is the principal minor of R obtained by deleting row and column α , etc. R is the true correlation determinant whose rank it is the object of this paper to test. If R is assumed to be of rank $n - t$, then all principal minors containing more than $n - t$ rows vanish and B reduces to

$$B = \sum_{\alpha_1 < \cdots < \alpha_t} \lambda_{\alpha_1} \lambda_{\alpha_2} \cdots \lambda_{\alpha_t} R_{\alpha_1 \alpha_2 \cdots \alpha_t} + \cdots + \lambda_1 \lambda_2 \cdots \lambda_n.$$

The tests (3) and (4) were designed to test hypothetical values of ϕ by means of the sample Z . Evidently the value of ϕ can be postulated by assigning hypothetical values to the λ 's, the σ 's, and the principal minors of R .

Assigning values to the λ 's does not curtail the degrees of freedom in these

⁶ S. S. Wilks, loc. cit., p. 477.

⁷ Ragnar Frisch, *Statistical Confluence Analysis by Means of Complete Regression Systems*, Oslo, 1934.

tests because they were derived on the basis of (1) which depends only on the m 's, σ 's, and ρ 's. The λ 's do restrict the range of the ρ 's, but not their degrees of freedom.

An inspection of the expression for ϕ shows that ϕ can be made to assume any desired value irregardless of the rank of R by merely assigning the σ 's properly. It is therefore necessary to make some assumption regarding the σ 's if the test is to serve the purpose for which it is intended. Here it will be sufficient to assume that the product of the population variances may be replaced by the product of the sample variances. This assumption will ordinarily be approximately fulfilled for the size samples for which it is legitimate to employ (3) or (4); consequently this assumption does not restrict the range of application of the test.

To postulate values of the principal minors of R beyond postulating the rank of R would introduce hypotheses and restrictions which are irrelevant to the fundamental purpose of the test. This difficulty will be avoided by replacing all non-vanishing minors of R by their upper bounds of unity. Since this will overestimate the value of B , and hence of ϕ , the usual significance level of .05 may be considered as decisive. Let the value of B when unity is inserted for all non-vanishing principal minors be denoted by D . Then

$$(10) \quad D = \sum_{\alpha_1 < \dots < \alpha_t}^n \lambda_{\alpha_1} \lambda_{\alpha_2} \dots \lambda_{\alpha_t} + \dots + \lambda_1 \lambda_2 \dots \lambda_n.$$

Since

$$\prod_1^n (1 + \lambda_i) = 1 + \sum_{\alpha=1}^n \lambda_\alpha + \sum_{\alpha_1 < \alpha_2}^n \lambda_{\alpha_1} \lambda_{\alpha_2} + \dots + \lambda_1 \lambda_2 \dots \lambda_n$$

it will often be convenient to write D in the form

$$(11) \quad D = \prod_1^n (1 + \lambda_i) - \left\{ 1 + \sum_{\alpha=1}^n \lambda_\alpha + \dots + \sum_{\alpha_1 < \dots < \alpha_{t-1}}^n \lambda_{\alpha_1} \lambda_{\alpha_2} \dots \lambda_{\alpha_{t-1}} \right\}.$$

As a consequence of all the above assumptions,

$$(12) \quad \frac{Z}{\phi} = \frac{|a_{ij}|}{\phi} = \frac{(1 + \lambda_1) \dots (1 + \lambda_n) |r_{ij}|}{B} \\ \cong \frac{(1 + \lambda_1) \dots (1 + \lambda_n) |r_{ij}|}{D}$$

where $|r_{ij}|$ is the sample correlation determinant.

All the essential material for testing the rank of the true correlation matrix is contained in (3), (4), (11), and (12). In summary, the hypothesis to be tested and the procedure to follow in performing the test are as follows.

The population of n variables from which the sample is supposed drawn is assumed to be such that (a) the true parts of the variables are normally distributed, (b) the error parts are normally distributed but are uncorrelated with the true parts and with each other, (c) the product of the variances may be replaced by the product of the sample variances, (d) the values of the λ 's

are postulated as judged by the accuracy in measurement of the variables, and (e) the rank of the true correlation matrix is $n - t$.

Given the value $|r_{ij}|$ of the sample correlation determinant, a lower bound for the value of Z/ϕ is calculated from (11) and (12). This lower bound is inserted in either (3) or (4), depending on the size of the sample. If (3) is used and if $P \leq .05$, or if (4) is used and $w \geq 2$, one may conclude, as judged by the sample variance, that it is very unlikely that the sample was drawn in random sampling from the population specified above. If one has reason to believe that the variables are sensibly normal as indicated above and that the postulated values of the λ 's are quite accurate, then the test shows quite definitely that the postulated rank of the true correlation matrix is unsubstantiated by the sample, and therefore a higher rank should be tested until a non-significant value is obtained. Because a lower bound rather than the value of Z/ϕ is used, the test can be used on minimum ranks only, and hence a value of $Z < \phi$ will not yield a test of significance. However, the test does handle the problem for which it was designed and which is of fundamental interest, and that is to see whether or not one is justified in assuming that a sample represents only a certain minimum number of components.

4. Applications

(a) Hotelling⁸ has used an example taken from other sources to illustrate his test on components. In order to compare results, this same example will be treated here under the assumptions outlined above. In this example the reliability coefficients are given. From the definition of a reliability coefficient r_i , it follows at once that $r_i = \frac{1}{1 + \lambda_i}$. The population values of the λ 's will be set equal to the values obtained from these sample reliability coefficients. The data for this problem are

$$|r_{ij}| = .235, N = 140, n = 4, \lambda_1 = .087, \lambda_2 = .119, \lambda_3 = .101, \lambda_4 = .773.$$

Assume that the true correlation matrix in the population is of rank two, that is, that two components are sufficient to describe the results on these tests. Since N is large compared with n^2 , it will be sufficient to use (4). The values of (11), (12), and (4) are found to be

$$D = \prod_1^4 (1 + \lambda_i) - \left\{ 1 + \sum_1^4 \lambda_i \right\} = .294$$

$$\frac{Z}{\phi} \geq \frac{\prod (1 + \lambda_i) |r_{ij}|}{D} = 1.90$$

$$w \geq \sqrt{\frac{140}{8}} [1.90 - 1] = 3.76$$

⁸ Loc. cit., p. 16.

Since the standard deviation of w is unity, this value demonstrates clearly that the hypothesis of only two components is untenable as judged by the sample correlation determinant. If one assumes three components, the test will be found to yield a non-significant value. Hence it may be concluded that under the hypotheses on which the test is based, the sample does not justify the assumption of less than three components. Hotelling's test indicated the necessity for two components but was uncertain about the third, the decision resting upon a variate value of 1.31 as against a standard deviation of unity.

(b) Thurstone, in his "Vectors of Mind," considers an example taken from a series of fifteen psychological tests. After applying his centroid method to the data, he inspects his results and concludes that four components are sufficient to account for everything except random errors. It is impossible to test his conclusions explicitly as above because the size of the sample is not given and the reliability coefficients are not known. Nevertheless, if it is legitimate to assume that the sample is sufficiently large to justify the use of this test, interesting conclusions can be obtained on the assumption that only four components are needed.

Suppose that $\lambda_i = \frac{1}{2}$, which implies that the variance of error is half as large as the true sampling variance for each variable. Here (10) is more convenient than (11) for computing the value of D . The values of (10) and (12) are found to be

$$D = {}_{15}C_3\left(\frac{1}{2}\right)^{12} + {}_{15}C_2\left(\frac{1}{2}\right)^{13} + {}_{15}C_1\left(\frac{1}{2}\right)^{14} + \left(\frac{1}{2}\right)^{15} = .125$$

$$\frac{Z}{\phi} \geq \frac{|r_{ij}|}{.0003}.$$

Evidently, the value of $|r_{ij}|$ must lie in the neighborhood of .0003 if the test is not to yield a significant result which contradicts the hypothesis. However, the correlations in $|r_{ij}|$ are given to only three decimal places, and therefore a legitimate value in the neighborhood of .0003 can not be realized. It is to be noted that the postulated values of the λ 's are equivalent to postulating that all reliability coefficients are equal to $\frac{2}{3}$, a value which should be considered as unusually low. It would seem reasonable to avoid using material in which the variance of error is larger than one-half the variance of random sampling, unless the variance of random sampling is exceedingly small.