

CENTRO PER LA FORMAZIONE IN ECONOMIA E POLITICA DELLO SVILUPPO RURALE
UNIVERSITÀ DI NAPOLI FEDERICO II - DIPARTIMENTO DI SCIENZE STATISTICHE
UNIVERSITÀ DI SALERNO - DIPARTIMENTO DI SCIENZE ECONOMICHE E STATISTICHE

Quaderni di STATISTICA

VOLUME 9 - 2007

LIGUORI EDITORE

Volume 9, anno 2007

ISSN 1594-3739 (edizione a stampa)

Registrazione al n. 5264 del 6/12/2001 presso il Tribunale di Napoli
ISBN-13 978 - 88 - 207 - 4209 - 6

Direttore responsabile: Gennaro Piccolo

© 2007 by Liguori Editore

Tutti i diritti sono riservati

Prima edizione italiana Dicembre 2007

Finito di stampare in Italia nel mese di Dicembre 2007 da OGL - Napoli

Questa opera è protetta dalla Legge sul diritto d'autore (Legge n. 633/1941).

Tutti i diritti, in particolare quelli relativi alla traduzione, alla citazione, alla riproduzione in qualsiasi forma, all'uso delle illustrazioni, delle tabelle e del materiale software a corredo, alla trasmissione radiofonica o televisiva, alla registrazione analogica o digitale, alla pubblicazione e diffusione attraverso la rete Internet sono riservati, anche nel caso di utilizzo parziale.

La riproduzione di questa opera, anche se parziale o in copia digitale, è ammessa solo ed esclusivamente nei limiti stabiliti dalla Legge ed è soggetta all'autorizzazione scritta dell'Editore. La violazione delle norme comporta le sanzioni previste dalla legge.

Il regolamento per l'uso dei contenuti e dei servizi presenti sul sito della Casa Editrice Liguori è disponibile al seguente indirizzo: http://www.liguori.it/politiche_contatti/default.asp?c=legal

L'utilizzo in questa pubblicazione di denominazioni generiche, nomi commerciali e marchi registrati, anche se non specificamente identificati, non implica che tali denominazioni o marchi non siano protetti dalle relative leggi o regolamenti.

Il C.F.E.P.S.R. si avvale per la stampa dei Quaderni di Statistica del contributo dell'Istituto Banco di Napoli - Fondazione.

La carta utilizzata per la stampa di questo volume è inalterabile, priva di acidi, a pH neutro, conforme alle norme UNI EN Iso 9706 X, realizzata con materie prime fibrose vergini provenienti da piantagioni rinnovabili e prodotti ausiliari assolutamente naturali, non inquinanti e totalmente biodegradabili.

Indice

R. ARBORETTI GIANCRISTOFARO, S. BONNINI, L. SALMASO, A performance indicator for multivariate data	1
D. PICCOLO, A general approach for modelling individual choices ...	31
S. M. PAGNOTTA, The behavior of the sphericity test when data are rank transformed	49
A. PALLINI, On variance reduction in some Bernstein-type approxi- mations	63
A. NACCARATO, Full Information Least Orthogonal Distance Esti- mator of structural parameters in simultaneous equation models	87
M. CORDUAS, Dissimilarity criteria for time series data mining	107
 FORUM	
S. PACILLO, Estimation of ARIMA models under non-normality	133
M. IANNARIO, A statistical approach for modelling Urban Audit Perception Surveys.....	149

The behavior of the sphericity test when data are rank transformed

Stefano M. Pagnotta

Dipartimento Persona, Mercato, Istituzioni, Università degli Studi del Sannio
E-mail: pagnotta@unisannio.it

Summary: In this paper we give empirical evidence of the behavior of the sphericity test when original data are transformed in ranks. The study is performed by an extensive Montecarlo simulation. Specifically the type I error probabilities the powers under different alternative hypotheses are evaluated. Also consider the robustness of the test when the population is not multivariate Gaussian distributed is investigated. Finally the selection of the principal components is discussed.

Keywords: Sphericity test; Rank transformation; Empirical asymptotic distribution.

1. Introduction

Principal Component Analysis (PCA) is a standard methodological tool adopted when a large set of p numerical variables $X_k, k = 1, 2, \dots, p$, is available. The original data are linearly transformed so that new variables $Y_j, j = 1, 2, \dots, p$, are computed. The Y_j 's are mutually uncorrelated and ordered according to their variances, i.e. $\text{var}[Y_1] \geq \text{var}[Y_2] \geq \dots \geq \text{var}[Y_p]$; moreover the identity $\sum_{j=1}^p \text{var}[Y_j] \equiv \sum_{k=1}^p \text{var}[X_k]$ is satisfied. The main problem of PCA is the selection of the principal components Y_j 's in order to reduce the p dimension of the original data to a lower one so that most part of $\sum_{k=1}^p \text{var}[X_k]$ is preserved. In literature many rules of selection (see for example Jolliffe (2002), chap. 6) are suggested and one of them is based on the sphericity test when the data are assumed to be drawn from a multivariate Gaussian probability law.

The test of sphericity concerns the null hypothesis $\mathcal{H}_0 : \Sigma = \sigma^2 \mathbf{I}$, where Σ is the covariance matrix a p -variate Gaussian population. Specifically, \mathbf{I} is the identity matrix and σ^2 is a positive value. The test-statistic was first derived by using the generalized likelihood ratio method (Mauchly, 1940). It involves the determinant and the trace of the ML estimate of the covariance matrix $\hat{\Sigma} = \frac{1}{n} \sum_i (\mathbf{x}_i - \bar{\mathbf{x}}) (\mathbf{x}_i - \bar{\mathbf{x}})^t$, where the $\mathbf{x}_i, i = 1, 2, \dots, n$, are p -dimensional column-vectors, $\bar{\mathbf{x}} = \frac{1}{n} \sum_i \mathbf{x}_i$ is the sample mean and n the sample size. Both the determinant and the trace of $\hat{\Sigma}$ can be evaluated starting from its eigenvalues $\hat{l}_j, j = 1, 2, \dots, p$; hence the test-statistics simplifies to

$$\lambda = \left(\frac{\frac{1}{p} \sum_{j=1}^p \hat{l}_j}{\sqrt[p]{\prod_{j=1}^p \hat{l}_j}} \right)^{-np/2}. \quad (1)$$

Mauchly provided the critical values of the test for $p = 2$ and 3 while other authors (Pillai and Nagarsenker, 1971; Nagarsenker and Pillai, 1973; Marques and Coelho, 2007) studied the sampling distribution of

$$W = -2 \log \lambda = np \log \left(\frac{\frac{1}{p} \sum_{j=1}^p \hat{l}_j}{\sqrt[p]{\prod_{j=1}^p \hat{l}_j}} \right). \quad (2)$$

The application of the Wilks theorem assures that W converges in distribution, as n increases, to a chi-square probability law with $q = (p - 1)(p + 2)/2$ degrees of freedom.

In order to use the sphericity test to select the optimal number of principal components, the null hypothesis has to be formulated as $\mathcal{H}_0 : l_1 = l_2 = \dots = l_{p-1} = l_p$, where $l_j, j = 1, 2, \dots, p$, are the eigenvalues of Σ . It can be shown that l_j is the variance $\text{var}[Y_j]$ of the population principal component $Y_j, j = 1, 2, \dots, p$, while $\hat{l}_j = \text{var}[\hat{Y}_j]$ is the variance of the corresponding sampling version. When \mathcal{H}_0 is not rejected then the principal component analysis is not useful for dimensionality reduction, whereas if the null hypothesis is rejected at least one of the Y_j 's retains much more variability than the others. Given that $l_1 \geq l_2 \geq \dots \geq l_p$, it is more likely that \mathcal{H}_0 is rejected because $l_1 > l_2 = \dots = l_{p-1} = l_p$. In this case a second hypothesis is considered $\mathcal{H}_0^{(1)} : l_2 = \dots = l_p$ where the

first eigenvalue l_1 is left out. Now if $\mathcal{H}_0^{(1)}$ is not rejected the only principal component retained in the analysis is \hat{Y}_1 , corresponding to the first eigenvalue l_1 . When $\mathcal{H}_0^{(1)}$ is rejected, the likely event is $l_2 > l_3 = \dots = l_p$ and the further hypothesis $\mathcal{H}_0^{(2)} : l_3 = \dots = l_p$ has to be tested. If we set $\mathcal{H}_0^{(0)} \equiv \mathcal{H}_0$ and consider the integer parameter $k = 0, 1, \dots, p - 2$, to each hypothesis $\mathcal{H}_0^{(k)} : l_{k+1} = l_{k+2} = \dots = l_{p-1} = l_p$ corresponds the test statistic:

$$W^{(k)} = n \cdot (p - k) \cdot \log \left(\frac{\frac{1}{p-k} \sum_{j=k+1}^p \hat{l}_j}{\sqrt[p-k]{\prod_{j=k+1}^p \hat{l}_j}} \right),$$

that is asymptotically distributed as a chi-square random variable with

$$q^{(k)} = (p - k + 2)(p - k - 1)/2 \quad (3)$$

degrees of freedom. For $k = 0$ we have the original test of sphericity.

In order to accelerate the convergence to the asymptotic distribution under the hypothesis that the population is a p -variate Gaussian, Bartlett (1937) suggested to replace n by $n - \frac{2p+1}{6}$. This modification introduced by Bartlett does not change the degree of freedom of the asymptotic distribution. Through the paper we then consider the following test statistic:

$$W^{(k)} = (p - k) \left(n - \frac{2p + 1}{6} \right) \log \left(\frac{\frac{1}{p-k} \sum_{j=k+1}^p \hat{l}_j}{\sqrt[p-k]{\prod_{j=k+1}^p \hat{l}_j}} \right). \quad (4)$$

When the population is not gaussian the significance level of the test generally degenerates. As a matter of fact, the test statistics $W^{(0)}$ for the hypothesis $\mathcal{H}_0 : \Sigma = \sigma^2 \mathbf{I}$ holds the nominal level of significance only when the population is a multivariate random variable whose fourth cumulants are zero as shown by Waternaux (1984). A part from this case the sphericity test is unreliable.

Data transformations are often considered in order to promote normality and/or to control large errors in the components of one or more observations. One of the non parametric transformation routinely applied to multidimensional data consists in replacing each of the observations of

a variable with its position in the ordered set of values. This scheme of rank transformation is labeled RT-2 by Conover and Iman (1981).

The rank transformation (RT) is suggested by different authors; some of them propose the RT just as one of the possible choice of data-transformations (Mardia et. al, 1979, pp. 235-; Jambu, 1989, p. 126); while some others provide justifications to adopt it (Jobson, 1992, p. 387; Baxter, 1995). In any case no study has been performed to compare the effect of the application of the RT-2 to data matrix in the field of PCA. When this transformation is adopted, only some information about the original probability law of the population can be restored from the data. Nothing of the original marginal probability law of the population is preserved, in fact the empirical distribution function of a RT observed variable reduces to a straight line starting from 0 to 1 in the range $[1, n]$. We will refer to this consequence of the RT on a single variable as flattening effect. Instead Borkowf (2002) has analytically proved that the Pearson correlation index is underestimated by the Spearman rank correlation coefficient, at least in the case of the bivariate Gaussian population. For multivariate case no theoretical result is available, but the proof of Borkowf supports the empirical evidence of a similar result in the case of a multivariate Gaussian population. At the end only some information about the population is preserved in the covariance structure. For this reason no automatic inferential procedure is allowed, although Conover and Iman (1981) suggest to transform the data in ranks and then apply the standard parametric analysis. In some cases this procedure can give reliable results, as found by Nath and Pavur (1985) for MANOVA. However Headrick et al. (2001) show that in the framework of the multiple regression the inferential procedures fail when data are previously rank transformed. Thompson (1991) suggests that the inferential output of procedures applied to rank transformed data has to be used with extreme caution unless the distributional properties of the parametric tools applied to RT data are investigated.

When a dataset is rank transformed according the RT-2 scheme, we will find (see section 3) that the asymptotic distribution of the sphericity test statistics $W^{(0)}$ changes with respect to the degrees of freedom.

Consequently the *p-value* of the statistic (provided routinely by statistical softwares) is misleading.

In this paper we recompute the degrees of freedom and study the properties of the test statistic by simulation. We will show that the sphericity test is reliable both when the data are from a multivariate Gaussian and when the data matrix is previously RT-2, whenever the marginals of the variable are correlated and the degrees of freedom of the asymptotic sampling distribution are redefined. When the marginal distributions are independent, the sphericity test has a very low power. The sphericity test applied to RT-2 data works fine when the population is not a multivariate Gaussian. With respect to the selection of the principal components when the data matrix is RT-2, we can affirm that the sequential use of the sphericity test is unable to correctly select the components.

This paper contains three more sections. In the next we recompute the degree of freedom of the test statistic (4) taking into account that the data are RT-2. In section 3 the design and the output of the simulations are presented. Some concluding remarks are reported in section 4.

2. Theoretical and empirical remarks

In this section we show how the asymptotic distribution of (4), for $k = 0$, changes when the data are drawn from a multivariate Gaussian and then transformed according to the RT-2 scheme. The literature does not give any theoretical support to this changing, but the result of a preliminary Montecarlo simulation is promising in this direction. For this purpose 10000 samples of 210 observations have been drawn from a 5-variate Gaussian distribution, with $\mu = \mathbf{0}$, $\Sigma = \mathbf{I}$. For each sample the values of (4), for $k = 0$, has been computed by estimating the eigenvalues of $\hat{\Sigma}$ both by using the raw data, and the RT-2 transformed data, say W_{Nor} and W_{RT2} respectively. Finally 30 equally spaced quantiles, from 0.01 to 0.99, are plotted against the theoretical ones from a chi-square distribution with $q = 14$ (see Figure 1). The triangles describe the agreement of the empirical quantiles of the untransformed gaussian samples with the theoretical ones, while the squares are the empirical quantiles when the samples are RT-2 transformed. Both trends are linear but the squares are

shifted below the case of the normal samples. This configuration suggests that there is only a change in the degrees of freedom of the empirical distribution of (4) when the data are RT-2.

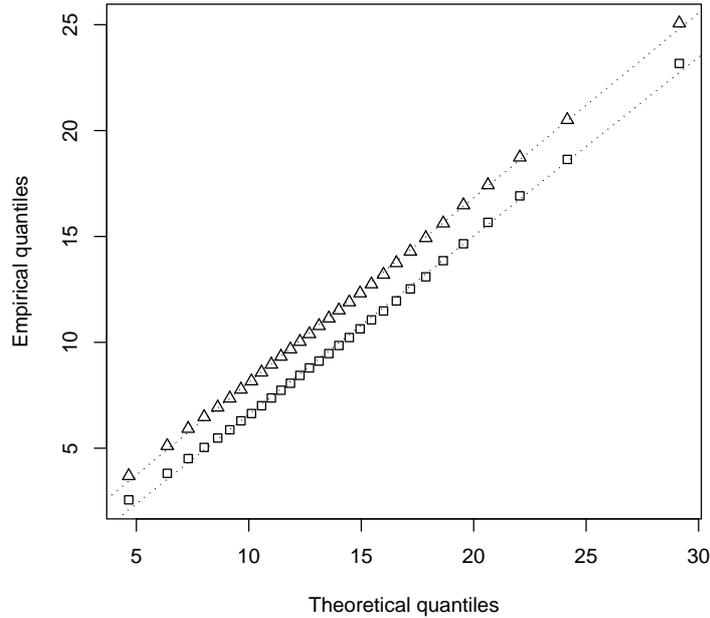


Figure 1. Empirical quantiles of W_{Nor} (triangles) and W_{RT2} (squares) plotted versus the theoretical quantiles of a chi-square distribution with 14 degrees of freedom.

The test statistics (2), as specified in the previous section, is given by the joint use of the generalized likelihood ratio method and the Wilks theorem when the null hypothesis is $\mathcal{H}_0 : \Sigma = \sigma^2 \mathbf{I}$, the population is a multivariate gaussian and then the degrees of freedom of the asymptotic chi-square distribution of (2) are given by the difference between the total number of parameters to be estimated and the number of parameters under the null hypothesis. For the sphericity test there are p parameters corresponding to the p means and $p(p+1)/2$ for the variances and the

covariances in Σ . The dimension of the space of the parameters under the null hypothesis is $p + 1$; since only p means and the variance σ^2 are estimated, given $\Sigma = \sigma^2 \mathbf{I}$.

When the data are rank transformed according to the RT-2 column-wise scheme, the number of parameters to be estimated changes. It is no longer necessary to estimate¹ the mean of each variable because it is equal to $(n + 1) / 2$, while the variances are equal to $(n^2 - 1) / 12$. It follows that the degrees of freedom, in this case, decrease to $q_{RT2} = p(p - 1) / 2$ that corresponds to the covariances in $\hat{\Sigma}$ out of the null hypothesis.

In order to use the sphericity test for selecting the principal components when data are RT, q_{RT2} has to be parameterized with respect to $k = 0, 1, \dots, p - 2$, for testing the sequence of hypothesis $\mathcal{H}_0^{(k)}$. After some algebra, it can be shown that if the data are RT-2 transformed, the degrees of freedom of the asymptotic distribution of (4) become

$$q_{RT2}^{(k)} = (p - k) \cdot (p - k - 1) / 2. \quad (5)$$

3. Empirical investigation

In this section we present the Montecarlo simulations which are useful to evaluate the reliability of the sphericity test.

The Montecarlo simulations we run share common factors. The data matrix \mathbf{X} has dimensions $[n, p]$, with $n = 15 \cdot q^{(0)}$, being $q^{(0)} = (p - 1)(p + 2) / 2$ the total number of the free parameters to be estimated, and $p = 3, 5, 10, 20, 50$. The values in \mathbf{X} are drawn by a p -normal distribution with zero mean vector and $\Sigma = \mathbf{I}$. Moreover the number of the replications is $M = 10000$. For each observed data-matrix \mathbf{X}_m , $m = 1, 2, \dots, M$, the eigenvalues are computed before (the *raw data* columns) and after the rank transformation (the *RT-2 data* columns) and then used to compute the sampling distribution of (4), for $k = 0$. Table 1 shows $\hat{\alpha}$, the estimated empirical level of significance of the sphericity test with $\alpha = 0.05$. The $\hat{\alpha}$'s of the *raw data* are computed with respect to a chi-square distribution with $q^{(0)}$ degrees of freedom. The other values, as $\hat{\alpha}$ for the *RT-2 data*,

¹ Here we assume that raw data are drawn by a continuous random variable where no ties are allowed.

are instead evaluated with respect to a chi-square distribution with $q_{RT2}^{(0)}$ degrees of freedom. The values of the $\hat{\alpha}$'s, both for raw and RT2 data, show that the theoretical level of significance holds (for each of the $\hat{\alpha}$'s, the test $\mathcal{H}_0 : \alpha = 0.05$ is performed).

Furthermore, in Table 1 we report the estimated power of the test, $\hat{\pi}(\cdot)$, evaluated with respect to three different alternative hypotheses $\mathcal{H}_a : \Sigma \neq \sigma^2 \mathbf{I}$, $a = 1, 2, 3$. The alternative hypotheses differ for the structure of the matrix Σ and try to mirror frequent cases met in analyzing data with principal components.

Under \mathcal{H}_1 , Σ is a diagonal matrix where:

$$\Sigma_{11} = 0.8 \cdot p \quad \text{and} \quad \Sigma_{jj} = 0.2 \cdot p / (p - 1), \quad j = 2, 3, \dots, p. \quad (6)$$

In this case the marginals of the multivariate Gaussian are mutually independent. PCA is not addressed hence no dimensionality reduction is possible when the original data came from an multivariate population with uncorrelated marginals.

For the \mathcal{H}_2 , Σ is a non diagonal matrix with eigenvalues forced to be equal to the diagonal elements Σ_{jj} , $j = 2, 3, \dots, p$, set under \mathcal{H}_1 . The matrix Σ is obtained by using the product $\mathbf{Q}\Gamma\mathbf{Q}^t$, where \mathbf{Q} is a random orthogonal matrix held fixed for all the runs of the simulation, and Γ is a diagonal matrix equal to the matrix Σ defined in (6). With reference to PCA, this is the case in which the variables are correlated and the first principal component collects the most percentage of the variability of the data set.

Finally, under \mathcal{H}_3 , a non diagonal matrix Σ , built according to \mathcal{H}_2 , is still considered but its eigenvalues are now taken as the integers from 1 to p , i.e. $\Gamma_{jj} = j$, $j = 1, 2, \dots, p$. This is the case where the eigenvalues are linearly decreasing and the choice of the principal components is ambiguous.

As in the case of the $\hat{\alpha}$'s, the eigenvalues in (4), for $k = 0$, are computed for the raw and RT-2 data. From Table 1 we see that the test has the same empirical power $\hat{\pi}$ both for the alternatives \mathcal{H}_2 and \mathcal{H}_3 , so that no difference can be pointed out by the use of the rank transformation applied to the raw data. In this respect we recall that under the RT-2 data columns the chi-square distribution has $q_{RT2}^{(0)}$ degrees of freedom. As far

it concerns the hypothesis \mathcal{H}_1 , the corresponding empirical power $\hat{\pi}(\mathcal{H}_1)$ equals the empirical type I error probability $\hat{\alpha}$ when the data are RT-2. This result is expected and it is due to the flattening effect of the rank transformation. When a RT is applied to all marginals of an observed multivariate variable, the data matrices \mathbf{X} , drawn under \mathcal{H}_0 and \mathcal{H}_1 , are the same. It follows that the corresponding estimates of Σ are structurally similar.

When PCA is applied to real data set, very often the hypothesis that the data are generated from a multivariate gaussian random variable is not satisfactory. For this reason it is important to explore the reliability of the test when the distributional assumption about the population is violated. In other words we are going to explore the robustness of the test.

We adopt the same data generation schemes for all the models of Σ under \mathcal{H}_0 , \mathcal{H}_1 , \mathcal{H}_2 and \mathcal{H}_3 as illustrated for the previous set of simulations, but we generate the first marginal of the multivariate distribution from a t -distribution with 3 degrees of freedom, divided by $\sqrt{3}$, so that its theoretical variance is 1. The results are collected in Table 2. As expected, looking at the $\hat{\alpha}$'s corresponding to the raw data, the test does not recognize the shape of the covariance matrix under \mathcal{H}_0 because the data are not from a multivariate normal distribution. On the contrary, the RT-2 data preserve the nominal level of the test. The reason is the same we gave as remark to the power under \mathcal{H}_1 in the previous numerical experiment. When we consider the powers $\hat{\pi}(\cdot)$'s, the results are the same for the raw data and the RT-2 transformed. The results in the sub-table of the raw data show that the test is unreliable when the data come from a non Gaussian multivariate population. Instead, it continues to be reliable for RT-2 data when the population has a non diagonal variance/covariance matrix.

The last simulation concerns the use of the sphericity test to correctly select the principal components. The population follows a p -variate Gaussian distribution where the variance-covariance matrix is non diagonal and its eigenvalues are set as the Σ_{jj} in (6). Under this Σ -structure we expect that the null hypothesis $\mathcal{H}_0^{(0)} : l_1 = l_2 = \dots = l_p$ is rejected, while the hypothesis $\mathcal{H}_0^{(1)} : l_2 = l_3 = \dots = l_p$ is not rejected. In the frame-

Table 1. Simulations results assuming a multivariate Normal population.

p	raw data				RT-2 data			
	$\hat{\alpha}$	$\hat{\pi}(\mathcal{H}_1)$	$\hat{\pi}(\mathcal{H}_2)$	$\hat{\pi}(\mathcal{H}_3)$	$\hat{\alpha}$	$\hat{\pi}(\mathcal{H}_1)$	$\hat{\pi}(\mathcal{H}_2)$	$\hat{\pi}(\mathcal{H}_3)$
3	0.0472	1.0000	1.0000	0.9694	0.0510	0.0502	1.0000	0.9393
5	0.0486	1.0000	1.0000	1.0000	0.0482	0.0532	1.0000	1.0000
10	0.0497	1.0000	1.0000	1.0000	0.0495	0.0459	1.0000	1.0000
20	0.0512	1.0000	1.0000	1.0000	0.0516	0.0518	1.0000	1.0000
50	0.0475	1.0000	1.0000	1.0000	0.0486	0.0477	1.0000	1.0000

Table 2. Simulations results assuming non Normal populations.

p	raw data				RT-2 data			
	$\hat{\alpha}$	$\hat{\pi}(\mathcal{H}_1)$	$\hat{\pi}(\mathcal{H}_2)$	$\hat{\pi}(\mathcal{H}_3)$	$\hat{\alpha}$	$\hat{\pi}(\mathcal{H}_1)$	$\hat{\pi}(\mathcal{H}_2)$	$\hat{\pi}(\mathcal{H}_3)$
3	0.2636	1.0000	1.0000	0.9780	0.0530	1.0000	1.0000	0.9373
5	0.3035	1.0000	1.0000	1.0000	0.0513	1.0000	1.0000	1.0000
10	0.3022	1.0000	1.0000	1.0000	0.0513	1.0000	1.0000	1.0000
20	0.2882	1.0000	1.0000	1.0000	0.0511	1.0000	1.0000	1.0000
50	0.2504	1.0000	1.0000	1.0000	0.0512	1.0000	1.0000	1.0000

work of the PCA this means that only the first linear combination has to be retained for interpretation, while the others are meaningless. The test statistic for $\mathcal{H}_0^{(0)}$ is specified in (4) for $k = 0$, and $\hat{\alpha}_0$ is the corresponding empirical level of significance; for $k = 1$, (4) corresponds to the hypothesis $\mathcal{H}_0^{(1)}$ with empirical level $\hat{\alpha}_1$. The $\hat{\alpha}_k$, $k = 0, 1$, for the RT-2 data are computed by using the $q_{RT2}^{(k)}$ as in (5). The empirical results are in Table 3. These estimates highlight that test-statistic (4) is not able to correctly select the first linear combination when data are rank transformed.

Table 3. Empirical levels of significance $\hat{\alpha}_0$ and $\hat{\alpha}_1$, corresponding to $\mathcal{H}_0^{(0)}$ and $\mathcal{H}_0^{(1)}$, when data are drawn from a multivariate Normal population.

p	raw data		RT-2 data	
	$\hat{\alpha}_0$	$\hat{\alpha}_1$	$\hat{\alpha}_0$	$\hat{\alpha}_1$
3	1.000	0.050	1.000	0.986
5	1.000	0.048	1.000	1.000
10	1.000	0.049	1.000	1.000
20	1.000	0.048	1.000	1.000
50	1.000	0.051	1.000	1.000

4. Concluding remarks

This paper concerns the reliability of the sphericity test with respect to the Principal Component Analysis when the observed data are rank transformed. It is in fact a common practice to routinely transform in ranks a data matrix when, for example, some observations appear to be very far from the bulk of the data; the transformed data are then used as input for a principal component analysis. The performance of the test has been investigated under the assumption that data are drawn from a multivariate Gaussian and when they are generated from a non Normal multivariate random variable and then transformed in ranks.

When the data matrix is RT according to the RT-2 scheme, it has been given empirical evidence that the sampling distribution of the test-

statistics is still in the family of the chi-square distributions, but the degrees of freedom decrease with respect to the case when population is a multivariate Normal. For this reason we recomputed the degrees of freedom.

The results from simulation study demonstrate that the sphericity test works as theoretically expected only when the population behind the data is a multivariate Gaussian. When rank transformation RT-2 is applied to the data, the test continues to be reliable only in case data are correlated and the degrees of freedom are computed according to (5). The power of the test remarkably decreases when the variables are independent. However the use of RT when the population is not Gaussian makes the test still reliable.

Finally, the sequential use of the sphericity test can no longer be advised when the data are RT-2 transformed. In fact, the empirical level of significance for $\mathcal{H}_0^{(0)} : l_2 = l_3 = \dots = l_p$ does not comply the theoretical one.

Acknowledgements: This research was supported by the 2006 research funds of the Università degli Studi del Sannio.

References

- Bartlett, M.S. (1937), Properties of sufficiency and statistical tests, *Proceedings of the Royal Society of London, Series A*, 160, 268-282.
- Baxter M.J. (1995), Standardization and transformation in principal component analysis, with applications to archeometry, *Applied Statistics*, 44, 513-527.
- Borkowf C.B. (2002), Computing the nonnull asymptotic variance and the asymptotic relative efficiency of spearman's rank correlation, *Computational Statistics & Data Analysis*, 39, 271-286.
- Conover W.J., Iman R.L. (1981), Rank transformation as a bridge between parametric and nonparametric statistics, *The American Statistician*, 35, 124-129.
- Headrick T.C., Rotou O. (2001), An investigation of the rank transformation in multiple regression, *Computational Statistics & Data Analysis*, 38, 203-215.

Jambu M. (1989), *Exploratory and multivariate data analysis*, Academic Press, San Diego.

Jobson J.D. (1992), *Applied multivariate analysis - volume II: Categorical and multivariate methods*, Springer Verlag, New York.

Jolliffe M. (2002), *Principal component analysis (2nd ed.)*, Springer Verlag, New York.

Mardia K.V., Kent J.T., Bibby J.M. (1979), *Multivariate analysis*, Academic Press, London.

Marques F.J., Coelho C.A. (2007), Near-exact distributions for the sphericity likelihood ratio test statistics, *Journal of Statistical Planning and Inference*, 137, 1560-1575.

Mauchly J.W. (1940), Significance test for sphericity of a normal n -variate distribution, *Annals of Mathematical Statistics*, 11, 204-209.

Nagarsenker B.N., Pillai K.C.S. (1973), Distribution of the likelihood ratio criterion for testing a hypothesis specifying a covariance matrix, *Biometrika*, 60, 359-364.

Nath R., Pavur R. (1985), A new statistic in the one-way multivariate analysis of variance, *Computational Statistics & Data Analysis*, 2, 297-315.

Pillai K.C.S., Nagarsenker B.N. (1971), On the distribution of the sphericity test criterion in classical and complex normal population having unknown covariance matrices, *Annals of Mathematical Statistics*, 42, 764-767.

Thompson G.L. (1991), A unified approach to rank tests for multivariate and repeated measure designs, *Journal of the American Statistical Association*, 42, 410-419.

Waternaux C.M. (1984), Principal components in the nonnormal case: the test of equality of Q roots, *Journal of Multivariate Analysis*, 14, 232-335.