# A better approximation of moments of the eigenvalues and eigenvectors of the sample covariance matrix

A. Enguix-González *, J.L. Moreno-Rebollo, J.M. Muñoz-Pichardo

*Department of Statistics and Operations Research, Faculty of Mathematics, University of Seville, C/ Tarfia s/n, 41012 Seville, Spain*

### ABSTRACT

Lawley (Lawley, 1956) obtained an approximation, through the first terms of a series expansion, of certain moments of an eigenvalue of the sample covariance matrix. The aim of this paper is to improve that approximation and to calculate a similar approximation for certain moments of the associated eigenvector. The results have practical applications in certain fields of Statistics, such as Influence Analysis.

© 2015 Elsevier Inc. All rights reserved.

## 1. Introduction

Eigenvalues and eigenvectors of the sample covariance matrix are the basic statistics in some statistical techniques, such as Principal Component Analysis, and their moments play a major role in the inferential process.

However, the exact distributions of the eigenvalues and eigenvectors of the sample covariance matrix remain unknown. In the literature various approximations have been obtained. Lawley [4] obtained the first terms of a series expansion of a sample eigenvalue related to a simple population eigenvalue. Moreover, under the assumption of normality, approximated values for the mean, variance and covariance matrix were obtained from the expansion series. Approximations for the moments of the sample eigenvectors are given in Seber [6].

In this paper, we add a term to the series expansion obtained by Lawley [4] and improve the approximations to their moments. Furthermore, series expansions for the eigenvectors and their moments are obtained. The use of symbolic computation with MAPLE has proved to be very helpful in dealing with the cumbersome calculations.

In Section 2, the notation to be used throughout the paper and some basic results are set out. In Section 3, the procedure applied to obtain the terms of the series expansions is laid out and the expressions of the five first terms are provided. In Section 4, approximations of a number of moments of the eigenvalues, which improve those given by Lawley [4], are obtained. Furthermore, some moments of the eigenvectors are calculated with the same degree of accuracy. In Section 5, a numeric evaluation through a Monte Carlo simulation has been carried out. In the last section, the most relevant conclusions are stated.

---

\* Corresponding author.
  *E-mail address:* aenguix@us.es (A. Enguix-González).

## 2. Notation and basic results

Let $\underline{X}$ be a $p$-dimensional random vector with mean vector $\underline{\mu} = (\mu_1, \mu_2, \ldots, \mu_p)'$ and covariance matrix $\mathbf{\Sigma} = (\sigma_{st})$. Let $\lambda_1, \lambda_2, \ldots, \lambda_p$ be the eigenvalues of $\mathbf{\Sigma}$, $\underline{\alpha}_1, \underline{\alpha}_2, \ldots, \underline{\alpha}_p$ orthonormal associated eigenvectors, respectively, and $\mathbf{A} = [\underline{\alpha}_1 \, \underline{\alpha}_2 \, \cdots \, \underline{\alpha}_p]'$. Let $\underline{Y} = (Y_1, \ldots, Y_p)' = \mathbf{A}(\underline{X} - \underline{\mu})$. The (population) principal components of $\underline{X}$ are $Y_k = \underline{\alpha}'_k (\underline{X} - \underline{\mu})$, $k = 1, \ldots, p$. Obviously, the covariance matrix of $\underline{Y}$ is given by $\mathbf{\Sigma}_Y = \mathbf{A}\mathbf{\Sigma}\mathbf{A}' = \mathbf{\Lambda} = diag\{\lambda_1, \ldots, \lambda_p\}$, where $diag$ denotes the diagonal matrix.

In practice, $\underline{\mu}$ and $\mathbf{\Sigma}$ are unknown parameters. Therefore, they have to be estimated from a random sample, $\underline{X}_1, \ldots, \underline{X}_n$, with $\underline{X}_i = (X_{i1}, \ldots, X_{ip})'$, $i = 1, \ldots, n$. These are usually estimated by the sample mean vector, $\overline{X} = (\overline{X}_1, \overline{X}_2, \ldots, \overline{X}_p)'$, where $\overline{X}_j = \frac{1}{n} \sum_{i=1}^{n} X_{ij}$, $j = 1, \ldots, p$, and the sample covariance matrix $\mathbf{S} = (S_{lm}) = \frac{1}{n-1} \sum_{i=1}^{n} (\underline{X}_i - \overline{X})(\underline{X}_i - \overline{X})'$, respectively. Furthermore, the eigenvalues, $\widetilde{\lambda}_1, \widetilde{\lambda}_2, \ldots, \widetilde{\lambda}_p$, and orthonormal associated eigenvectors, $\widetilde{\underline{\alpha}}_1, \widetilde{\underline{\alpha}}_2, \ldots, \widetilde{\underline{\alpha}}_p$, of $\mathbf{S}$ can be considered as estimators of the population eigenvalues and eigenvectors, respectively. The sample principal components of $\underline{X}_i$ are given by $\widetilde{Y}_{ik} = \widetilde{\underline{\alpha}}'_k (\underline{X}_i - \overline{X})$, $k = 1, \ldots, p, i = 1, \ldots, n$. Moreover, the sample covariance matrix of $\underline{Y}$ is given by $\mathbf{S}_Y = \mathbf{A}\mathbf{S}\mathbf{A}' = \mathbf{T} = (T_{lm})$.

Since $\mathbf{T}$ is an orthonormal transformation of $\mathbf{S}$, the eigenvalues of $\mathbf{S}$ coincide with the eigenvalues of $\mathbf{T}$ and the eigenvectors of $\mathbf{T}$ are given by $\widetilde{\beta}_k = \mathbf{A}\widetilde{\alpha}_k$.

By using the orthonormal transformation of $\mathbf{S}$, $\mathbf{T} = \mathbf{A}\mathbf{S}\mathbf{A}'$, and taking into account that $\mathbf{T}$ converges almost surely to $\mathbf{\Lambda} = diag\{\lambda_1, \ldots, \lambda_p\}$, Lawley [4] calculated the first terms of the series expansion of an eigenvalue of $\mathbf{S}$. Assuming that $\lambda_1, \ldots, \lambda_p > 0$, if the sample size is sufficiently large and $\lambda_k$ is simple, the convergent series expansion of the sample eigenvalue $\widetilde{\lambda}_k$ is given by

$$\widetilde{\lambda}_k = \lambda_k + (T_{kk} - \lambda_k) - \sum_{l \neq k} \frac{T_{lk}^2}{\lambda_l - \lambda_k} - (T_{kk} - \lambda_k) \sum_{l \neq k} \frac{T_{lk}^2}{(\lambda_l - \lambda_k)^2} + \sum_{l \neq k} \sum_{m \neq k} \frac{T_{lk}(T_{lm} - \lambda_l \delta_{lm}) T_{mk}}{(\lambda_l - \lambda_k)(\lambda_m - \lambda_k)} + R_k, \tag{1}$$

where $\delta$ represents Kronecker's delta and $R_k$ is a sum of terms with a product of four or more factors of the type $T_{lm} - \lambda_l \delta_{lm}$.

Since the rest of a convergent series tends to zero, $\widetilde{\lambda}_k$ can be approximated by $\widetilde{\lambda}_k - R_k$. From $\widetilde{\lambda}_k - R_k$ approximations to the moments of $\widetilde{\lambda}_k$ can be obtained.

Assuming that $\underline{X} \sim \mathcal{N}_p (\underline{\mu}, \mathbf{\Sigma})$ and $\lambda_k$ is a simple eigenvalue, Lawley [4] obtains:

$$E(\widetilde{\lambda}_k) = \lambda_k - \frac{1}{n-1} \sum_{j \neq k} \frac{\lambda_j \lambda_k}{\lambda_j - \lambda_k} + O(n^{-2}),$$

$$var(\widetilde{\lambda}_k) = \frac{2\lambda_k^2}{n-1} \left[1 - \frac{1}{n-1} \sum_{j \neq k} \frac{\lambda_j^2}{(\lambda_j - \lambda_k)^2}\right] + O(n^{-3}),$$

$$cov[\widetilde{\lambda}_j, \widetilde{\lambda}_k] = \frac{2}{(n-1)^2} \frac{\lambda_j^2 \lambda_k^2}{(\lambda_j - \lambda_k)^2} + O(n^{-3}), \quad j \neq k.$$

Regarding the eigenvectors associated to simple eigenvalues, it is known that, see Seber [6],

$$E(\widetilde{\underline{\alpha}}_k) = \underline{\alpha}_k + O(n^{-1}),$$

$$var(\widetilde{\underline{\alpha}}_k) = \frac{\lambda_k}{n-1} \sum_{j \neq k} \frac{\lambda_j}{(\lambda_j - \lambda_k)^2} \underline{\alpha}_j \underline{\alpha}'_j + O(n^{-2}),$$

$$cov(\widetilde{\underline{\alpha}}_j, \widetilde{\underline{\alpha}}_k) = -\frac{1}{n-1} \frac{\lambda_j \lambda_k}{(\lambda_j - \lambda_k)^2} \underline{\alpha}_j \underline{\alpha}'_k + O(n^{-2}), \quad \text{for } j \neq k.$$

## 3. Series expansion of eigenvalues and eigenvectors

In this section, the necessary tools to obtain the series expansions of the sample eigenvalues and eigenvectors are provided, and the first five terms are explicitly obtained.

Let $\mathcal{D}_p$ be the set of symmetric positive definite matrices of dimension $p$. Let $\mathbf{\Sigma}_0 \in \mathcal{D}_p$ and let $\mathcal{N}(\mathbf{\Sigma}_0) \subseteq \mathcal{D}_p$ be a neighborhood of $\mathbf{\Sigma}_0$. Let $\lambda_k^{\mathbf{\Sigma}_0}, \underline{\alpha}_k^{\mathbf{\Sigma}_0}$, $k = 1, \ldots, p$, be functions satisfying:

- $\lambda_k^{\mathbf{\Sigma}_0} : \mathcal{N}(\mathbf{\Sigma}_0) \longrightarrow \mathbb{R}^+$; $\underline{\alpha}_k^{\mathbf{\Sigma}_0} : \mathcal{N}(\mathbf{\Sigma}_0) \longrightarrow \mathbb{R}^p$, $k = 1, \ldots, p$.
- If $\mathbf{\Sigma} \in \mathcal{N}(\mathbf{\Sigma}_0)$, then $\lambda_1^{\mathbf{\Sigma}_0}(\mathbf{\Sigma}) > \lambda_2^{\mathbf{\Sigma}_0}(\mathbf{\Sigma}) > \cdots > \lambda_p^{\mathbf{\Sigma}_0}(\mathbf{\Sigma}) > 0$.