



Taylor & Francis
Taylor & Francis Group

Criteria for Evaluating Dimension-Reducing Components for Multivariate Data

Author(s): Daniel Gervini and Valentin Rousson

Source: *The American Statistician*, Vol. 58, No. 1 (Feb., 2004), pp. 72-76

Published by: Taylor & Francis, Ltd. on behalf of the American Statistical Association

Stable URL: <http://www.jstor.org/stable/27643501>

Accessed: 17-10-2017 05:21 UTC

REFERENCES

Linked references are available on JSTOR for this article:

http://www.jstor.org/stable/27643501?seq=1&cid=pdf-reference#references_tab_contents

You may need to log in to JSTOR to access the linked references.

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <http://about.jstor.org/terms>



JSTOR

American Statistical Association, Taylor & Francis, Ltd. are collaborating with JSTOR to digitize, preserve and extend access to *The American Statistician*

Criteria for Evaluating Dimension-Reducing Components for Multivariate Data

Daniel GERVINI and Valentin ROUSSON

Principal components are the benchmark for linear dimension reduction, but they are not always easy to interpret. For this reason, some alternatives have been proposed in recent years. These methods produce components that, unlike principal components, are correlated and/or have nonorthogonal loadings. This article shows that the criteria commonly used to evaluate principal components are not adequate for evaluating such alternatives, and proposes two new criteria that are more suitable for this purpose.

KEY WORDS: Linear prediction; Principal components; Rotated components; Simple components.

1. INTRODUCTION

In multivariate datasets, it is often the case that the variables are highly correlated and provide redundant information. The number of variables is then unnecessarily large, and essentially the same information can be conveyed by fewer dimensions if the variables are wisely combined. In many cases, a lower dimension also helps to visualize patterns in the data that would otherwise go unnoticed.

For these reasons, dimension-reduction techniques have played an important role in multivariate analysis. Many of these techniques construct a system of q variables which are linear combinations of the original p variables; the new variables are called components. The most popular of these methods is principal component analysis, which was originally proposed by Hotelling (1933); a comprehensive and up-to-date reference is Jolliffe (2002). The idea is to sequentially construct a system of components that are uncorrelated and have maximal variance. The component coefficients (called loadings) obtained in this way turn out to be orthogonal. For a given target dimension q , the first q principal components are then the optimal dimension-reducing system because they extract the maximal variability, and they are both statistically nonredundant (uncorrelated) and geometrically nonredundant (orthogonal loadings). For this reason, they are considered the benchmark for linear dimension reduction. Principal components are optimal under different cri-

teria that evaluate uncorrelated or orthogonal components (Rao 1964; Okamoto 1969; McCabe 1984).

In most cases, however, the researcher not only wants to reduce the dimension of the dataset but also wants to obtain components that are interpretable in the context of his research. Unfortunately, principal component loadings sometimes show complicated patterns and are not easy to interpret; see Cadima and Jolliffe (1995) for interesting examples. To improve interpretability, alternative methods that produce components with simpler loadings' patterns have been proposed over the years—for example, Neuhaus and Wrigley (1954), Kaiser (1958), Hausmann (1982), Kiers (1991), Jolliffe and Uddin (2000), Vines (2000), and Rousson and Gasser (2004); see also Jolliffe (2002, chap. 11) for other proposals. These methods produce components that are no longer uncorrelated and have nonorthogonal loadings, so they are less efficient than principal components at dimension reduction. But how can we quantify precisely the loss of dimension-reducing efficiency of such components? How can we compare the performance of different methods? Remember that most of the optimality criteria mentioned in the preceding paragraph assume that the components are either uncorrelated or have orthogonal loadings, which is no longer the case with the alternative methods.

To illustrate the problem with a real dataset, let us consider the audiometric example analyzed by Jackson (1991, chap. 5) and reanalyzed by Vines (2000). The data consist of measurements of lower hearing threshold on 100 men. Observations were obtained, on each ear, at frequencies 500, 1,000, 2,000, and 4,000 Hz, so that eight variables were recorded for each individual. The sample variance of measurements at 4,000 Hz turned out to be about nine times higher than those at 500 Hz, so the variables were standardized before computing the principal components—that is, the principal components were computed on the correlation matrix rather than the covariance matrix. The covariance–correlation matrix is given in Table 1. Jackson (1991) argued that the first four principal components, which explain 87% of the total variability, provide a good approximation to the data. The principal component loadings are given in Table 2. They can be interpreted as follows: the first component is an indicator of average hearing loss, the second one is a contrast between high- and low-frequency hearing loss, the third one is a contrast between hearing loss at the two highest frequencies, and the fourth one is a contrast between the two ears. Note that some of the components (the third one in particular) are somewhat difficult to interpret at first glance, because some loadings are relatively small but not really close to zero, so it is not clear whether they are significant or not. More clear-cut loadings are obtained with the methods of Vines (2000) and Rousson and Gasser (2004), which produce the same components for this dataset; the loadings are given in Table 2. The simplicity of the loadings' patterns allows unequivocal interpretation of these components, but before deciding to use this system instead of the

The authors are Postdoctorates, Department of Biostatistics, University of Zürich, Sumatrastrasse 30, CH-8006 Switzerland (E-mail: gervini@ifspm.unizh.ch; rousson@ifspm.unizh.ch). Their research was supported by the Swiss National Science Foundation (grants BE 20.63579.00 and 3200-064047.00/1). The authors thank Theo Gasser for many fruitful discussions, the referees for their constructive criticism, and the editor and an associate editor for suggestions to improve the presentation.

Table 1. Covariance–Correlation Matrix of Hearing Loss Data

	Left ear				Right ear			
	500	1,000	2,000	4,000	500	1,000	2,000	4,000
L, 500	41.07	(0.78)	(0.40)	(0.26)	(0.70)	(0.64)	(0.24)	(0.20)
L, 1,000	37.73	57.32	(0.54)	(0.27)	(0.55)	(0.71)	(0.36)	(0.22)
L, 2,000	28.13	44.44	119.70	(0.42)	(0.24)	(0.45)	(0.70)	(0.33)
L, 4,000	32.10	40.83	91.21	384.78	(0.18)	(0.26)	(0.32)	(0.71)
R, 500	31.79	29.75	18.64	25.01	50.75	(0.66)	(0.16)	(0.13)
R, 1,000	26.30	34.24	31.21	33.03	30.23	40.92	(0.41)	(0.22)
R, 2,000	14.12	25.30	71.26	57.67	10.52	24.62	86.30	(0.37)
R, 4,000	25.28	31.74	68.99	269.12	18.19	27.22	67.26	373.66

NOTE: Correlations are given in parenthesis.

principal components, the statistician should know how much is lost in terms of dimension-reducing power. The simpler system is no longer uncorrelated, so it does not make sense to simply add up the variances and compare it with the total variance of the original dataset, as it is done with the principal components. Vines (2000) compared the variance of each simple component with the variance of the corresponding principal component, concluding that “little explanatory power (in terms of variance) is lost by this radical simplification. Furthermore the highest correlation between the first four simple components is only 0.151” (Vines 2000, p. 448). This is rather vague, however; it would be nice to have a criterion that indicates unambiguously, with a single number (ranging from 0 to 1, say, with 1 being optimal) the dimension-reducing power of the system.

Some of the authors mentioned earlier (Neuhaus and Wrigley 1954; Kaiser 1958; Kiers 1991; Jolliffe and Uddin 2000) have already proposed nonstandard criteria to evaluate components. The problem is that these authors aimed at simplicity, so they proposed criteria that measure some sort of simplicity of the system rather than its dimension-reducing power. As a result, it does not make much sense to use, for example, the quartimax criterion of Neuhaus and Wrigley (1954) to evaluate varimax components, which by definition maximize the different simplicity criterion of Kaiser (1958). Consider, for example, the varimax components for the hearing loss example, given in Table 2. They are harder to interpret than either the principal or the simple components, and are highly correlated. Varimax components offer only disadvantages in this example, yet they are considered optimal (by definition) under the criterion of Kaiser (1958), while principal and simple components are considered suboptimal. It is clear, then, that we need criteria that evaluate the

dimension-reducing power of components independently of the notion of simplicity. In this article we are going to review some of the existing criteria and propose two new ones, because we found that none of the existing criteria is completely adequate for this task.

2. NECESSARY AND DESIRABLE PROPERTIES OF CRITERIA FOR EVALUATING COMPONENTS

From the discussion in Section 1, we conclude that criteria for evaluating components should assign optimal value to the principal components, since they are the most efficient dimension-reducing system in terms of variability extraction and nonredundancy of information. These criteria should also be applicable to systems of components that may not be uncorrelated and may not have orthogonal loadings, because most alternatives to principal components do not. For example, both simple and varimax components for the hearing loss example are correlated, so we cannot simply add up the variances of the components and divide it by the sum of variances of the principal components; a more elaborate criterion, that takes correlations into account, is necessary. Specifically, correlations between components should be penalized, because they imply redundancy of information.

A criterion for evaluating dimension-reducing components, then, should satisfy at least two conditions:

1. *Generality.* The criterion has to be applicable to a broad range of components, with the only restriction of unit-norm and linearly independent loadings—the least restrictive assumptions that rule out artificial cases. Under these general conditions, and for a given target dimension q , the criterion must be maximized by the first q principal components.

Table 2. Component Loadings for Hearing Loss Data

Variable	Principal components				Simple components				Varimax components			
	1st	2nd	3rd	4th	1st	2nd	3rd	4th	1st	2nd	3rd	4th
L, 500	0.40	−0.32	0.16	−0.33	0.35	−0.35	0.00	−0.35	0.60	0.03	−0.09	0.15
L, 1,000	0.42	−0.23	−0.05	−0.48	0.35	−0.35	0.00	−0.35	0.67	−0.03	0.11	−0.03
L, 2,000	0.37	0.24	−0.47	−0.28	0.35	0.35	−0.50	−0.35	0.29	0.02	0.61	−0.19
L, 4,000	0.28	0.47	0.43	−0.16	0.35	0.35	0.50	−0.35	0.13	0.70	−0.01	−0.10
R, 500	0.34	−0.39	0.26	0.49	0.35	−0.35	0.00	0.35	0.03	0.02	−0.16	0.74
R, 1,000	0.41	−0.23	−0.03	0.37	0.35	−0.35	0.00	0.35	0.07	−0.02	0.15	0.58
R, 2,000	0.31	0.32	−0.56	0.39	0.35	0.35	−0.50	0.35	−0.26	−0.01	0.75	0.21
R, 4,000	0.25	0.51	0.43	0.16	0.35	0.35	0.50	0.35	−0.13	0.71	0.02	0.09

2. *Uniqueness.* The criterion must be maximized *only* by the principal components under the conditions mentioned above.

The Uniqueness condition might seem too strong, but it guarantees that correlations between components and deviations from orthogonality are penalized.

Other properties may be useful, even desirable, but we do not think that they are strictly necessary. For instance:

- *Additivity.* Many criteria can be naturally expressed as a sum of q terms, indicating the contribution of each component towards the overall dimension reduction. This is a good thing, but we are mainly interested in evaluating systems of components as a whole, rather than individual contributions of the components.

- *Invariance under permutation of components.* Because we are evaluating systems as a whole, a criterion that assigns different values to two systems of components which are just a permutation of one another is not very appealing; therefore, permutation invariance is desirable. In practice, however, the components are computed in a sequential way, so that a natural ordering is given by construction and alternatives consisting merely on permutations are normally not contemplated.

The next section reviews some existing criteria, focusing on those that satisfy the property of Generality. It turns out that none of them satisfies the property of Uniqueness. This motivates our introduction in Section 4 of two new criteria that satisfy both properties.

3. EXISTING CRITERIA

Before we start reviewing the existing criteria, let us introduce some notation. Consider a random vector $\mathbf{x} \in \mathbb{R}^p$, that without loss of generality will be assumed to have zero mean. A linear dimension reduction technique will produce a system of components $\mathbf{y} = \mathbf{A}^\top \mathbf{x}$, where $\mathbf{A} \in \mathbb{R}^{p \times q}$ is called the loading matrix and $q \leq p$. The principal components are defined as follows. Let $\Sigma = \text{cov}(\mathbf{x})$ and $\Sigma = \Gamma \Lambda \Gamma^\top$ be the eigenvalue decomposition of Σ , where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ with $\lambda_1 \geq \dots \geq \lambda_p > 0$ and $\Gamma \in \mathcal{O}(p)$, where $\mathcal{O}(p)$ denotes the family of $p \times p$ orthogonal matrices ($\mathcal{O}(p, q)$ will denote the family of $p \times q$ orthogonal matrices). The elements of $\mathbf{z} = \Gamma^\top \mathbf{x}$ are called the principal components of \mathbf{x} . Note that $\text{cov}(\mathbf{z}) = \Lambda$, so that the principal components are uncorrelated. Γ_q will indicate the loading matrix consisting of the first q columns of Γ , and $\Lambda_q = \text{diag}(\lambda_1, \dots, \lambda_q)$. Note that the loading matrix Γ , and consequently Γ_q , is only determined up to column sign reversal and exchange of columns with identical eigenvalues, so it is not unique. But to avoid unnecessary complications in phraseology, we will refer to Γ_q as “the unique” loading matrix.

There are essentially three approaches to dimension reduction: prediction (find $\mathbf{y} = \mathbf{A}^\top \mathbf{x}$ that provides the best linear prediction of \mathbf{x}), variability maximization (find \mathbf{y} with the largest possible variance among linearly independent combinations of \mathbf{x}), and correlation (find \mathbf{y} that is maximally correlated with \mathbf{x}). The review that follows is organized in three subsections corresponding to these approaches.

3.1 Prediction Approach

The best linear predictor of \mathbf{x} based on $\mathbf{y} = \mathbf{A}^\top \mathbf{x}$, in the sense of minimizing $E(\|\mathbf{x} - \mathbf{B}\mathbf{y}\|^2)$, is $\hat{\mathbf{B}}\mathbf{y}$ with $\hat{\mathbf{B}} = \Sigma \mathbf{A}(\mathbf{A}^\top \Sigma \mathbf{A})^{-1}$. Therefore, the matrix \mathbf{A} producing the optimal predictor is $\hat{\mathbf{A}}$ that minimizes

$$E(\|\mathbf{x} - \Sigma \mathbf{A}(\mathbf{A}^\top \Sigma \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{x}\|^2) = \text{tr}(\Sigma) - \text{tr}(\Sigma \mathbf{A}(\mathbf{A}^\top \Sigma \mathbf{A})^{-1} \mathbf{A}^\top \Sigma),$$

or equivalently, $\hat{\mathbf{A}}$ that maximizes $\text{tr}(\Sigma \mathbf{A}(\mathbf{A}^\top \Sigma \mathbf{A})^{-1} \mathbf{A}^\top \Sigma)$. Rao (1964) showed that $\hat{\mathbf{A}} = \Gamma_q$ is the maximizer in $\mathcal{O}(p, q)$, and the maximum is $\text{tr}(\Lambda_q) = \sum_{k=1}^q \lambda_k$. Then our first criterion is

$$\text{BLP}(\mathbf{A}) = \frac{\text{tr}(\Sigma \mathbf{A}(\mathbf{A}^\top \Sigma \mathbf{A})^{-1} \mathbf{A}^\top \Sigma)}{\sum_{k=1}^q \lambda_k}, \quad (1)$$

where BLP stands for “best linear prediction” (all criteria in this article are standardized so that the optimum is 1). This can be rewritten as an additive criterion if one so wishes.

The value of BLP depends only on the subspace spanned by the columns of \mathbf{A} , rather than on the actual matrix \mathbf{A} . Consequently, any full-rank transformation of \mathbf{A} is equivalent for this criterion. Gervini and Rousson (2003) proved that any matrix of the form $\hat{\mathbf{A}} = \Gamma_q \mathbf{C}$ with nonsingular $\mathbf{C} \in \mathbb{R}^{q \times q}$ maximizes (1). Therefore BLP satisfies the property of Generality, but not Uniqueness; it fails to discriminate between systems of components that are obviously not equivalent from a practical point of view. Therefore, this criterion is not adequate for our purposes.

3.2 Variability Maximization Approach

Finding a system of uncorrelated components with largest variance is probably the most familiar approach to dimension reduction. The q components with largest variance are the ones that carry most of the information of the original data, while the others vary little about zero. In fact, if \mathbf{x} lies on a q -dimensional subspace of \mathbb{R}^p with probability 1, then the variance of the last $p - q$ components is exactly zero. The total variance of the system can be defined as either the trace or the determinant of the covariance matrix (the latter is usually known as generalized variance). Using $\text{tr}(\text{cov}(\mathbf{y})) = \text{tr}(\mathbf{A}^\top \Sigma \mathbf{A})$ is more common, and if the components are assumed to be uncorrelated, the maximization can be carried out in a sequential way, by maximizing $\text{var}(y_j) = \mathbf{a}_j^\top \Sigma \mathbf{a}_j$ subject to the restrictions $\|\mathbf{a}_j\| = 1$ and $\text{cov}(y_j, y_k) = \mathbf{a}_j^\top \Sigma \mathbf{a}_k = 0$ for all $k < j$. The optimal \mathbf{y} turns out to be the vector of the first q principal components. It is interesting to note that the optimal loading matrix comes out orthogonal, although this was not an explicit restriction. If one imposes the restriction $\mathbf{A} \in \mathcal{O}(p, q)$ instead of uncorrelation, then the maximizers of $\text{tr}(\mathbf{A}^\top \Sigma \mathbf{A})$ turn out to be the matrices of the form $\Gamma_q \mathbf{R}$ with $\mathbf{R} \in \mathcal{O}(q)$, that is, the rotations of the first q principal components. It is clear, however, that maximizing the trace only makes sense when either one of the restrictions of uncorrelation or orthogonality is imposed, which violates the property of Generality.

On the other hand, the generalized variance $\det(\text{cov}(\mathbf{y}))$ is maximized by Γ_q under the unique restriction of unit-norm loadings, without assuming orthogonality or uncorrelation (see

Okamoto 1969). Then the criterion

$$GV(\mathbf{A}) = \left(\frac{\det(\mathbf{A}^\top \Sigma \mathbf{A})}{\prod_{k=1}^q \lambda_k} \right)^{\frac{1}{q}}$$

satisfies the property of Generality, in contrast with the trace criterion. But unfortunately GV is invariant under rotation of the components and then it does not satisfy the property of Uniqueness. This implies, for example, that GV cannot discriminate between principal components and the varimax rotation. Although GV is more informative than BLP , it is still not good enough for our purposes.

3.3 Correlation Approach

The third approach to dimension reduction consists in finding components that are maximally correlated with the data, using measures of matrix correlation based on the sample data matrix. Given a sample $\mathbf{x}_1, \dots, \mathbf{x}_n$, let \mathbf{X} be the $n \times p$ data matrix and $\mathbf{Y} = \mathbf{X}\mathbf{A}$ be the $n \times q$ matrix of components. Robert and Escoufier (1976) measured the closeness between “data configurations” $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{Y}\mathbf{Y}^\top$ as $\text{corr}(\mathbf{X}\mathbf{X}^\top, \mathbf{Y}\mathbf{Y}^\top)$, where $\text{corr}(\mathbf{A}, \mathbf{B})$ is the inner-product matrix correlation $\langle \mathbf{A}, \mathbf{B} \rangle / (\langle \mathbf{A}, \mathbf{A} \rangle^{\frac{1}{2}} \langle \mathbf{B}, \mathbf{B} \rangle^{\frac{1}{2}})$ with $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}^\top \mathbf{B})$. If $\mathbf{S} = \mathbf{X}^\top \mathbf{X} / n$ denotes the sample covariance matrix, it is not difficult to see that

$$\text{corr}(\mathbf{X}\mathbf{X}^\top, \mathbf{Y}\mathbf{Y}^\top) = \frac{\text{tr}(\mathbf{S}\mathbf{A}\mathbf{A}^\top \mathbf{S})}{\{\text{tr}(\mathbf{S}^2) \text{tr}((\mathbf{A}^\top \mathbf{S}\mathbf{A})^2)\}^{\frac{1}{2}}}. \quad (2)$$

This is known as the RV -coefficient of Robert and Escoufier (1976), who showed that (2) is uniquely maximized by the first q sample principal components among all uncorrelated components. It is possible to relax the restriction of uncorrelation at the price of losing uniqueness (this trade-off is unavoidable, since $\text{corr}(\mathbf{X}\mathbf{X}^\top, \mathbf{Y}\mathbf{Y}^\top)$ is invariant under rotation of components). Assuming only that the loadings have norm one, all the maximizers of $\text{corr}(\mathbf{X}\mathbf{X}^\top, \mathbf{Y}\mathbf{Y}^\top)$ are the rotations of the first q principal components (the proof was given by Gervini and Rousson 2003). Therefore, replacing the sample covariance matrix by the population covariance matrix, from (2) we deduce a criterion that

satisfies the property of Generality, but not Uniqueness:

$$RV(\mathbf{A}) = \frac{\text{tr}(\mathbf{A}^\top \Sigma^2 \mathbf{A})}{\{\sum_{j=1}^q \lambda_j^2\}^{\frac{1}{2}} \{\text{tr}((\mathbf{A}^\top \Sigma \mathbf{A})^2)\}^{\frac{1}{2}}}.$$

Other measures of matrix correlation were used by Cadima and Jolliffe (2001) in the context of variable selection based on principal components, but none of them can be turned into a criterion better than RV , so we do not elaborate on this. The interested reader is referred to Gervini and Rousson (2003) for more details.

4. NEW CRITERIA

We saw in the previous section that the existing criteria do not satisfy the condition of Uniqueness. None of them can discriminate between rotations of principal components, and BLP does not even discriminate between arbitrary full-rank transformations. This section proposes two new criteria that satisfy the Uniqueness property.

Our first proposal is, essentially, a sum of variances corrected for correlations. The idea is that if a new component $y_k = \mathbf{a}_k^\top \mathbf{x}$ is added to a system of $k-1$ components, an indicator of the *real* contribution of y_k to the total variance of the system is the residual variance of the linear prediction of y_k given the first $k-1$ components. Adding all these residual variances together gives

$$\sum_{k=1}^q (\mathbf{a}_k^\top \Sigma \mathbf{a}_k - \mathbf{a}_k^\top \Sigma \mathbf{A}_{(k-1)} (\mathbf{A}_{(k-1)}^\top \Sigma \mathbf{A}_{(k-1)})^{-1} \mathbf{A}_{(k-1)}^\top \Sigma \mathbf{a}_k), \quad (3)$$

where $\mathbf{A}_{(k)} = (\mathbf{a}_1, \dots, \mathbf{a}_k)$. Note that (3) is just the sum of variances of the components if they are uncorrelated (because $\mathbf{A}_{(k-1)}^\top \Sigma \mathbf{a}_k = 0$), otherwise it is strictly smaller. Therefore (3) penalizes correlations, as we wanted. Moreover, the unique maximizer of (3) among full-rank matrices with unit-norm columns is Γ_q (see Gervini and Rousson 2003 for a proof, which is not trivial). Then the criterion

$$CSV(\mathbf{A}) = \frac{\sum_{k=1}^q (\mathbf{a}_k^\top \Sigma \mathbf{a}_k - \mathbf{a}_k^\top \Sigma \mathbf{A}_{(k-1)} (\mathbf{A}_{(k-1)}^\top \Sigma \mathbf{A}_{(k-1)})^{-1} \mathbf{A}_{(k-1)}^\top \Sigma \mathbf{a}_k)}{\sum_{k=1}^q \lambda_k},$$

where CSV stands for “corrected sum of variances,” satisfies both the Generality and the Uniqueness properties (and is also additive).

However, CSV is not invariant under permutation of components. In practice this is not very problematic, but it is not hard to construct an invariant criterion if one wishes to. One possibility is to simply take the maximum of CSV among all permutations of the components. Another possibility is to define a “symmetrically corrected sum of variances”

$$SCSV(\mathbf{A}) = \frac{\sum_{k=1}^q (\mathbf{a}_k^\top \Sigma \mathbf{a}_k - \mathbf{a}_k^\top \Sigma \mathbf{A}_{-k} (\mathbf{A}_{-k}^\top \Sigma \mathbf{A}_{-k})^{-1} \mathbf{A}_{-k}^\top \Sigma \mathbf{a}_k)}{\sum_{k=1}^q \lambda_k}, \quad (4)$$

where \mathbf{A}_{-k} is the $p \times (q-1)$ matrix obtained after deleting the k th column of \mathbf{A} . The numerator of (4) is the sum of the residual variances of the linear predictors of y_k given the other $q-1$ components. Note that $SCSV(\mathbf{A}) = CSV(\mathbf{A})$ if the system is uncorrelated and $SCSV(\mathbf{A}) < CSV(\mathbf{A})$ otherwise. Then, $SCSV$ is also uniquely maximized by Γ_q among full-rank matrices with unit-norm columns. This criterion also satisfies the properties of Generality and Uniqueness, plus additivity and invariance under permutation of components. But it penalizes correlations more strongly than CSV and then it can be overly pessimistic in some situations. Besides, it is not a sequential criterion: it “looks into the future,” subtracting from $\text{var}(y_k)$ correlations with components y_j with $j > k$. It must be noted, however, that the properties of invariance under permutation of components

Table 3. Covariance–Correlation Matrices for Components of Hearing Loss Data

Principal components				Simple components				Varimax components			
3.93	0	0	0	3.86	(−0.12)	(−0.09)	(−0.14)	1.85	(0.27)	(0.42)	(0.73)
0	1.62	0	0	−0.30	1.59	(0.15)	(0.05)	0.49	1.71	(0.41)	(0.24)
0	0	0.98	0	−0.18	0.19	0.99	(0.03)	0.75	0.71	1.74	(0.35)
0	0	0	0.46	−0.18	0.04	0.02	0.45	1.30	0.40	0.60	1.68

NOTE: Correlations are given in parenthesis.

and sequentiality are at odds with each other (except, of course, in the case of uncorrelated components, where one just takes the sum of variances). At this point we cannot envisage a criterion that is simultaneously permutation invariant, sequential, and penalizes correlations so as to satisfy the Uniqueness property. But in practice sequentiality seems to be preferable over permutation invariance, so that we tend to favor the CSV criterion.

5. EXAMPLE

Let us apply the criteria reviewed in Section 3 and the new criteria proposed in Section 4 to the hearing loss data presented in the Introduction. Remember that we have three alternative systems of components, shown in Table 2: the optimal principal components, the less optimal but better interpretable simple components, and the highly correlated and not very meaningful varimax components.

For this dataset, the varimax rotation evenly redistributes the total variance among components and reintroduces high correlations, as shown in Table 3. The varimax components are clearly unattractive in this example. Yet BLP, GV, and RV criteria assign maximum optimality to those components (remember that this is always so, because these criteria are invariant under rotations). On the other hand, the proposed CSV and SCSV criteria assign values 0.79 and 0.61 to these components, which is far from optimal and a more realistic evaluation of the system's performance.

That the simple components are better at dimension reduction than the varimax components is evident from Table 3. The variances are closer to the principal component variances and the correlations are relatively small. BLP, GV, and RV values are high for this system (0.99, 0.97, and 0.99, respectively), but this is hardly surprising, since these criteria tend to err on the optimistic side, assigning high values to any reasonable system of components. What is more interesting, the more demanding CSV and SCSV criteria also assign high values to this system: 0.98 and 0.95, respectively. We conclude that the correlated system of simple components incurs only a 2% (respectively 5%) loss of dimension-reducing power compared to principal components. This indicates that simple components are a good alternative for this dataset.

6. CONCLUSION

A number of alternatives to principal components have been proposed recently, that sacrifice some of the dimension-reducing power of the principal components in exchange for simplicity of the loadings and better interpretability. This calls for criteria that are able to evaluate the performance of correlated and/or nonorthogonal systems of components. We have shown in this

article that the existing criteria are not appropriate for this, because they do not handle correlations and lack of orthogonality in adequate ways. The example in Section 5 and other examples analyzed by Gervini and Rousson (2003) reveal that these criteria often assign full or almost full optimality to systems that are too far from the principal components. In contrast, the new criteria proposed in Section 4 can discriminate well between “good” and “bad” suboptimal systems. Of these two criteria we tend to favor CSV, but the examples in Gervini and Rousson (2003) show that both criteria are consistent in their evaluations if the systems are not too far from optimal. For these reasons, we think that our proposals are a significant improvement over existing criteria.

[Received May 2003. Revised September 2003.]

REFERENCES

- Cadima, J., and Jolliffe, I. T. (1995), “Loadings and Correlations in the Interpretation of Principal Components,” *Journal of Applied Statistics*, 22, 203–214.
- (2001), “Variable Selection and the Interpretation of Principal Subspaces,” *Journal of Agricultural, Biological, and Environmental Statistics*, 6, 62–79.
- Gervini, D., and Rousson, V. (2003), “Criteria for Evaluating Dimension-Reducing Components for Multivariate Data,” unpublished technical report.
- Hausmann, R. (1982), “Constrained Multivariate Analysis,” in *Optimisation in Statistics*, eds. S. H. Zacks and J. S. Rustagi, Amsterdam: North-Holland, pp. 137–151.
- Hotelling, H. (1933), “Analysis of a Complex of Statistical Variables into Principal Components,” *Journal of Educational Psychology*, 24, 417–441 and 498–520.
- Jackson, J. E. (1991), *A User's Guide to Principal Components* (1st ed.), New York: Wiley.
- Jolliffe, I. T. (1995), “Rotation of Principal Components: Choice of Normalization Constraints,” *Journal of Applied Statistics*, 22, 29–35.
- (2002), *Principal Component Analysis* (2nd ed.), New York: Springer.
- Jolliffe, I. T., and Uddin, M. (2000), “The Simplified Component Technique: An Alternative to Rotated Principal Components,” *Journal of Computational and Graphical Statistics*, 9, 689–710.
- Kaiser, H. F. (1958), “The Varimax Criterion for Analytic Rotation in Factor Analysis,” *Psychometrika*, 23, 187–200.
- Kiers, H. A. L. (1991), “Simple Structure in Component Analysis Techniques for Mixtures of Qualitative and Quantitative Variables,” *Psychometrika*, 56, 197–212.
- McCabe, G. P. (1984), “Principal Variables,” *Technometrics*, 26, 137–145.
- Neuhaus, J., and Wrigley, C. (1954), “The Quartimax Method: An Analytical Approach to Orthogonal Simple Structure,” *British Journal of Statistical Psychology*, 7, 81–91.
- Okamoto, M. (1969), “Optimality of Principal Components,” in *Multivariate Analysis II*, ed. P. R. Krishnaiah, New York: Academic Press, pp. 673–685.
- Rao, C. R. (1964), “The Use and Interpretation of Principal Component Analysis in Applied Research,” *Sankhya*, Ser. A, 26, 329–358.
- Robert, P., and Escoufier, Y. (1976), “A Unifying Tool for Linear Multivariate Statistical Methods: The RV-Coefficient,” *Applied Statistics*, 25, 257–265.
- Rousson, V. and Gasser, T. (2004), “Simple Component Analysis,” *Applied Statistics*.
- Vines, S. K. (2000), “Simple Principal Components,” *Applied Statistics*, 49, 441–451.