# On the usage of joint diagonalization in multivariate statistics

Klaus Nordhausen [a,*], Anne Ruiz-Gazen [b]

[a] *Department of Mathematics and Statistics, University of Jyväskylä, Finland*
[b] *Toulouse School of Economics, Université de Toulouse Capitole, France*

## ARTICLE INFO

## ABSTRACT

Scatter matrices generalize the covariance matrix and are useful in many multivariate data analysis methods, including well-known principal component analysis (PCA), which is based on the diagonalization of the covariance matrix. The simultaneous diagonalization of two or more scatter matrices goes beyond PCA and is used more and more often. In this paper, we offer an overview of many methods that are based on a joint diagonalization. These methods range from the unsupervised context with invariant coordinate selection and blind source separation, which includes independent component analysis, to the supervised context with discriminant analysis and sliced inverse regression. They also encompass methods that handle dependent data such as time series or spatial data.

## 1. Introduction

Classical multivariate analysis, such as that presented in Anderson [3] assumes that the data at hand follow a multivariate normal model. This is very convenient as the multivariate normal distribution is fully specified by its mean vector and covariance matrix and these two statistics suffice to develop tractable and optimal inference tools for this model. Early on, it was known that statistical methods based on the mean vector and covariance matrix are very sensitive to atypical observations in the data and are not very efficient for observations coming from a heavy-tailed distribution. To alleviate these problems, the normal model is commonly broadened to the elliptical model, which keeps the shape of the probability contours but allows for one additional kurtosis parameter, thus allowing heavier and lighter tails than the normal model. For robustness and optimality reasons, alternative location measures for the mean vector and alternative dispersion measures for the covariance matrix were developed in the elliptical framework. These measures are often expected to have certain properties under affine transformations of the data, in which case they are called location functionals $T$ and scatter functionals $S$. It can then be shown that in an elliptical model, all location functionals, including the mean vector, correspond to the symmetry centre and that all scatter functionals are proportional to the covariance matrix, if they exist [85,87]. Thus, scatter functionals measure the same population quantity in the elliptical model, and it

---

* Corresponding author.
 *E-mail address:* klaus.k.nordhausen@jyu.fi (K. Nordhausen).

is sufficient to use one location functional and one scatter functional for inference purposes (A slightly larger model where location and scatter functionals measure the same population quantities is also discussed in [85,87].) Since approximately the turn of the last century, interest in the simultaneous use of two or more scatter matrices, which is of course most interesting when these functionals do not measure the same population quantities, has increased.

In the present paper we will show how two or more scatter functionals are jointly used in multivariate statistics and in which models this is of interest. For this purpose we recall first in Section 2 the concept of scatter functionals in detail and discuss some of their properties. Section 3 gives details on the simultaneous and joint diagonalization of scatter functionals which is the main tool used in our context. Invariant coordinate selection (ICS) is discussed in Section 4, blind source separation (BSS) is discussed in Section 5 and the use of joint diagonalization in the context of supervised dimension reduction (SDR) methods is discussed in Section 6. Finally, the paper is concluded in Section 7.

## 2. Scatter matrices

Joint diagonalization has been used in unsupervised and supervised contexts and for independent and dependent data.

In the unsupervised case with independent data, the definition of a scatter matrix, sometimes also called pseudo-covariance, is a generalization of the covariance matrix definition (see [26,35,54,85,102] among others).

Following [85], let us first define the functional version of a scatter estimator. For a $p$-dimensional vector $\boldsymbol{X}$ with distribution function $F_{\boldsymbol{X}}$, a functional $\boldsymbol{S}(F_{\boldsymbol{X}})$ also denoted by $\boldsymbol{S}(\boldsymbol{X})$ is called a scatter functional if it is a $p \times p$ symmetric positive semidefinite and affine equivariant matrix. Note that in [102], the definition is more stringent than that in [85], and assumes that a scatter matrix is positive definite. We recall that an affine equivariant matrix $\boldsymbol{S}(\boldsymbol{X})$ is such that

$$\boldsymbol{S}(\boldsymbol{AX} + \boldsymbol{b}) = \boldsymbol{AS}(\boldsymbol{X})\boldsymbol{A}^{\top},$$

where $^{\top}$ denotes the transpose operator, $\boldsymbol{A}$ is a full rank $p \times p$ matrix and $\boldsymbol{b}$ a $p$-vector.

For distributions $F_{\boldsymbol{X}}$ with finite second moments, the covariance functional is defined by:

$$\mathrm{Cov}(\boldsymbol{X}) = E\left[(\boldsymbol{X} - E(\boldsymbol{X}))(\boldsymbol{X} - E(\boldsymbol{X}))^{\top}\right]$$

and is affine equivariant.

Let us now consider the empirical version of a scatter functional. This means that we have a $p$-variate dataset $\boldsymbol{X}_n = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)^{\top}$ and the scatter functional $\boldsymbol{S}(F_n)$ for the empirical distribution $F_n$. A scatter matrix statistic or estimator is thus a $p \times p$ symmetric positive semidefinite and affine equivariant matrix. In this framework, an affine equivariant matrix $\boldsymbol{S}(\boldsymbol{X}_n)$ is such that

$$\boldsymbol{S}(\boldsymbol{X}_n\boldsymbol{A} + \boldsymbol{1}_n\boldsymbol{b}^{\top}) = \boldsymbol{A}^{\top}\boldsymbol{S}(\boldsymbol{X}_n)\boldsymbol{A},$$

where $\boldsymbol{A}$ is a full rank $p \times p$ matrix, $\boldsymbol{b}$ a $p$-vector and $\boldsymbol{1}_n$ an $n$-vector full of ones.

The empirical covariance matrix is defined by:

$$\mathrm{Cov}(\boldsymbol{X}_n) = \frac{1}{n}\sum_{i=1}^{n}(\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)(\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)^{\top}$$

where $\bar{\boldsymbol{x}}_n = 1/n\sum_{i=1}^{n}\boldsymbol{x}_i$ is the empirical mean. The mean is an affine equivariant location estimator $\boldsymbol{T}$ such that:

$$\boldsymbol{T}(\boldsymbol{AX} + \boldsymbol{b}) = \boldsymbol{AT}(\boldsymbol{X}) + \boldsymbol{b},$$

for the functional version and

$$\boldsymbol{T}(\boldsymbol{X}_n\boldsymbol{A}^{\top} + \boldsymbol{1}_n\boldsymbol{b}^{\top}) = \boldsymbol{AT}(\boldsymbol{X}_n) + \boldsymbol{b},$$

for the empirical version where $\boldsymbol{A}$ is a full rank $p \times p$ matrix and $\boldsymbol{b}$ a $p$-vector.

For elliptical distributions with second moments, scatter functionals are proportional to the covariance matrix (see, e.g., [9]).

Many scatter matrices have been defined with the objective of making the covariance matrix estimator more robust (see, e.g., [35,53]). Tyler et al. [102] divide the scatter matrices in three classes depending on their robustness properties. The first class includes scatter estimators with a zero breakdown point such as the usual covariance matrix but also the one-step M-estimators with a functional defined by:

$$\mathrm{Cov}_w(\boldsymbol{X}) = E\left[w(D^2(\boldsymbol{X}))(\boldsymbol{X} - E(\boldsymbol{X}))(\boldsymbol{X} - E(\boldsymbol{X}))^{\top}\right],$$

where $w$ is a non-negative and continuous weight function and $D^2(\boldsymbol{X}) = (\boldsymbol{X} - E(\boldsymbol{X}))^{\top}\mathrm{Cov}(\boldsymbol{X})^{-1}(\boldsymbol{X} - E(\boldsymbol{X}))$ is the Mahalanobis distance. The sample version of the one-step M-estimator is:

$$\mathrm{Cov}_w(\boldsymbol{X}_n) = \frac{1}{n}\sum_{i=1}^{n}w(D^2(\boldsymbol{x}_i))(\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)(\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)^{\top},$$

where $D^2(\boldsymbol{x}_i) = (\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)^{\top}\mathrm{Cov}(\boldsymbol{X}_n)^{-1}(\boldsymbol{x}_i - \bar{\boldsymbol{x}}_n)$.

The covariance matrix is obtained with $w(d) = 1$ while we get the $\text{Cov}_{-1}$ matrix defined by [25] when $w(d) = 1/d$. As noticed by [70], when $w(d) = d^\alpha$ with $\alpha < 0$, such estimators down-weight values with large Mahalanobis distance and so have a robust flavour even if they have a zero breakdown point. From the same class, the fourth-moment based estimator $\text{Cov}_4$ obtained with $w(d) = d$ is widely used in the blind source separation literature (see, e.g., [86,100]). It is highly nonrobust since it up-weights values with large Mahalanobis distances but it proves to be useful in particular situations.

The second class of estimators contains scatter matrices with a moderate breakdown point such as the class $(\boldsymbol{T}, \boldsymbol{S})$ of M-estimators that are defined (see, e.g., [52]) as solutions of systems of equations of the following form:

$$E\left[u_1\left[(\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))^\top \boldsymbol{S}(\boldsymbol{X})^{-1}(\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))\right](\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))\right] = \boldsymbol{0} \tag{1}$$

$$E\left[u_2\left[(\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))^\top \boldsymbol{S}(\boldsymbol{X})^{-1}(\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))\right](\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))(\boldsymbol{X} - \boldsymbol{T}(\boldsymbol{X}))^\top\right] = \boldsymbol{S}(\boldsymbol{X}) \tag{2}$$

where $u_1$ and $u_2$ are appropriate weight functions. The sampling version of M-estimators is easily derived from the previous equations.

The third class contains high breakdown point estimators such as S-estimators or minimum covariance determinant estimators (see [53] for details and other scatter estimators from the same class).

In some statistical methods, such as independent component analysis, a scatter matrix is also expected to verify the joint independence property or the block independence property. A scatter functional $\boldsymbol{S}(\boldsymbol{X})$ has the joint independence property if for a vector with mutually independent components $\boldsymbol{S}(\boldsymbol{X})$ is diagonal. In the class of one-step M-estimators with nonnegative and continuous weight function, [103] proves that the only scatter functionals with the independence property are the Cov and $\text{Cov}_4$ estimators and their nonnegative linear combinations. In the case where all components of $\boldsymbol{X}$ are not necessarily independent but consist of independent blocks of components, the block dependence property states that the mutually independent subvectors of $\boldsymbol{X}$ correspond to diagonal submatrices leading to a block diagonal scatter matrix (see [85,102] for more details).

Note that it is also possible to define some symmetrized version of the previous scatter matrices by considering $\boldsymbol{S}(\boldsymbol{U} - \boldsymbol{V})$, where $\boldsymbol{U}$ and $\boldsymbol{V}$ are independent copies of $\boldsymbol{X}$ (see [85,102]). In other words, the symmetrization is obtained by applying the scatter functional to pairwise differences. The symmetrized scatter matrices possess the joint and block independence property (see [85,88]).

In [46], the definition of a scatter matrix is extended to the context of supervised methods. In this context, in addition to the $p$-vector $\boldsymbol{X}$, a response variable $Y$ is available. Following [46], a supervised scatter functional $\boldsymbol{S}$ is a function of the joint distribution $F_{\boldsymbol{X},Y}$ of $(\boldsymbol{X}, Y)$ which is affine equivariant in the sense that

$$\boldsymbol{S}(F_{\boldsymbol{AX}+\boldsymbol{b},Y}) = \boldsymbol{A}\boldsymbol{S}(F_{\boldsymbol{X},Y})\boldsymbol{A}^\top,$$

for all full rank matrices $\boldsymbol{A}$ and all $p$-vectors $\boldsymbol{b}$. One example of such a supervised scatter functional is:

$$\boldsymbol{S}_{SIR}(F_{\boldsymbol{X},Y}) = \text{Cov}(E(\boldsymbol{X}|Y)).$$

Note that in the case of a discrete response variable, $\boldsymbol{S}_{SIR}$ corresponds to the between covariance matrix $\text{Cov}_{\boldsymbol{B}}$.

Thus far, we have focused on samples of independent data but as will be detailed below, joint diagonalization is also widely used in the context of time series and spatial data. In such contexts, affine equivariant estimators that are not necessarily positive semidefinite are considered and go beyond the scatter matrix definition above.

In the context of $p$-variate time series, let us consider a stochastic process $\boldsymbol{X}_T = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_T)$ measured at time $t \in \{1, \ldots, T\}$. For a given lag $\tau \in \{0, 1, \ldots\}$, the sample version of the cross-autocovariance matrix $\text{ACov}_\tau(\boldsymbol{X}_T)$ is given by

$$\text{ACov}_\tau(\boldsymbol{X}_T) = \frac{1}{T - \tau}\sum_{t=1}^{T-\tau}(\boldsymbol{x}_t - \bar{\boldsymbol{x}}_T)(\boldsymbol{x}_{t+\tau} - \bar{\boldsymbol{x}}_T)^\top$$

where $\bar{\boldsymbol{x}}_T = 1/T \sum_{t=1}^T \boldsymbol{x}_t$. Note that $\text{ACov}_0(\boldsymbol{X}_T) = \text{Cov}(\boldsymbol{X}_T)$.

The $\text{ACov}_\tau$ matrix is not necessarily symmetric and is sometimes symmetrized when it is expected to be symmetric for the model under consideration (see [61,101]). A symmetrized version of $\text{ACov}_\tau$ is defined by

$$\text{ACov}_\tau^S = \frac{1}{2}(\text{ACov}_\tau + \text{ACov}_\tau^\top).$$

Let us now consider multivariate data measured at spatial locations $\boldsymbol{s}_1, \ldots, \boldsymbol{s}_n$ in a domain $\mathcal{S} \subseteq \mathbb{R}^d$, $\boldsymbol{X}_n = (\boldsymbol{x}(\boldsymbol{s}_1), \ldots, \boldsymbol{x}(\boldsymbol{s}_n))$. [7] define local covariance, or scatter, matrices, by:

$$\text{LCov}_f(\boldsymbol{X}_n) = \frac{1}{n}\sum_{i=1}^n\sum_{j=1}^n f(\boldsymbol{s}_i - \boldsymbol{s}_j)(x(\boldsymbol{s}_i) - \bar{\boldsymbol{x}}_n)(x(\boldsymbol{s}_j) - \bar{\boldsymbol{x}}_n)^\top, \tag{3}$$

where $\bar{\boldsymbol{x}}_n = 1/n\sum_{i=1}^n \boldsymbol{x}(\boldsymbol{s}_i)$ and $f : \mathbb{R}^d \to \mathbb{R}$ is called the kernel function. Examples of kernels $f$ are the ball and ring kernels $B(h)(\boldsymbol{s}) = I(\|\boldsymbol{s}\| \leq h)$ with fixed $h \geq 0$ and $R(h_1, h_2)(\boldsymbol{s}) = I(h_1 \leq \|\boldsymbol{s}\| \leq h_2)$ with fixed $h_2 \geq h_1 \geq 0$ where $\|.\|$ denotes the euclidian norm and $I(\cdot)$ denotes the indicator function.

Note that the ACov and LCov are considered as scatter matrices but may not be semipositive definite.

## 3. Simultaneous and joint diagonalization

In PCA (see, e.g., [40]), the covariance $\text{Cov}(\boldsymbol{X}_n)$ (or correlation) matrix, which is a symmetric real valued matrix, is diagonalized. It means that the following transformation is calculated:

$$\boldsymbol{U}(\boldsymbol{X}_n)\text{Cov}(\boldsymbol{X}_n)\boldsymbol{U}(\boldsymbol{X}_n)^{\top} = \Lambda(\boldsymbol{X}_n)$$

where $\Lambda(\boldsymbol{X}_n) = \text{diag}(\lambda_1 \geq \cdots \geq \lambda_p)$ is the diagonal matrix containing the ordered eigenvalues of $\text{Cov}(\boldsymbol{X}_n)$ and $\boldsymbol{U}(\boldsymbol{X}_n) = (\boldsymbol{u}_1, \ldots, \boldsymbol{u}_p)^{\top}$ contains its corresponding orthonormal eigenvectors as rows. Because the matrix $\text{Cov}(\boldsymbol{X}_n)$ is symmetric, the matrix $\boldsymbol{U}(\boldsymbol{X}_n)$ is orthogonal with respect to the usual inner-product and can be chosen such that $\boldsymbol{U}(\boldsymbol{X}_n)\boldsymbol{U}(\boldsymbol{X}_n)^{\top} = \boldsymbol{U}(\boldsymbol{X}_n)^{\top}\boldsymbol{U}(\boldsymbol{X}_n) = \boldsymbol{I}_p$. This procedure is also called the spectral or the eigenvalue–eigenvector decomposition or eigendecomposition of $\text{Cov}(\boldsymbol{X}_n)$.

When considering two scatter matrices $\boldsymbol{S}_1(\boldsymbol{X}_n)$ and $\boldsymbol{S}_2(\boldsymbol{X}_n)$, it is possible to find a matrix $\boldsymbol{W}(\boldsymbol{X}_n)$ such that both matrices are transformed into diagonal matrices:

$$\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_1(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top} = \Lambda_1(\boldsymbol{X}_n) \text{ and } \boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_2(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top} = \Lambda_2(\boldsymbol{X}_n) \tag{4}$$

where $\Lambda_1(\boldsymbol{X}_n)$ and $\Lambda_2(\boldsymbol{X}_n)$ are diagonal matrices (see e.g., [102]).

This procedure is called simultaneous diagonalization (see, e.g., [95]). It leads to the diagonalization of $\boldsymbol{S}_1(\boldsymbol{X}_n)^{-1}\boldsymbol{S}_2(\boldsymbol{X}_n)$ which is not necessarily symmetric:

$$\boldsymbol{S}_1(\boldsymbol{X}_n)^{-1}\boldsymbol{S}_2(\boldsymbol{X}_n) = \boldsymbol{W}(\boldsymbol{X}_n)^{\top}\Lambda(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{-1}$$

with $\Lambda(\boldsymbol{X}_n) = \Lambda_1^{-1}(\boldsymbol{X}_n)\Lambda_2(\boldsymbol{X}_n)$.

Generally, $\Lambda_1(\boldsymbol{X}_n)$ is taken as the identity matrix and we focus on this particular case from now on. In this case, Problem (4) is equivalent to the diagonalization of $\boldsymbol{S}_2(\boldsymbol{X}_n)$ with a matrix of eigenvectors $\boldsymbol{W}(\boldsymbol{X}_n)$ that is orthogonal with respect to the inner product induced by $\boldsymbol{S}_1(\boldsymbol{X}_n)$ instead of the canonical inner product.

It is easy to see that Problem (4) is equivalent to the usual diagonalization of $\boldsymbol{S}_1(\boldsymbol{X}_n)^{-1/2}\boldsymbol{S}_2(\boldsymbol{X}_n)\boldsymbol{S}_1(\boldsymbol{X}_n)^{-1/2}$ which is a symmetric matrix with ordered eigenvalues given by $\Lambda_2(\boldsymbol{X}_n)$ and orthonormal eigenvectors given by $\boldsymbol{S}_1(\boldsymbol{X}_n)^{1/2}\boldsymbol{W}(\boldsymbol{X}_n)$.

Finally, Problem (4) is also equivalent to the problem of finding values $\lambda_i(\boldsymbol{X})$ and vectors $\boldsymbol{w}_i(\boldsymbol{X})$, $i \in \{1, \ldots, p\}$, such that

$$\boldsymbol{S}_2(\boldsymbol{X}_n)\boldsymbol{w}_i(\boldsymbol{X}) = \lambda_i(\boldsymbol{X})\boldsymbol{S}_1(\boldsymbol{X}_n)\boldsymbol{w}_i(\boldsymbol{X})$$

which is called the generalized eigendecomposition problem.

This simultaneous diagonalization procedure is used in different contexts and takes different names depending on the context. For instance, when using the scatter matrices Cov and $\text{Cov}_4$, the method is called FOBI in the signal processing literature (see e.g., [86] and Section 5.1). When considering the usual covariance matrix and one autocovariance matrix ACov, it is called AMUSE in the time series context (see e.g., [89] and Section 5.2).

Going beyond two scatter matrices is more challenging but has also been studied in the literature. We will call the procedure "joint diagonalization" as soon as the number of scatter matrices is larger than two. Let us consider $\boldsymbol{S}_0(\boldsymbol{X}), \boldsymbol{S}_1(\boldsymbol{X}), \ldots, \boldsymbol{S}_K(\boldsymbol{X})$, i.e., $K + 1$ scatter matrices associated with a random vector $\boldsymbol{X}$. It is known that, for such a collection of symmetric matrices, there exists a matrix $\boldsymbol{P}(\boldsymbol{X})$, such that $\boldsymbol{P}(\boldsymbol{X})\boldsymbol{S}_k(\boldsymbol{X})\boldsymbol{P}(\boldsymbol{X})^{\top}$ is diagonal, for each $k \in \{0, \ldots, K\}$, if and only if all pairs of scatter matrices commute (see [95]).

In the blind source separation model (see Section 5), the assumption that the scatter matrices commute is true, and the joint diagonalization is possible. However, when considering the sampling versions of the scatter matrices, the property is lost. In such a situation, we can try to make the matrices jointly "as diagonal as possible" (see, e.g., [21,32,58]).

One idea is to take one of the scatter matrices, let us say $\boldsymbol{S}_0(\boldsymbol{X}_n)$, as a reference and to find a transformation $\boldsymbol{W}(\boldsymbol{X}_n)$ such that $\boldsymbol{S}_0(\boldsymbol{X}_n)$ is diagonalized while the other scatter matrices, $\boldsymbol{S}_1(\boldsymbol{X}_n), \ldots, \boldsymbol{S}_K(\boldsymbol{X}_n)$ are only approximately diagonalized. A popular criterion is based on a least squares approach. More precisely, the criterion consists in looking for $\boldsymbol{W}(\boldsymbol{X}_n)$ which minimizes the sum of squares of the off-diagonal elements of all possible scatter matrices after transformation, and under the constraint that $\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_0(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top}$ is the identity

$$\min \sum_{k=1}^{K} \|\text{off}(\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_k(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top})\|^2, \text{ subject to } \boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_0(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top} = \boldsymbol{I}_p$$

where $\text{off}(\boldsymbol{A}) = \boldsymbol{A} - \text{diag}(\boldsymbol{A})$, for a square matrix $\boldsymbol{A}$, and $\|\cdot\|$ denotes the matrix Frobenius norm.

This criterion is equivalent to maximizing the following sum

$$\sum_{k=1}^{K} \|\text{diag}(\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_k(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top})\|^2, \text{ subject to } \boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_0(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^{\top} = \boldsymbol{I}_p. \tag{5}$$

Popular algorithms for this approximate joint diagonalization are based on Jacobi rotations (see, e.g., [13,21,58]). In the multivariate time series context, with $\boldsymbol{S}_0(\boldsymbol{X}_n) = \text{Cov}(\boldsymbol{X}_n)$ and, for $\boldsymbol{S}_1(\boldsymbol{X}_n), \ldots, \boldsymbol{S}_K(\boldsymbol{X}_n)$, autocovariance matrices with different lags, this algorithm is called SOBI [59]. Other possible algorithms are also discussed in [36]. In particular, [62]

introduced a deflation-based algorithm such that the single rows of the matrix $\boldsymbol{W}(\boldsymbol{X}_n)$ are calculated one after the other by looking at a maximization problem similar to (5) but for each row of $\boldsymbol{W}(\boldsymbol{X}_n)$. The existence and the uniqueness of the solution are discussed.

There exist also several other proposals (see, e.g., [36,58]). If there is no reason that $\boldsymbol{S}_0(\boldsymbol{X}_n)$ plays a special role, it is possible to replace the constraint $\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_0(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^\top = \boldsymbol{I}_p$ by other constraints (see, e.g., [107,108]). Moreover, the square in the maximization criterion can be replaced by other functions allowing to weight differently the involved scatter matrices (see [58,62] for details).

In the rest of this article, we present many existing multivariate data analysis methods that use a simultaneous or approximate joint diagonalization of scatter matrices.

## 4. Invariant coordinate selection

There are two popular data transformations based on a single scatter matrix. The first one is principal component analysis, which uses the transformation matrix as the orthogonal matrix $\boldsymbol{U}$ obtained via the eigenvalue–eigenvector decomposition of $\boldsymbol{S}(\boldsymbol{X}_n) = \boldsymbol{U}(\boldsymbol{X}_n)^\top \boldsymbol{D}(\boldsymbol{X}_n)\boldsymbol{U}(\boldsymbol{X}_n)$ where the $i$th row of $\boldsymbol{U}(\boldsymbol{X}_n)$ contains the $i$th eigenvector of $\boldsymbol{S}(\boldsymbol{X}_n)$ and the diagonal matrix $\boldsymbol{D}(\boldsymbol{X}_n)$ contains on its diagonal the corresponding eigenvalues for which we assume that they are ordered in descending order. The principal components are the observations projected along the principal vectors, i.e. $\boldsymbol{z}_i = \boldsymbol{U}(\boldsymbol{X}_n)\boldsymbol{x}_i$, $i \in \{1, \ldots, n\}$, where it is often assumed that the observations are centred. The principal components then have the property that they are uncorrelated with respect to $\boldsymbol{S}(\boldsymbol{X}_n)$, i.e., $\boldsymbol{S}(\boldsymbol{Z}_n) = \boldsymbol{D}(\boldsymbol{X}_n)$. Traditional PCA is based on the regular covariance matrix (see, e.g., [40] for details); however, within an elliptical distribution framework any scatter matrix can be used for the same purpose.

Another transformation is the so-called whitening transformation which, besides the scatter $\boldsymbol{S}(\boldsymbol{X}_n)$, needs a location $\boldsymbol{T}(\boldsymbol{X}_n)$. Whitened observations are obtained as

$$\boldsymbol{x}_i^{st} = \boldsymbol{S}(\boldsymbol{X}_n)^{-1/2}(\boldsymbol{x}_i - \boldsymbol{T}(\boldsymbol{X}_n)),$$

$i \in \{1, \ldots, n\}$. These observations have the properties that $\boldsymbol{T}(\boldsymbol{X}_n^{st}) = \boldsymbol{0}$ and $\boldsymbol{S}(\boldsymbol{X}_n^{st}) = \boldsymbol{I}_p$ which means that compared to PCA, whitened transformations not only uncorrelate the components but also give them equal scales. Note however, that the whitened components are not necessarily just the scaled principle components but might have undergone an additional rotation. Actually, Ilmonen et al. [39] mention five alternative ways to compute $\boldsymbol{S}(\boldsymbol{X}_n)^{-1/2}$ which might all differ by an orthogonal rotation. If not specified otherwise, we consider the symmetric variant $\boldsymbol{S}(\boldsymbol{X}_n)^{-1/2} = \boldsymbol{U}(\boldsymbol{X}_n)\boldsymbol{D}(\boldsymbol{X}_n)^{-1/2}\boldsymbol{U}(\boldsymbol{X}_n)^\top$ with $\boldsymbol{U}(\boldsymbol{X}_n)$ and $\boldsymbol{D}(\boldsymbol{X}_n)$ as above. Again, this transformation usually uses the regular mean vector and the covariance matrix but other locations and scatter functionals can also be used, in which case an elliptical model is tacitly assumed.

One of the first ideas regarding the use of two scatter matrices was then based on performing these two transformations one after the other, but using a different scatter matrix for each one. The idea is then, ignoring the location for a moment, that the data are first whitened with respect to a scatter $\boldsymbol{S}_1$ and then PCA is performed on the whitened data using another scatter $\boldsymbol{S}_2$. This can be formulated as the simultaneous diagonalization problem of finding the transformation matrix $\boldsymbol{W}(\boldsymbol{X}_n)$ such that

$$\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_1(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^\top = \boldsymbol{I}_p, \quad \boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{S}_2(\boldsymbol{X}_n)\boldsymbol{W}(\boldsymbol{X}_n)^\top = \boldsymbol{D}(\boldsymbol{X}_n),$$

where $\boldsymbol{D}(\boldsymbol{X}_n)$ is a diagonal matrix with decreasing elements. Note that in the following, when the context is clear, we drop the dependence on $\boldsymbol{X}_n$ for $\boldsymbol{W}, \boldsymbol{D}, \boldsymbol{S}_1$ and $\boldsymbol{S}_2$. Based on Section 3, it is clear that this is a generalized eigenvalue–eigenvector problem and $\boldsymbol{W}$ and $\boldsymbol{D}$ can be computed accordingly.

Thus, in a model-free context, this transformation can be considered as an investigation if, after removing the second order information as measured by $\boldsymbol{S}_1, \boldsymbol{S}_2$ can still find any structure in the data, which is, for example, not the case when the observations follow an elliptical distribution.

This transformation was first denoted generalized PCA in [14–16] but the more commonly acknowledged name at present is invariant coordinate selection (ICS) as established in [102]. Note that some special scatter combinations are considered under specific names. For example the combination $\boldsymbol{S}_1 = \text{Cov}$ and $\boldsymbol{S}_2 = \text{Cov}_{-1}$ is known as principal axis analysis [25] and the combination $\boldsymbol{S}_1 = \text{Cov}$ and $\boldsymbol{S}_2 = \text{Cov}_4$ is known as fourth order blind identification (FOBI) [12] which may be one of the most popular combinations.
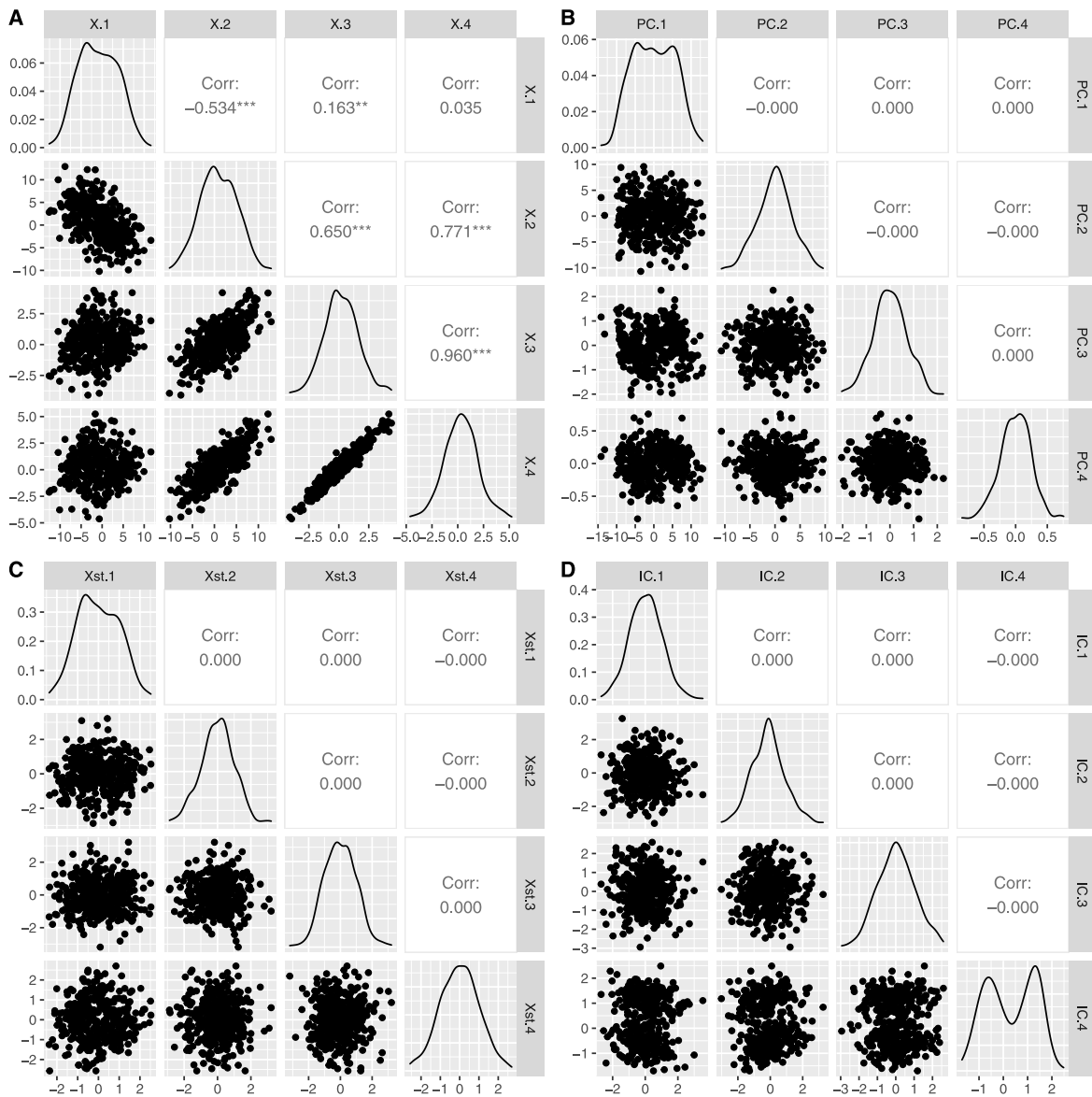
The name ICS is motivated based on the following equivariance property which holds when the elements of $\boldsymbol{D}$ are all distinct:

$$\boldsymbol{W}(\boldsymbol{X}_n)\boldsymbol{A}^{-1} = \boldsymbol{J}\boldsymbol{W}(\boldsymbol{X}_n\boldsymbol{A}^\top + \boldsymbol{1}_n\boldsymbol{b}^\top),$$

where $\boldsymbol{A}$ is a $p \times p$ matrix and $\boldsymbol{b}$ a p-vector. $\boldsymbol{J}$ denotes a sign change matrix, i.e., a diagonal matrix with $\pm 1$ on its diagonal. Thus, in connection with a location functional $\boldsymbol{T}$,

$$\boldsymbol{W}(\boldsymbol{X}_n)\,(\boldsymbol{x}_i - \boldsymbol{T}(\boldsymbol{X}_n)) = \boldsymbol{J}\boldsymbol{W}(\boldsymbol{X}_n\boldsymbol{A}^\top + \boldsymbol{1}_n\boldsymbol{b}^\top)\big((\boldsymbol{A}\boldsymbol{x}_i + \boldsymbol{b}) - \boldsymbol{T}(\boldsymbol{X}_n\boldsymbol{A}^\top + \boldsymbol{1}_n\boldsymbol{b}^\top)\big),$$

which means that the so called ICS-components $\boldsymbol{z}_i = \boldsymbol{W}(\boldsymbol{X}_n)\,(\boldsymbol{x}_i - \boldsymbol{T}(\boldsymbol{X}_n))$ are affine invariant under linear transformations up to their signs. Therefore one can argue that the ICS components show the intrinsic structure of the data independent of the coordinate system in which the data were originally presented. This is quite different from the principal components

**Fig. 1.** Comparison of different data transformations. Panel A shows a matrix scatterplot with density estimators on the diagonal and correlations for the original data with $p = 4$, Panel B the corresponding principal components, Panel C the whitened components and Panel D the invariant coordinates where $S_1 = \text{Cov}$ and $S_2 = \text{Cov}_4$. Therefore all four panels give different views of the same dataset where in the ICS representation the two groups are best visible.

and the whitened observations which do not have this type of invariance property. The differences between the transformations are illustrated in Fig. 1 where the observable 4-variate data are shown together with the corresponding principal components, whitened components and invariant coordinates, which are in this case based on FOBI. In this artificial example, the two latent clusters are best visible in the invariant coordinates (see the two modes of the density estimator on the 4th plot of the diagonal of panel D).

In their seminal paper Tyler et al. [102] also provide an interpretation of the eigenvalues contained in the diagonal of $D$. Let $w$ be a $p$-vector. Then, one can consider $w^{\top} S_1(X) w$ as a squared measure of the scale of $X$ in the direction of $w$. As the ratio of two squared scale measures can be seen as a kurtosis measure [94],

$$\kappa(w^{\top} X) = \frac{w^{\top} S_1(X) w}{w^{\top} S_2(X) w}$$

is therefore the kurtosis measured in the $S_1 - S_2$-sense of $w^\top X$. Hence, the diagonal elements in $D$ are the (ordered) kurtosis values of the invariant coordinates. Tyler et al. [102] show that the maximal/minimal kurtosis in the $S_1 - S_2$-sense that can be obtained for $X$, corresponds to the first/last eigenvalue. Thus Tyler et al. [102] state ICS can be "viewed as a projection pursuit without the pursuit effort". If $S_1$ and $S_2$ are normalized for the Gaussian distribution, then an eigenvalue $d_j = 1, j \in \{1, \ldots, p\}$, can be taken as an indicator that the component $z_j$ follows a Gaussian distribution.

Based on the invariance property of ICS and the properties of the eigenvalues ICS has been used for many purposes, mainly in an exploratory data analysis way.

### 4.1. ICS for descriptive statistics

As the ICS components show the intrinsic nature of the data, they are a natural start to describe the basic data features. Nordhausen et al. [78] actually suggest using an additional second location $T_2$ and fixing the sign of the $j$th component of $z_i$ so that $(T_2(Z_n))_j \geq 0$, which means that the difference in the locations is a measure of skewness of the components. Thus, if $T_2(Z_n) \approx 0$ the data are symmetric and if all eigenvalues of $D$ are the same, it is an indication of ellipticity. Similarly the eigenvalues can give indications, together with the skewness measure, for other multivariate models such as skew-elliptical models as discussed, for example, in [47,78]. A more formal inference framework is given in Ilmonen et al. [37] where the limiting distributions of $W$, $D$ and $T_1 - T_2$ are derived, based on the location and scatter functionals, especially when they are moment based. Kankainen et al. [41] developed tests for multivariate normality based on these ideas when using the pair of scatter matrices Cov and $\mathrm{Cov}_4$ (FOBI).

### 4.2. ICS for dimension reduction and outlier detection

Recently, datasets have been increasing in dimension and in sample size. Standard assumptions in modern multivariate statistics are such that datasets containing considerable noise and relevant features can be concentrated in a much smaller signal subspace. The goal of dimension reduction is then to estimate the signal subspace whose dimension is usually unknown. The question is how to define what makes a signal. For example, PCA says that the signal subspace is the one that contains most variation in the data, and there are many rules how to choose the subspace dimension (see for example [40]). PCA is likely the most commonly used dimension reduction method and seems to be quite successful in practice. It is from a theoretical point, however difficult, to argue why the directions in which, for example, groups are to be separated or outliers are to be identified, should be those with large variation. The construction of counterexamples is quite easy. Kurtosis, on the other hand, is a natural indicator of non-Gaussianity and, is one of the most popular projection pursuit indices [34]. In a mixture model framework, the classical kurtosis measure depends on the mixing proportion. If the group sizes are approximately equal, the kurtosis is small, while for unbalanced groups, with the extreme case of outliers, the kurtosis will be large. Thus, the eigenvalues contained in the matrix $D$ give an indicator of interestingness of the components and allow, for example, the search for groups. In the above example, as shown in Fig. 1, the last invariant coordinate is most interesting as the groups in this example are of equal size. This makes component selection slightly more challenging, as first and last components might be of interest. This is different from PCA where usually only the first few components are of interest. Tyler et al. [102] show that in a general framework with a mixture of elliptical distributions, ICS will find Fisher's linear discriminant without knowing the class labels, and ICS was considered a method for dimension reduction prior to group identification, for example, in [2,27,28,90,102].

Similarly, the reduction of the dimension to make outlier detection easier via ICS was considered in [5,6,80], especially in the context of reliability when it is known that the proportion of outliers is small. [5] show that it is easier to identify the outliers when they can be captured in a few invariant coordinates. If all invariant coordinates need to be selected for outlier detection, then the method corresponds to Mahalanobis-type outlier detection approach where the Mahalanobis distances are computed with respect to $S_1$.

The determination of which and how many components to retain is still often done visually or based on heuristics. However, when assuming a non-Gaussian component analysis framework where it is assumed that the data can be decomposed into a non-Gaussian (signal) subspace that is independent of the remaining (noise) Gaussian subspace, formal inference about the subspace dimensions was discussed in the context of FOBI in [48,49,81,82], and for general scatter combinations, it was discussed in [93].

### 4.3. ICS as a transformation–retransformation method

As discussed above, in multivariate statistics, it is of key interest that the results of the analysis do not depend on the co-ordinate system used. Thus, estimates should have an appropriate equivariance property under affine transformations, and tests, for example, should be invariant. However, there are multivariate methods that are not affine equivariant/invariant, which is considered a major flaw. For example, multivariate methods based on marginal signs and ranks [see, for example, 91] suffer from this disadvantage. ICS can help in this context as a transformation–retransformation approach. This means that multivariate methods are applied to the invariant coordinates and, if required, retransformed to the original scale. This was discussed, for example, in [79,80].

As ICS is used with very different purposes in mind, Tyler et al. [102] argued that there is no general best scatter combination for $\boldsymbol{S}_1$ and $\boldsymbol{S}_2$. Depending on the problem and data at hand, the scatter matrices might require different properties. In general, however, Tyler et al. [102] argued that it would be advisable not to use two highly robust scatter matrices, as interesting features will not be detected when both scatter matrices focus too much on the same "inner" part of the data. Alashwali and Kent [2] argued that it might be advisable that both scatter matrices are computed with respect to the same location functional, especially when subsequent clustering is the goal. Which scatter functional is used first and which second is also of minor consequence. The effect of changing the order is to invert the eigenvalues and reversing the order of the components, which then also have different scales. A common convention is, for example, to choose the order so that $\boldsymbol{S}_1$ is more robust than $\boldsymbol{S}_2$ and that the ICS components are centred with respect to the location functional which goes most naturally with $\boldsymbol{S}_1$. For further discussions regarding invariant transformations, we refer to [39,96,97].

To apply ICS and related methods in R [92], packages ICS[80], ICSOutlier [6], ICSShiny [4] and ICtest [83] are available.

## 5. Blind source separation

ICS is often seen as a mainly exploratory tool for multivariate analysis. A more model based approach where joint diagonalization plays a major role is blind source separation (BSS). The basic BSS model is

$$\boldsymbol{X} = \boldsymbol{A}\boldsymbol{Z} + \boldsymbol{\mu},$$

where $\boldsymbol{X}$ is a $p$-variate observable phenomenon that is seen as a linear mixture of a somewhat standardized latent $p$-variate source $\boldsymbol{Z}$, where the mixing is represented by the full rank $p \times p$ matrix $\boldsymbol{A}$ and the location of $\boldsymbol{X}$ is specified by the $p$-vector $\boldsymbol{\mu}$. Standard assumptions are that $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$ which indicates that the components of $\boldsymbol{Z}$ are at least uncorrelated. The goal in BSS is to estimate $\boldsymbol{Z}$ based on a realized sample $\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n$ of $\boldsymbol{X}$ alone. The location $\boldsymbol{\mu}$ in this case is mainly considered as a nuisance parameter and, in the following, we assume for simplicity $\boldsymbol{\mu} = \boldsymbol{0}$.

Clearly, without further assumptions, it is not possible to solve the BSS problem, and there must be at least one further structural component given for $\boldsymbol{Z}$ that can be exploited. Different BSS models have been suggested in the literature, with different additional assumptions. The ones we will consider here are: (i) The observations are independent and identically distributed (iid), and the components of $\boldsymbol{Z}$ are independent and non-Gaussian. This case is known as independent component analysis (ICA). (ii) The observed data are a $p$-variate time series, and the components of the latent time series are uncorrelated or independent. Then, additional information that can be exploited is serial dependence. (iii) The observed data come from a $p$-variate spatial random field where the $p$ latent fields are again uncorrelated and independent, and the additional structure to be exploited is the spatial dependence.

Before going into detail, we point out that general overviews for BSS are, for example, [1,20,22,75] and that BSS approaches that are based on joint diagonalization are often called algebraic BSS methods.

All approaches make use of the following key result [64]. Let $\boldsymbol{X}^{st} = \mathrm{Cov}(\boldsymbol{X})^{-1/2}(\boldsymbol{X} - \boldsymbol{E}(\boldsymbol{X}))$ be the standardized version of $\boldsymbol{X}$, then

$$\boldsymbol{X}^{st} = \boldsymbol{U}^\top \boldsymbol{Z},$$

where $\boldsymbol{U}$ is some orthogonal $p \times p$ matrix. Thus, after whitening, the BSS problem can be reduced to the problem of finding an orthogonal matrix.

The strategy of all algebraic BSS methods described below makes use of the generalized concept of a scatter functional which only requires affine equivariance but relaxes the positive definiteness requirement. Then, the approach is to select $K \geq 1$ scatter functionals $\boldsymbol{S}_1, \ldots, \boldsymbol{S}_K$ for which

$$\boldsymbol{S}_i(\boldsymbol{Z}) = \boldsymbol{D}_i, \quad i \in i, \ldots, K,$$

holds with $\boldsymbol{D}_i$ being diagonal matrices. Thus, all the scatter matrices used are diagonal when computed for the sources. Then, the approach is to find the orthogonal matrix $\boldsymbol{U}$ such that:

$$\boldsymbol{U}\boldsymbol{S}_i(\boldsymbol{X}^{st})\boldsymbol{U}^\top = \boldsymbol{D}_i, \quad i \in i, \ldots, K,$$

which yields the unmixing matrix $\boldsymbol{W} = \boldsymbol{U}\mathrm{Cov}(\boldsymbol{X})^{-1/2}$. Thus, algebraic BSS methods consists of a joint diagonalization problem, where the unmixing matrix $\boldsymbol{W}$ diagonalizes the $K + 1$ scatter matrices $\mathrm{Cov}(\boldsymbol{X}), \boldsymbol{S}_1(\boldsymbol{X}), \ldots, \boldsymbol{S}_K(\boldsymbol{X})$ under the constraint that $\boldsymbol{W}\mathrm{Cov}(\boldsymbol{X})\boldsymbol{W}^\top = \boldsymbol{I}_p$.

In the following, we will discuss the different additional structural requirements made on the latent components and which scatter functionals are suitable.

Before this, let us discuss why we should perform BSS:

1. Often, it is assumed that the latent components have either physical meanings (BSS was suggested first in the signal processing literature) or that they are easier to interpret than the original components.
2. Another motivation is that often only a few components are considered interesting and the remainder noise, thus, it can be used for dimension reduction.
3. As the latent components are assumed uncorrelated or even independent, each component can be modelled in a univariate way, and instead of fitting a $p$-variate model, one could fit $p$ univariate models, which is often considered much simpler. For spatial data, such a benefit is demonstrated in [69] in the context of prediction.

### 5.1. Independent component analysis

ICA is the best known BSS approach. ICA methods are designed for iid data but are also often applied for dependent observations in which case, however, not all available information is exploited. In the ICA, it is assumed that:

**(IC1):** $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$

**(IC2):** The components of $\boldsymbol{Z} = (Z_1, \ldots, Z_p)$, are mutually independent.

**(IC3):** At most one component of $\boldsymbol{Z}$ is Gaussian.

The first algebraic ICA approach is FOBI [12], as described above, which means that $K = 1$ and $\boldsymbol{S}_1 = \mathrm{Cov}_4$ in the BSS framework. FOBI yields an unmixing matrix if all independent components have distinct kurtosis values. The statistical properties of FOBI, in the ICA model, were derived in Miettinen et al. [64]. [88] generalized then FOBI by showing that Cov and $\boldsymbol{S}_1$ in FOBI can be replaced by any scatter functionals that have the independence property, which is also further investigated in [77]. Therefore, one can say, given a suitable choice of scatter functionals, that ICS can also solve the ICA problem when kurtosis values of $\boldsymbol{Z}$, in the sense of the two scatter functionals involved, are distinct.

In ICA, however, assuming that the components have distinct kurtosis (in the sense of the involved scatter matrices) is considered as a strong constraint. Therefore, [72] suggested using $K + 1$ scatter functionals that have all the independence property, and then using joint diagonalization as described in Section 3 to obtain the unmixing matrix. This is more flexible in the sense that this solves the ICA problem, where for each component there is at least one scatter combination $\boldsymbol{S}_0, \boldsymbol{S}_j, j \in \{1, \ldots, K\}$ with a distinct kurtosis compared with other components. Thus, the K-Scatter ($K > 1$) approach is more flexible than the 2-scatter approach. However, this approach still cannot separate components that have the same distribution. An ICA approach that also works for identical distributed components and is based on joint diagonalization is JADE (joint diagonalization of eigenmatrices), which uses cumulant matrices. Therefore, the method does not fully fit in the framework presented here, as no scatter functionals are diagonalized but certain cumulant matrices. We refer the reader to [13,64] for further details. The 2-scatter ICA approach is often not optimal and FOBI is, for example, always less efficient than JADE [64]. The 2-scatter ICA approach, compared to the $K$-scatter approach or JADE, is usually much easier to compute, especially FOBI which is often the start of an ICA analysis.

An extension of ICA is independent subspace analysis (ISA), which is also known as multivariate ICA. In this framework $\boldsymbol{Z}$ does not have $p$ "univariate" independent components but only $c$ components which may be multivariate. Thus, $\boldsymbol{Z} = (\boldsymbol{Z}_1^\top, \ldots, \boldsymbol{Z}_c^\top)^\top$ where $\boldsymbol{Z}_i$ has dimension $p_i$, $i \in \{1, \ldots, c\}$ with $p_1 + \cdots + p_c = p$. In ISA, the individual components cannot be recovered but only their subspaces. An approach based on joint diagonalization is first to perform 2-scatter ICA, then compute a third scatter that has the block independence property for the obtained components, and finally blockdiagonalize this third scatter. For details, see for example [74].

### 5.2. BSS for time series

In ICA, the extra feature exploited was non-Gaussianity for BSS. In time series, serial dependence can be exploited, which in turn allows multiple Gaussian components. Different BSS approaches imply different assumptions regarding the time series. For simplicity, we continue using $\boldsymbol{X}$ for the stochastic process $\boldsymbol{X}_t$.

1. Second order source separation (SOS) makes the following model assumptions:

    **(SOS:)** $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$.
    **(SOS2:)** $\mathrm{ACov}_\tau(\boldsymbol{Z}) = \boldsymbol{D}_\tau$ for all $\tau \geq 1$ where $\boldsymbol{D}_\tau$ are all diagonal matrices.

    Therefore, in the SOS model, the latent components are uncorrelated throughout time and it is usually assumed that the latent components are linear processes.
    The first SOS method is known as AMUSE (algorithm for multiple unknown signals extraction) [101], which chooses $K = 1$ and $\boldsymbol{S}_1 = \mathrm{ACov}_\tau^S$ for some lag $\tau$, which is very often 1. AMUSE can solve the SOS problem if all autocorrelations at the used lag are distinct. AMUSE is very sensitive to the chosen lag and was extended to SOBI (second order blind identification) in [8]. The method consists in choosing $K$ symmetrized autocovariance matrices with different lags $\tau_1, \ldots, \tau_K$, which are then jointly diagonalized. This is again more flexible, and different autocovariance matrices can contribute to the separation of different lags. The statistical properties of AMUSE are discussed in [61] and those of SOBI in [36,59,62]. Note, however, that SOBI is not always better than AMUSE but it is in most cases. The choice of lags is, however, an open question that has a large impact in practice, where the default is often to use simply the first 12 lags. For more sophisticated considerations for lag selection, see, for example, [98,99]. Replacing Cov and ACov's in AMUSE or SOBI with robust alternatives is discussed, for example, in [38] but requires that the time series be symmetric.
2. Independent component time series (IC time series) model. In this model, one makes the assumptions:

    **(IC time series 1):** $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$
    **(IC time series 2):** The latent time series contained in $\boldsymbol{Z}$ are all independent.

Note that (IC time series 1) implies stationarity which is assumed in the methods presented in the following. It could, however, be relaxed and there is some overlap with the BSS approach considering nonstationary data described subsequently.

The main difference when assuming the IC time series model compared with the SOS model is that independence between components is required, and one usually has components with stochastic volatility, such as GARCH components. Statistics measuring second order information do not necessarily carry information, and higher order information is therefore commonly used. The main method, in our context, is the generalized FOBI (gFOBI) [57] which extends FOBI by defining the scatter functional for the mean zero process

$$\text{Cov}_{4,\tau}(\boldsymbol{X}) = E(\boldsymbol{X}_{t+\tau}\boldsymbol{X}_t^\top \text{Cov}(\boldsymbol{X}_t)^{-1}\boldsymbol{X}_t\boldsymbol{X}_{t+\tau}^\top),$$

which therefore can be seen as a lagged fourth moment matrix. For gFOBI, one selects a set of lags $\tau_1, \ldots, \tau_K$ used in the joint diagonalization approach for $\boldsymbol{S}_j = \text{Cov}_{4,\tau_j}, j \in \{1, \ldots, j\}$. If the set consists only of $\tau_1 = 0$ the method reduces to FOBI. gFOBI therefore can solve the BSS problem if all fourth moments are finite and the set of lags contains a lag $\tau$ for which the $i$th and $j$th diagonal elements of $\text{Cov}_{4,\tau}(\boldsymbol{X})$ are distinct, for all pairs $i \neq j$. Note that JADE was similarly extended for this setting to gJADE in [57].

3. Nonstationary source separation (NSS). Thus far, stationary data are considered. In the NSS models this is relaxed slightly and the assumptions are:

(**NSS 1**): $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\text{Cov}(\boldsymbol{Z}_t) = \boldsymbol{D}_t$, where $\boldsymbol{D}_t$ is a diagonal matrix depending on $t$.
(**NSS 2**): $\text{ACov}_\tau(\boldsymbol{Z}_t) = \boldsymbol{D}_{\tau,t}$ for all $\tau \geq 1$ where $\boldsymbol{D}_{\tau,t}$ are all diagonal matrices.

In this model, the mean is stationary, but not the second moment. To make the model formulation easier, it is often assumed that, for the observed time series $\boldsymbol{X}_T$, the latent components are scaled such that $\text{Cov}(\boldsymbol{Z}_T) = \boldsymbol{I}_p$. The main idea for NSS is to divide the observed time span $T$ into $K$ non-overlapping intervals $T_1, \ldots, T_K$, and then compute the scatter matrices separately for each interval, and jointly diagonalize them.

The NSS method.SD [19] uses $K = 2$ and the unmixing matrix $W$ corresponds to the matrix that simultaneously diagonalizes $\text{Cov}(\boldsymbol{X}_{T_1})$ and $\text{Cov}(\boldsymbol{X}_{T_2})$. Similar to 2-scatter ICA and AMUSE, the performance of this approach is sensitive to the division, and requires that each component has a distinct variance in the intervals. A straightforward extension is NSS.JD [19] that chooses $K > 2$ and jointly diagonalizes $\text{Cov}(\boldsymbol{X}_{T_1}^{st}), \ldots, \text{Cov}(\boldsymbol{X}_{T_K}^{st})$, where $\boldsymbol{X}^{st} = \text{Cov}(\boldsymbol{X}_T)^{-1/2}(\boldsymbol{X}_T - \boldsymbol{1}_T\bar{\boldsymbol{X}}_T^\top)$. NSS.SD and NSS.JD only require a temporal ordering of the observations but do not otherwise exploit the serial dependence. The approach NSS.TD.JD [18] therefore chooses $K$ intervals and a set of $L$ lags $\tau_1, \ldots, \tau_L$ and jointly diagonalizes the $K \times L$ autocovariance matrices $\text{ACov}_{\tau_i}^S(\boldsymbol{X}_{T_j}^{st}), i \in \{1, \ldots, L\} j \in \{1, \ldots, K\}$. Thus, for $K = 1$, this approach reduces to SOBI, and the general idea is that the data follow a block stationary model. NSS was first considered in the context of audio signals, and $K$ was chosen such that there are sufficient observations within an interval, so that the scatter matrices can be computed with sufficient precision. Another framework is that, on $K$ subjects, the same experiment was performed and produces for each subject a $p$-variate time series. Then, assuming that for all subjects the same "mixing" occurred, one can concatenate the $K$ time series and apply an NSS approach, where the intervals correspond to the concatenation points. Such an approach is often referred to as groupICA. NSS with robust scatter functionals, was, for example considered in [70].

As it is not always clear which of the three time series BSS models is suitable, there exist generalizations combining different approaches. In general, approaches such as gSOBI [60], cannot be expressed in a joint diagonalization framework. Nordhausen et al. [71] used almost all scatter matrices, as described above, including the subdivision into intervals, for joint diagonalization to cover all three models. This is, however, very challenging as the different scatter functionals are of different magnitudes, and it is not very clear how to weight them. This is still an area for further research. More details about general BSS approaches for time series are reviewed in [22,89].

### 5.3. BSS for spatial data

Most areas where BSS was applied to date produced time series data. Therefore, the focus of BSS was mainly on time series methods. However, recently, BSS was also considered in the context of spatial data. In that case, $\boldsymbol{X} = \boldsymbol{X}(\boldsymbol{s})$ is a $p$-variate random field specified on the domain $\mathcal{S}$, where the domain can be 1, 2 or 3 dimensional. To estimate the latent components, one has a sample of $n$ points $\boldsymbol{X}(\boldsymbol{s}_i)$, sampled at the distinct locations $\boldsymbol{s}_i \in \mathcal{S}, i \in \{1, \ldots, n\}$. $\boldsymbol{X}_S$ denotes then the data matrix with the sampled observations.

Two spatial settings that have been considered so far in a BSS framework, can be seen as spatial counterpart to the SOS and NSS time series models where the role of the autocovariance matrices will be taken on by local covariance matrices (see the definition in Section 2).

1. Spatial blind source separation model (SBSS) makes the following model assumptions:

(**SBSS1**): $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\text{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$.
(**SBSS2**): $E(\boldsymbol{Z}(\boldsymbol{s}), \boldsymbol{Z}(\boldsymbol{s}')^\top) = \boldsymbol{D}_h$, where $h = \|\boldsymbol{s} - \boldsymbol{s}'\|$ for all $\boldsymbol{s}$ and $\boldsymbol{s}' \in \mathcal{S}$ with $\boldsymbol{s} \neq \boldsymbol{s}'$ and the diagonal matrix $\boldsymbol{D}_h$ contains the univariate covariance functions corresponding to the latent fields.
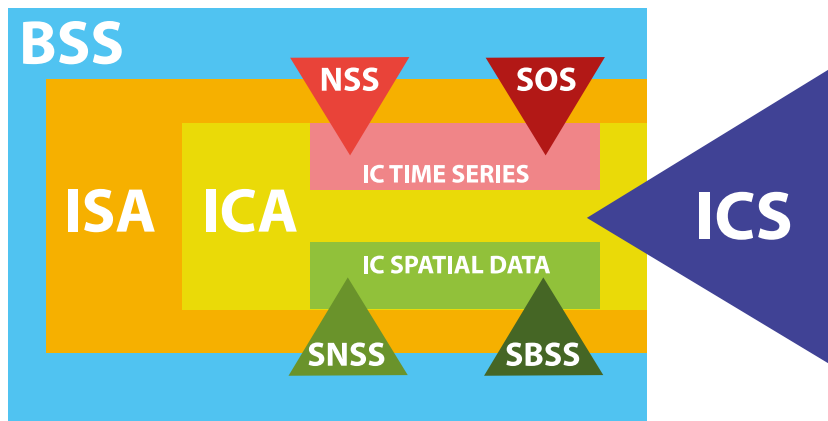
**Fig. 2.** Schematic overview of the different BSS models and ICS. For the definition of the models see Section 5.

Thus, in SBSS, the latent fields are assumed to be uncorrelated/independent stationary random fields.

Nordhausen et al. [76] suggested a 2-scatter approach that jointly diagonalizes Cov and one LCov$_f$ for SBSS. The performance of this approach depends again heavily on the chosen kernel $f$. Bachoc et al. [7] then suggested a joint diagonalization approach with Cov and $\boldsymbol{S}_1 = \mathrm{LCov}_{f_1}, \ldots, \boldsymbol{S}_K = \mathrm{LCov}_{f_K}$, for $K \geq 2$, where the so-called ring kernels, with different radii, are the most natural kernels considered so far. The statistical properties of the two approaches are given in [7] in the case of latent Gaussian random fields and show again that the joint diagonalization approach seems preferable.

2. Spatial nonstationary source separation (SNSS). This model assumes

(**SNSS 1**): $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}(\boldsymbol{s})) = \boldsymbol{D}_{\boldsymbol{s}}$, where $\boldsymbol{D}_{\boldsymbol{s}}$ is a diagonal matrix depending on $\boldsymbol{s} \in \mathcal{S}$.

(**SNSS 2**): $E(\boldsymbol{Z}(\boldsymbol{s}), \boldsymbol{Z}(\boldsymbol{s}')^\top) = \boldsymbol{D}_{\boldsymbol{s},\boldsymbol{s}'}$ or all $\boldsymbol{s}$ and $\boldsymbol{s}' \in \mathcal{S}$ with $\boldsymbol{s} \neq \boldsymbol{s}'$ and the diagonal matrix $\boldsymbol{D}_{\boldsymbol{s},\boldsymbol{s}'}$ depends on the locations.

Thus, the location is stationary and all latent fields are uncorrelated or independent. However, for all latent fields, the spatial covariance is non-stationary.

The idea for algebraic BSS in this model is quite similar to NSS in the time series case. The domain is divided into $K$ non-overlapping subdomains $\mathcal{S}_1, \ldots, \mathcal{S}_K$ such that $\mathcal{S}_1 \cup \cdots \cup \mathcal{S}_K = \mathcal{S}$. Then, analogously to the time series setting, [65] suggested the methods SNSS.SD, SNSS.JD and SNSS.SJD.

For SNSS.SD, $K = 2$ and the covariance matrices of the two domains are simultaneously diagonalized, which is again very sensitive to division into subdomains. SNSS.JD, which whitens the data using $\mathrm{Cov}(\boldsymbol{X}_S)$ and then jointly diagonalizes the $K > 2$ covariance matrices obtained for the subdomains, i.e., $\mathrm{Cov}(\boldsymbol{X}_{S_1}^{st}), \ldots, \mathrm{Cov}(\boldsymbol{X}_{S_K}^{st})$, is less sensitive to division into subdomains. Both, SNSS.SD and SNSS.JD, are based only on the spatial ordering of the points. If it is assumed that there would be some kind of block-stationary model underlying, SNSS.SJD suggests to compute $L$ local covariance matrices $\mathrm{LCov}_{f_j}(\boldsymbol{X}_{S_i}^{st})$, $i \in \{1, \ldots, K\}, j \in \{1, \ldots, L\}$, for all subdomains, and then jointly diagonalizes these $K \times L$ matrices.

All the above BSS methods, either simultaneously diagonalize two scatter functionals or jointly diagonalize $K+1$ scatter matrices with one scatter playing a special role. There exist many other BSS models or methods where joint diagonalization plays a role. Some are for example summarized in [17,100]. BSS is still an active research area, and spatial BSS is currently actively developed. A schematic overview of the BSS models covered in the present review and of ICS is given in Fig. 2.

Note that, as mentioned earlier, the goal of BSS is to estimate the latent components, and for that purpose, we used simultaneous and joint diagonalization (see Section 3) to obtain an unmixing matrix $\boldsymbol{W}$ such that

$$\boldsymbol{z}_i = \boldsymbol{W}(\boldsymbol{x}_i - \boldsymbol{T}(\boldsymbol{X}_n)), \quad i \in \{i, \ldots, n\}.$$

This is however rather imprecise, as none of the models described above are identifiable in a strict sense. In all the BSS models described above, one can write

$$\boldsymbol{X} = \boldsymbol{A}\boldsymbol{Z} = (\boldsymbol{A}\boldsymbol{J}\boldsymbol{P})(\boldsymbol{P}^\top\boldsymbol{J}\boldsymbol{Z}) = \boldsymbol{A}^*\boldsymbol{Z}^*,$$

where $\boldsymbol{J}$ is $p \times p$ sign-change matrix and $\boldsymbol{P}$ a $p \times p$ permutation matrix. Thus, the signs and order of the components cannot be fixed. Consequently, for any unmixing matrix $\boldsymbol{W}$, the matrix $\boldsymbol{J}\boldsymbol{P}\boldsymbol{W}$ is also an unmixing matrix for all permutation matrices $\boldsymbol{P}$ and all sign-change matrices $\boldsymbol{J}$. However, these identifiability issues are usually not considered to be a problem and the order of the components is, for example, usually fixed based on the diagonal elements of $\boldsymbol{D}$ in the case of simultaneous diagonalization, and based on $\mathrm{diag}(\sum_{i=1}^{K} \boldsymbol{W}\boldsymbol{S}_i(\boldsymbol{X}_n)\boldsymbol{W}^\top)$ in the case of joint diagonalization.

In performance studies, these indeterminacies have naturally to be taken into account and an overview of BSS performance measures is given for example in [84].

Most algebraic BSS methods described above are implemented in R via the R packages ICS, BSSasymp, JADE [63] , tsBSS [73] and SpatialBSS [68] where JADE contains also some performance measures.

## 6. Joint diagonalization in the context of supervised multivariate methods

PCA, ICS and BSS are often used as dimension reduction methods, which means that some components are selected and used in further modelling. However, if there is a response $Y$ to be modelled, no direct information is used when computing the new directions in an unsupervised manner. Such dimension reduction methods are therefore called unsupervised dimension reduction methods. When information about the target is used in the dimension reduction process, one refers to it as supervised dimension reduction (SDR). Surprisingly, many SDR methods can be seen within a joint diagonalization framework.

For example, linear discriminant analysis (LDA) [30] can be seen as a supervised method in a classification context. Let us consider a dataset $\boldsymbol{X} = (\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n)$, with $n$ observations and $p$ variables, which is partitioned into $K$ subpopulations or groups. Fisher's idea was to look for the best linear function of the $p$ variables which maximized the ratio of the between-groups covariance to the within-groups covariance. The intuition is that groups are more easily visible when the between-groups variability is large in comparison with the within-groups variability. We use the index $i$ for the group and $j$ for the observation in each group, so that $\boldsymbol{x}_{ij}$ denotes the $j$th observation in group $i$, for $i \in \{1, \ldots, K\}, j \in \{1, \ldots, n_i\}$, where $n_i$ denotes the number of observations in group $i$. Using the analysis of variance equation, we can decompose the total covariance matrix Cov, which does not take into account the groups, into the between scatter matrix $\mathrm{Cov}_{\boldsymbol{B}}$ and the within scatter matrix $\mathrm{Cov}_{\boldsymbol{W}}$ defined by:

$$\mathrm{Cov}_{\boldsymbol{B}} = \frac{1}{n} \sum_{i=1}^{K} n_i (\bar{\boldsymbol{x}}_i - \bar{\boldsymbol{x}})(\bar{\boldsymbol{x}}_i - \bar{\boldsymbol{x}})^\top \tag{6}$$

$$\mathrm{Cov}_{\boldsymbol{W}} = \frac{1}{n} \sum_{i=1}^{K} \sum_{j=1}^{n_i} (\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)(\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)^\top \tag{7}$$

where $\bar{\boldsymbol{x}}_i$ denotes the mean of the $i$th group and $\bar{\boldsymbol{x}}$ the overall mean. The Fisher's linear discriminant vectors are the eigenvectors of $\mathrm{Cov}_{\boldsymbol{W}}^{-1}\mathrm{Cov}_{\boldsymbol{B}}$ (see [51] for more details). Both $\mathrm{Cov}_{\boldsymbol{B}}$ and $\mathrm{Cov}_{\boldsymbol{W}}$ are scatter matrices in the sense that they are affine equivariant and semipositive definite. Note however that the between-matrix is of rank $K-1$ and thus is generally singular. As a consequence, in general, there are $K-1$ nontrivial linear discriminant vectors. As detailed in Section 3, the Fisher's discriminant vectors are obtained as the solution of a generalized eigendecomposition.

Also canonical correlation analysis (CCA) [33] can be formulated as the simultaneous diagonalization of two scatter functionals when using

$$\boldsymbol{S}_1(\boldsymbol{X}) = \mathrm{Cov}(\boldsymbol{X}) \quad \text{and} \quad \boldsymbol{S}_2(\boldsymbol{X}) = \mathrm{Cov}_{CCA} = \mathrm{Cov}(\boldsymbol{X}, \boldsymbol{Y})\mathrm{Cov}(\boldsymbol{Y})^{-1/2}\mathrm{Cov}(\boldsymbol{Y}, \boldsymbol{X}).$$

Most SDR methods are developed in a regression context where, for simplicity, we assume from now on that the response $Y$ is univariate. In the spirit of a BSS model, we assume that:

$$\boldsymbol{X} = \boldsymbol{A}\boldsymbol{Z} + \boldsymbol{\mu},$$

where $\boldsymbol{A}$ is the $p \times p$ full rank mixing matrix and $\boldsymbol{\mu}$ the $p$-variate location vector. For the latent $p$-vector, we assume there exists a partition $\boldsymbol{Z} = \left(\boldsymbol{Z}^{(1)\top}, \boldsymbol{Z}^{(2)\top}\right)^\top$ with respective dimensions $k$ and $p-k$. The assumptions are:

**(SDR 1):** $E(\boldsymbol{Z}) = \boldsymbol{0}$ and $\mathrm{Cov}(\boldsymbol{Z}) = \boldsymbol{I}_p$.

**(SDR 2):** $(Y, \boldsymbol{Z}^{(1)\top})^\top$ and $\boldsymbol{Z}^{(2)}$ are independent.

Thus all information on $Y$ is contained in $\boldsymbol{Z}^{(1)}$. Note that there are naturally many partitions of $\boldsymbol{Z}$ fulfilling these assumptions. But the partition of interest is the one with the smallest value $k$ that needs to be estimated together with the unmixing matrix $\boldsymbol{W}$. Note also that $\boldsymbol{Z}^{(1)}$ is only identifiable up to an orthogonal transformation, which means that the "unmixing" matrix can only recover the subspace of interest, which is however sufficient.

Liski et al. [46] defined supervised invariant coordinate selection (SICS) as the joint diagonalization of one unsupervised scatter functional ($\boldsymbol{S}_1$) and one supervised scatter functional ($\boldsymbol{S}_2$). Many well established supervised dimension reduction methods, like sliced inverse regression (SIR) [42], sliced average variance estimation (SAVE) [23], principal Hessian directions (pHd) [43], or directional regression (DR) [45], can be seen as special cases of SICS. All these methods use $\boldsymbol{S}_1 = \mathrm{Cov}$ and differ regarding $\boldsymbol{S}_2$. SIR, for example, uses $\boldsymbol{S}_{SIR}$ as defined in Section 2. For the exact forms of $\boldsymbol{S}_{SAVE}$, $\boldsymbol{S}_{pHd}$ and $\boldsymbol{S}_{DR}$ we refer to [46], where many other possibilities for supervised scatter functionals are listed. The advantage of an SDR approach over unsupervised methods is demonstrated in Fig. 3, where there is a response $Y$ which is to be explained by four possible predictors. Panel A gives the original data where no clear relationship between any of the predictors and $Y$ is visible. The PCs based on $\boldsymbol{X}$ given in panel B are not more informative regarding their relationship with $Y$. The invariant

**Fig. 3.** Comparison of different data transformations. Panel A shows a matrix scatterplot with density estimators on the diagonal and correlations for the original data and where $Y$ is the response to be modelled by the 4 predictors, Panel B the principal components based on $\boldsymbol{X}_n$, Panel C the invariant coordinates (FOBI) based on $\boldsymbol{X}_n$ and Panel D the supervised invariant coordinates based on $\boldsymbol{X}_n$ and $Y$. Clearly the supervised components make it easiest to see a relationship between response and predictors.

coordinates in panel C give some idea about the relationship when looking at IC.3. But the relationship is very clearly visible in panel D where SICS is displayed.

The performance of the SDR methods depends a lot on the true relationship between response and predictors, and different methods are more suitable to recognize certain types of dependencies than others. For detailed discussions about SDR methods, we refer, for example, to [44,50]. Note that many of these methods also can make weaker assumptions than (SDR2).

In the example in Fig. 3, the true $k$ is 1. In practice this needs to be estimated. In SDR, depending on the scatter used, the theoretical value of eigenvalues which correspond to $\boldsymbol{Z}^{(2)}$, i.e., the values in $\boldsymbol{D}$, are known, and therefore tests and estimators can be based on these eigenvalues, as for example discussed in [10,11,48,49,81].

Supervised dimension reduction methods are of course also of interest in the context of time series and spatial data. The effect of the predictors on the response might, for example, be delayed in the time series case, or depend on neighbouring values in the spatial setting. However, it is straightforward to adjust the SDR assumptions from above, and to formulate an appropriate BSS-SDR framework for dependent data. To take the temporal delay and spatial proximity into account,

**Table 1**

Multivariate methods which are based on the joint diagonalization of two or more scatter matrices.

| Name | Family | Primary data type | Scatters used |
|---|---|---|---|
| ICS [102] | ICS | non-elliptical iid data | two different scatter matrices |
| PAA [25] | ICS | non-elliptical iid data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathrm{Cov}_{-1}$ |
| LDA [30] | SDR | multigroup iid data | $\mathbf{S}_1 = \mathrm{Cov}_{\mathbf{W}}$, $\mathbf{S}_2 = \mathrm{Cov}_{\mathbf{B}}$ |
| CCA [33] | SDR | two group data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathrm{Cov}_{CCA}$. |
| FOBI [12,86] | ICS, ICA | non-elliptical/ICA iid data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathrm{Cov}_4$ |
| 2-Scatter-ICA [88] | ICA | ICA iid data | Two scatters with independence property |
| k-Scatter-ICA [72] | ICA | ICA iid data | k scatters with independence property |
| 3-scatter-ISA [74] | ISA | ISA iid data | Three different scatter matrices with (block) independence property |
| SICS [46] | SDR | regression data | an unsupervised and a supervised scatter matrix |
| SIR [42] | SDR | regression data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathbf{S}_{SIR}$ |
| SAVE [23] | SDR | regression data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathbf{S}_{SAVE}$ |
| pHd [43] | SDR | regression data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathbf{S}_{pHd}$ |
| DR [45] | SDR | regression data | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathbf{S}_{DR}$ |
| TSIR [55] | SDR | time series regression data | $\mathbf{S}_1 = \mathrm{Cov}$, and $K$ $\mathbf{S}_{TSIR,\tau}$ s |
| TSAVE [56] | SDR | time series regression data | $\mathbf{S}_1 = \mathrm{Cov}$, and $K$ $\mathbf{S}_{TSAVE,\tau}$ s |
| SSIR [67] | SDR | spatial regression data | $\mathbf{S}_1 = \mathrm{Cov}$, and $K$ $\mathbf{S}_{SIR,\tau}$ s |
| AMUSE [61,101] | SOS | stationary time series | $\mathbf{S}_1 = \mathrm{Cov}$, $\mathbf{S}_2 = \mathrm{ACov}_\tau^S$ |
| SOBI [8,59,64] | SOS | stationary time series | $\mathbf{S}_1 = \mathrm{Cov}$, $K \geq 2$ $\mathrm{ACov}_\tau^S$ s |
| gFOBI [57] | IC-time series | time series with for example stochastic volatility | $\mathbf{S}_1 = \mathrm{Cov}$, $K$ lagged 4th moment matrices |
| NSS.SD [19] | NSS | non-stationary time series | 2 Cov s |
| NSS.JD [19] | NSS | non-stationary time series | K + 1 Cov s |
| NSS.TD.JD [18] | NSS | block stationary time series | $K \times L$ $\mathrm{ACov}_\tau^S$ s |
| SBSS [7,76] | SBSS | stationary spatial data | Cov and K $\mathrm{LCov}_f$ s |
| SNSS.SD [65] | SNSS | non-stationary spatial data | 2 Cov s |
| SNSS.JD [65] | SNSS | non-stationary spatial data | K + 1 Cov s |
| SNSS.TD.JD [65] | SNSS | block stationary spatial data | $K \times L$ $\mathrm{LCov}_f$ s |

supervised temporal and spatial scatter functionals should be used, and more than two scatter matrices might be used. Matilainen et al. [55] define time series SIR (TSIR), which is based on $\mathbf{S}_{TSIR,\tau}(\mathbf{X}) = \mathrm{Cov}(E(\mathbf{X}_t|Y_{t+\tau}))$, where $\tau$ is some lag. Then, TSIR whitens the data using Cov and jointly diagonalizes $\mathbf{S}_{TSIR,\tau_i}(\mathbf{X}^{st})$ with $\tau_i \in \{\tau_1, \ldots, \tau_K\}$. Time series SAVE (TSAVE) is suggested in [56], and jointly diagonalizes Cov, and $K$ so-called time series SAVE matrices $\mathbf{S}_{TSAVE,\tau_i}$, using $K$ different lags. Spatial SIR (SSIR) was so far only considered for lattice data in [67], and jointly diagonalizes Cov and $\mathbf{S}_{SSIR,\tau}(\mathbf{X}) = \mathrm{Cov}(E(\mathbf{X}_\mathbf{s}|Y_{\mathbf{s}+\tau}))$, where $\tau$ is now a $d$-dimensional lag. These approaches are all fairly new and inference tools are still missing.

Various SDR approaches discussed here are, for example, implemented in R in the packages dr [106], ICS and tsBSS.

## 7. Conclusions

Many multivariate statistical methods make use of the joint diagonalization of several scatter matrices as illustrated in the previous sections. Table 1 summarizes the different models, methods and scatter functionals that are jointly diagonalized. However, the overview we propose in this paper is far from exhaustive. For example, Chabriel et al. [17], Theis and Inouye [100] give an overview of algebraic BSS methods that contains models and approaches not mentioned here. Then, there are also completely different multivariate statistical methods that use joint diagonalization, but that were not considered here, such as for example common principal component analysis [31]. Additionally, the methods have been extended in several directions to tackle more complex data, such as tensors, functional data or composition data [see for example 66,104,105], and similarities and different approaches are discussed in [24,29,81]. Finally, let us mention the problem of high-dimensional data and the sparsity question that needs further development.

## CRediT authorship contribution statement

**Klaus Nordhausen:** Conceptualization, Methdodology, Visualization, Writing. **Anne Ruiz-Gazen:** Conceptualization, Methodology, Writing.

**Declaration of competing interest**

**References**

[1] T. Adali, M. Anderson, G.-S. Fu, Diversity in independent component and vector analyses: Identifiability, algorithms, and applications in medical imaging, IEEE Signal Process. Mag. 31 (2014) 18–33.

[2] F. Alashwali, J.T. Kent, The use of a common location measure in the invariant coordinate selection and projection pursuit, J. Multivariate Anal. 152 (2016) 145–161.

[3] T. Anderson, An Introduction To Multivariate Statistical Analysis, third ed., Wiley, New York, 2003.

[4] A. Archimbaud, J. May, K. Nordhausen, A. Ruiz-Gazen, ICSShiny: ICS via a shiny application, 2018, R package version 0.5.

[5] A. Archimbaud, K. Nordhausen, A. Ruiz-Gazen, ICS for multivariate outlier detection with application to quality control, Comput. Statist. Data Anal. 128 (2018) 184–199.

[6] A. Archimbaud, K. Nordhausen, A. Ruiz-Gazen, Unsupervi019zed outlier detection with ICSOutlier, R Journal 10 (1) (2018) 234–250.

[7] F. Bachoc, M.G. Genton, K. Nordhausen, A. Ruiz-Gazen, J. Virta, Spatial blind source separation, Biometrika 107 (2020) 627–646.

[8] A. Belouchrani, K. Abed Meraim, J.-F. Cardoso, E. Moulines, A blind source separation technique based on second order statistics, IEEE Trans. Signal Process. 45 (1997) 434–444.

[9] M. Bilodeau, D. Brenner, Theory of Multivariate Statistics, Springer, New York, 2008.

[10] E. Bura, R. Cook, Extending sliced inverse regression: The weighted chi-squared test, J. Amer. Statist. Assoc. 96 (2001) 996–1003.

[11] E. Bura, J. Yang, Dimension estimation in sufficient dimension reduction: A unifying approach, J. Multivariate Anal. 102 (1) (2011) 130–142.

[12] J.-F. Cardoso, Source separation using higher order moments, in: Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing, IEEE, 1989, pp. 2109–2112.

[13] J.-F. Cardoso, A. Souloumiac, Jacobi angles for simultaneous diagonalization, SIAM J. Matrix Anal. Appl. 17 (1996) 161–164.

[14] H. Caussinus, M. Fekri, S. Hakam, A. Ruiz-Gazen, A monitoring display of multivariate outliers, Comput. Statist. Data Anal. 44 (1) (2003) 237–252.

[15] H. Caussinus, A. Ruiz, Interesting projections of multidimensional data by means of generalized principal component analyses, in: K. Momirović, V. Mildner (Eds.), Compstat, Physica-Verlag HD, Heidelberg, 1990, pp. 121–126.

[16] H. Caussinus, A. Ruiz-Gazen, Classification and generalized principal component analysis, in: P. Brito, G. Cucumel, P. Bertrand, F. de Carvalho (Eds.), Selected Contributions in Data Analysis and Classification, Springer, Berlin, 2007, pp. 539–548.

[17] G. Chabriel, M. Kleinsteuber, E. Moreau, H. Shen, P. Tichavsky, A. Yeredor, Joint matrices decompositions and blind source separation: A survey of methods, identification, and applications, IEEE Signal Process. Mag. 31 (3) (2014) 34–43.

[18] S. Choi, A. Cichocki, Blind separation of nonstationary and temporally correlated sources from noisy mixtures, in: Neural Networks for Signal Processing X. Proceedings of the 2000 IEEE Signal Processing Society Workshop, Vol. 1, IEEE, 2000, pp. 405–414.

[19] S. Choi, A. Cichocki, Blind separation of nonstationary sources in noisy mixtures, Electron. Lett. 36 (2000) 848–849.

[20] A. Cichocki, S.-I. Amari, Adaptive Blind Signal and Image Processing: Learning Algorithms and Applications, John Wiley & Sons, New York, 2002.

[21] D. Clarkson, A least squares version of algorithm AS 211: The F-G diagonalization algorithm, Appl. Stat. 37 (1988) 317–321.

[22] P. Comon, C. Jutten, Handbook of Blind Source Separation: Independent Component Analysis and Applications, Academic Press, Oxford, 2010.

[23] R. Cook, SAVE: A method for dimension reduction and graphics in regression, Comm. Statist. Theory Methods 29 (2000) 2109–2121.

[24] R.D. Cook, A slice of multivariate dimension reduction, J. Multivariate Anal. (2021) 104812, (online first).

[25] F. Critchley, A. Pires, C. Amado, Principal Axis Analysis, Technical Report, (06/14) The Open University Milton Keynes, 2006.

[26] C. Croux, G. Haesbroeck, Principal component analysis based on robust estimators of the covariance or correlation matrix: influence functions and efficiencies, Biometrika 87 (3) (2000) 603–618.

[27] M. Fekri, A. Ruiz-Gazen, A B-robust non-iterative scatter matrix estimator: Asymptotics and application to cluster detection using invariant coordinate selection, in: K. Nordhausen, S. Taskinen (Eds.), Modern Nonparametric, Robust and Multivariate Methods: Festschrift in Honour of Hannu Oja, Springer International Publishing, Cham, 2015, pp. 395–423.

[28] D. Fischer, M. Honkatukia, M. Tuiskula-Haavisto, K. Nordhausen, D. Cavero, R. Preisinger, J. Vilkki, Subgroup detection in genotype data using invariant coordinate selection, BMC Bioinformatics 18 (2017) 173–181.

[29] D. Fischer, K. Nordhausen, H. Oja, On linear dimension reduction based on diagonalization of scatter matrices for bioinformatics downstream analyses, Heliyon 6 (2020) e05732.

[30] R. Fisher, The use of multiple measurements in taxonomic problems, Ann. Eugen. 7 (2) (1936) 179–188.

[31] B. Flury, Common Principal Components & Related Multivariate Models, John Wiley & Sons, Chichester, 1988.

[32] B.N. Flury, W. Gautschi, An algorithm for simultaneous orthogonal transformation of several positive definite symmetric matrices to nearly diagonal form, SIAM J. Sci. Stat. Comput. 7 (1) (1986) 169–184.

[33] H. Hotelling, Relations between two sets of variates, Biometrika 28 (1936) 321–377.

[34] P. Huber, Projection pursuit, Ann. Statist. 13 (1985) 435–475.

[35] P. Huber, E. Ronchetti, Robust Statistics, Wiley, Hoboken, 2011.

[36] K. Illner, J. Miettinen, C. Fuchs, S. Taskinen, K. Nordhausen, H. Oja, F. Theis, Model selection using limiting distributions of second-order blind source separation algorithms, Signal Process. 113 (2015) 95–103.

[37] P. Ilmonen, J. Nevalainen, H. Oja, Characteristics of multivariate distributions and the invariant coordinate system, Statist. Probab. Lett. 80 (23) (2010) 1844–1853.

[38] P. Ilmonen, K. Nordhausen, H. Oja, F. Theis, An affine equivariant robust second-order BSS method, in: E. Vincent, A. Yeredor, Z. Koldovský, P. Tichavský (Eds.), Latent Variable Analysis and Signal Separation. LVA/ICA 2015. Lecture Notes in Computer Science, Vol. 9237, Springer, Cham, 2015, pp. 328–335.

[39] P. Ilmonen, H. Oja, R. Serfling, On invariant coordinate system (ICS) functionals, Internat. Statist. Rev. 80 (2012) 93–110.

[40] I. Jolliffe, Principal Component Analysis, second ed., Springer, New York, 2002.

[41] A. Kankainen, S. Taskinen, H. Oja, Tests of multinormality based on location vectors and scatter matrices, Stat. Methods Appl. 16 (2007) 357–379.

[42] K.-C. Li, Sliced inverse regression for dimension reduction, J. Amer. Statist. Assoc. 86 (414) (1991) 316–327.

[43] K.-C. Li, On principal hessian directions for data visualization and dimension reduction: Another application of Stein's lemma, J. Amer. Statist. Assoc. 87 (420) (1992) 1025–1039.

[44] B. Li, Sufficient Dimension Reduction Methods and Applications with R, Chapman and Hall/CRC, Boca Raton, 2018.

[45] B. Li, S. Wang, On directional regression for dimension reduction, J. Amer. Statist. Assoc. 102 (479) (2007) 997–1008.

[46] E. Liski, K. Nordhausen, H. Oja, Supervised invariant coordinate selection, Statistics 4 (2014) 711–731.

[47] N. Loperfido, Some theoretical properties of two kurtosis matrices, with application to invariant coordinate selection, J. Multivariate Anal. (2021) 104809.

[48] W. Luo, B. Li, Combining eigenvalues and variation of eigenvectors for order determination, Biometrika 103 (4) (2016) 875–887.

[49] W. Luo, B. Li, On order determination by predictor augmentation, Biometrika 108 (2021) 557–574.

[50] Y. Ma, L. Zhu, A review on dimension reduction, Internat. Statist. Rev. 81 (1) (2013) 134–150.

[51] K. Mardia, J. Kent, J. Bibby, Multivariate Analysis, Academic Press, London, 1979.

[52] R.A. Maronna, Robust M-estimators of multivariate location and scatter, Ann. Statist. (1976) 51–67.

[53] R.A. Maronna, R.D. Martin, V.J. Yohai, M. Salibián-Barrera, Robust statistics: Theory and methods (with R), John Wiley & Sons, New York, 2019.

[54] R.A. Maronna, V.J. Yohai, Robust estimation of multivariate location and scatter, in: N. Balakrishnan, T. Colton, B. Everitt, W. Piegorsch, F. Ruggeri, J. Teugels (Eds.), Wiley StatsRef: Statistics Reference Online, Wiley, 2016, pp. 1–12.

[55] M. Matilainen, C. Croux, K. Nordhausen, H. Oja, Supervised dimension reduction for multivariate time series, Econometr. Stat. 4 (2017) 57–69.

[56] M. Matilainen, C. Croux, K. Nordhausen, H. Oja, Sliced average variance estimation for multivariate time series, Statistics 53 (2019) 630–655.

[57] M. Matilainen, K. Nordhausen, H. Oja, New independent component analysis tools for time series, Statist. Probab. Lett. 105 (2015) 80–87.

[58] J. Miettinen, Alternative diagonality criteria for SOBI, in: K. Nordhausen, S. Taskinen (Eds.), Modern Nonparametric, Robust and Multivariate Methods: Festschrift in Honour of Hannu Oja, Springer, Cham, 2015, pp. 455–469.

[59] J. Miettinen, K. Illner, K. Nordhausen, H. Oja, S. Taskinen, F. Theis, Separation of uncorrelated stationary time series using autocovariance matrices, J. Time Series Anal. 37 (2016) 337–354.

[60] J. Miettinen, M. Matilainen, K. Nordhausen, S. Taskinen, Extracting conditionally heteroskedastic components using independent component analysis, J. Time Series Anal. 41 (2020) 293–311.

[61] J. Miettinen, K. Nordhausen, H. Oja, S. Taskinen, Statistical properties of a blind source separation estimator for stationary time series, Statist. Probab. Lett. 82 (11) (2012) 1865–1873.

[62] J. Miettinen, K. Nordhausen, H. Oja, S. Taskinen, Deflation-based separation of uncorrelated stationary time series, J. Multivariate Anal. 123 (2014) 214–227.

[63] J. Miettinen, K. Nordhausen, S. Taskinen, Blind source separation based on joint diagonalization in R: The packages JADE and BSSasymp, J. Stat. Softw. 76 (2017) 1–31.

[64] J. Miettinen, S. Taskinen, K. Nordhausen, H. Oja, Fourth moments and independent component analysis, Statist. Sci. 30 (3) (2015) 372–390.

[65] C. Muehlmann, F. Bachoc, K. Nordhausen, Spatial nonstationary source separation, 2021, https://arxiv.org/abs/2107.01916, Arxiv.

[66] C. Muehlmann, A. Fačevicová, A. Gardlo, H. Janečková, K. Nordhausen, Independent component analysis for compositional data, in: A. Daouia, A. Ruiz-Gazen (Eds.), Advances in Contemporary Statistics and Econometrics: Festschrift in Honor of Christine Thomas-Agnan, Springer, Cham, 2021, pp. 525–545.

[67] C. Muehlmann, K. Nordhausen, H. Oja, Sliced inverse regression for spatial data, in: E. Bura, B. Li (Eds.), Festschrift in Honor of R. Dennis Cook: Fifty Years of Contribution To Statistical Science, Springer, Cham, 2021, pp. 87–107.

[68] C. Muehlmann, K. Nordhausen, J. Virta, SpatialBSS: Blind source separation for multivariate spatial data, 2021, R package version 0.11-0.

[69] C. Muehlmann, K. Nordhausen, M. Yi, On cokriging, neural networks, and spatial blind source separation for multivariate spatial prediction, IEEE Geosci. Remote Sens. Lett. (2020) 1–5.

[70] K. Nordhausen, On robustifying some second order blind source separation methods for nonstationary time series, Statist. Papers 55 (1) (2014) 141–156.

[71] K. Nordhausen, G. Fischer, P. Filzmoser, Blind source separation for compositional time series, Math. Geosci. 53 (2021) 905–924.

[72] K. Nordhausen, H.W. Gutch, H. Oja, F.J. Theis, Joint diagonalization of several scatter matrices for ICA, in: F. Theis, A. Cichocki, A. Yeredor, M. Zibulevsky (Eds.), Latent Variable Analysis and Signal Separation: 10th International Conference, Springer, Berlin, 2012, pp. 172–179.

[73] K. Nordhausen, M. Matilainen, J. Miettinen, J. Virta, S. Taskinen, Dimension reduction for time series in a blind source separation context using R, J. Stat. Softw. 98 (2021) 1–30.

[74] K. Nordhausen, H. Oja, Scatter matrices with independent block property and ISA, in: 2011 19th European Signal Processing Conference, IEEE, 2011, pp. 1738–1742.

[75] K. Nordhausen, H. Oja, Independent component analysis: A statistical perspective, WIREs: Comput. Stat. 10 (2018) e1440.

[76] K. Nordhausen, H. Oja, P. Filzmoser, C. Reimann, Blind source separation for spatial compositional data, Math. Geosci. 47 (7) (2015) 753–770.

[77] K. Nordhausen, H. Oja, E. Ollila, Robust independent component analysis based on two scatter matrices, Aust. J. Stat. 37 (2008) 91–100.

[78] K. Nordhausen, H. Oja, E. Ollila, Multivariate models and the first four moments, in: D.R. Hunter, D.S.R. Richards, J.L. Rosenberger (Eds.), Nonparametric Statistics and Mixture Models, World Scientific, Hackensack, 2011, pp. 267–287.

[79] K. Nordhausen, H. Oja, D.E. Tyler, On the efficiency of invariant multivariate sign and rank test, in: E.P. Liski, J. Isotalo, J. Niemelä, S. Puntanen, G.P.H. Styan (Eds.), Festschrift for Tarmo Pukkila on His 60th Birthday, University of Tampere, Tampere, 2006, pp. 217–231.

[80] K. Nordhausen, H. Oja, D.E. Tyler, Tools for exploring multivariate data: The package ICS, J. Stat. Softw. 28 (2008) 1–31.

[81] K. Nordhausen, H. Oja, D. Tyler, Asymptotic and bootstrap tests for subspace dimension, J. Multivariate Anal. (2021) 104830, (online first).

[82] K. Nordhausen, H. Oja, D. Tyler, J. Virta, Asymptotic and bootstrap tests for the dimension of the non-Gaussian subspace, IEEE Signal Process. Lett. 24 (2017) 887–891.

[83] K. Nordhausen, H. Oja, D.E. Tyler, J. Virta, ICtest: Estimating and testing the number of interesting components in linear dimension reduction, 2021, R package version 0.3-4.

[84] K. Nordhausen, E. Ollila, H. Oja, On the performance indices of ICA and blind source separation, in: 2011 IEEE 12th International Workshop on Signal Processing Advances in Wireless Communications, 2011, pp. 486–490.

[85] K. Nordhausen, D.E. Tyler, A cautionary note on robust covariance plug-in methods, Biometrika 102 (3) (2015) 573–588.

[86] K. Nordhausen, J. Virta, An overview of properties and extensions of FOBI, Knowl.-Based Syst. 173 (2019) 113–116.

[87] H. Oja, Multivariate Nonparametric Methods with R. an Approach Based on Spatial Signs and Ranks, Springer, New York, 2010.

[88] H. Oja, S. Sirkiä, J. Eriksson, Scatter matrices and independent component analysis, Austrian J. Stat. 35 (2006) 175–189.

[89] Y. Pan, M. Matilainen, S. Taskinen, K. Nordhausen, A review of second-order blind identification methods, WIREs Comput. Stat. n/a (2021) e1550.

[90] D. Peña, F.J. Prieto, J. Viladomat, Eigenvectors of a kurtosis matrix as interesting directions to reveal cluster structure, J. Multivariate Anal. 101 (9) (2010) 1995–2007.

[91] M.L. Puri, P.K. Sen, Nonparametric Methods in Multivariate Analysis, John Wiley & Sons, New York, USA, 1971.

[92] R Core Team, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, 2021.

[93] U. Radojicic, K. Nordhausen, Non-Gaussian component analysis: Testing the dimension of the signal subspace, in: M. Maciak, M. Pesta, M. Schindler (Eds.), Analytical Methods in Statistics, AMISTAT 2019, Springer, Cham, 2020, pp. 101–123.

[94] U. Radojicic, K. Nordhausen, H. Oja, Notion of information and independent component analysis, Appl. Math. 65 (2020) 311–330.

[95] J.R. Schott, Matrix Analysis for Statistics, John Wiley & Sons, Hoboken, 2005.

[96] R. Serfling, Equivariance and invariance properties of multivariate quantile and related functions, and the role of standardisation, J. Nonparametr. Stat. 22 (2010) 915–936.

[97] R. Serfling, On invariant within equivalence coordinate system (IWECS) transformations, in: K. Nordhausen, S. Taskinen (Eds.), Modern Nonparametric, Robust and Multivariate Methods: Festschrift in Honour of Hannu Oja, Springer International Publishing, Cham, 2015, pp. 445–457.

[98] A.C. Tang, J.-Y. Liu, M.T. Sutherland, Recovery of correlated neuronal sources from EEG: The good and bad ways of using SOBI, NeuroImage 28 (2005) 507–519.

[99] S. Taskinen, J. Miettinen, K. Nordhausen, A more efficient second order blind identification method for separation of uncorrelated stationary time series, Statist. Probab. Lett. 116 (2016) 21–26.

[100] F. Theis, Y. Inouye, On the use of joint diagonalization in blind signal processing, in: IEEE International Symposium on Circuits and Systems, IEEE, 2006, pp. 3589–3593.

[101] L. Tong, V. Soon, Y. Huang, R. Liu, AMUSE: A new blind identification algorithm, in: Proceedings of IEEE International Symposium on Circuits and Systems, IEEE, 1990, pp. 1784–1787.

[102] D.E. Tyler, F. Critchley, L. Dümbgen, H. Oja, Invariant coordinate selection, J. R. Stat. Soc. Ser. B Stat. Methodol. 71 (3) (2009) 549–592.

[103] J. Virta, One-step M-estimates of scatter and the independence property, Statist. Probab. Lett. 110 (2016) 133–136.

[104] J. Virta, B. Li, K. Nordhausen, H. Oja, Independent component analysis for tensor-valued data, J. Multivariate Anal. 162 (2017) 172–192.

[105] J. Virta, B. Li, K. Nordhausen, H. Oja, Independent component analysis for multivariate functional data, J. Multivariate Anal. 176 (2020) 104568.

[106] S. Weisberg, Dimension reduction regression in R, J. Stat. Softw. 7 (1) (2002) 1–22.

[107] A. Yeredor, Non-orthogonal joint diagonalization in the least-squares sense with application in blind source separation, IEEE Trans. Signal Process. 50 (7) (2002) 1545–1553.

[108] A. Ziehe, P. Laskov, G. Nolte, K.-R. Müller, A fast algorithm for joint diagonalization with non-orthogonal transformations and its application to blind source separation, J. Mach. Learn. Res. 5 (Jul) (2004) 777–800.