

# groupICA: Independent component analysis for grouped data

**Niklas Pfister\***

Seminar for Statistics  
ETH Zürich, Switzerland

**Sebastian Weichwald\***

MPI for Intelligent Systems  
Tübingen, Germany

**Peter Bühlmann**

Seminar for Statistics  
ETH Zürich, Switzerland

**Bernhard Schölkopf**

MPI for Intelligent Systems  
Tübingen, Germany

We introduce **groupICA**, a novel independent component analysis (ICA) algorithm which decomposes linearly mixed multivariate observations into independent components that are corrupted (and rendered dependent) by hidden group-wise confounding. It extends the ordinary ICA model in a theoretically sound and explicit way to incorporate group-wise (or environment-wise) structure in data and hence provides a justified alternative to the use of ICA on data blindly pooled across groups. In addition to our theoretical framework, we explain its causal interpretation and motivation, provide an efficient estimation procedure and prove identifiability of the unmixing matrix under mild assumptions. Finally, we illustrate the performance and robustness of our method on simulated data and run experiments on publicly available EEG datasets demonstrating the applicability to real-world scenarios. We provide a scikit-learn compatible pip-installable Python package **groupICA** as well as R and Matlab implementations accompanied by a documentation and an audible example at <https://sweichwald.de/groupICA>.

## 1. Introduction

The analysis of multivariate data is often complicated by high dimensionality and complex inter-dependences between the observed variables. In order to identify patterns in such data it is therefore desirable and often necessary to reduce the dimensionality and separate different aspects of the data. In multivariate statistics, for example, principal component analysis (PCA) is a common preprocessing step that decomposes the data into orthogonal principle components which are sorted according to how much variance of the original data each component explains. There are two important applications of this. Firstly, one can reduce the dimensionality of the data by projecting it onto the lower dimensional space spanned by the

---

\*Authors contributed equally.

leading principal components which maximize the explained variance. Secondly, since the principle components are orthogonal, they separate in some sense different (uncorrelated) parts of the data. In many situations this enables a better interpretation and representation.

Often, however, PCA may not be sufficient to separate the data in a desirable way due to more complex inter-dependences in the multivariate data (see e.g., Section 1.3.3 in [Hyvärinen et al. \(2001\)](#) for an instructive example). This observation motivates the development of independent component analysis (ICA), formally introduced in its current form by [Cardoso \(1989\)](#) and [Comon \(1994\)](#). ICA is a widely used unsupervised blind source separation technique that aims at decomposing an observed linear mixture of independent source signals. More precisely, assuming that the observed data is a linear mixture of underlying independent variables, one seeks the unmixing matrix that maximizes the independence between the signals it extracts. There has been a large amount of research into different types of ICA procedures and their interpretations, e.g., [Bell and Sejnowski \(1995, Infomax\)](#) who maximize the entropy, [Hyvärinen \(1999, FastICA\)](#) who maximizes the kurtosis or [Belouchrani et al. \(1997, SOBI\)](#) who minimize time-lagged dependences, to name only three of the more widespread examples.

ICA has applications in many fields, for example in finance (e.g., [Back and Weigend, 1997](#)), the study of functional magnetic resonance imaging data (e.g., [McKeown et al., 1998a,b](#); [Calhoun et al., 2003](#)), and notably in the analysis of electroencephalography (EEG) data (e.g., [Makeig et al., 1996, 1997](#); [Delorme and Makeig, 2004](#)). The latter is motivated by the common assumption that the signals recorded at EEG electrodes are a (linear) superposition of cortical dipole signals ([Nunez and Srinivasan, 2006](#)). Indeed, ICA-based preprocessing has become the de facto standard for the analysis of EEG data where the extracted components are interpreted as corresponding to cortical sources (e.g., [Ghahremani et al., 1996](#); [Zhukov et al., 2000](#); [Makeig et al., 2002](#)) or used for artifact removal by dropping components that are dominated by ocular or muscular activity (e.g., [Jung et al., 2000](#); [Delorme et al., 2007](#)).

In many applications, the data at hand is heterogeneous and often parts of the samples can be grouped by the different settings (or environments) under which the observations were taken, e.g., we can group those samples of a multi-subject EEG recording that belong to the same subject. For the analysis and interpretation of such data across different groups, it is desirable to extract one set of common features or signals instead of obtaining individual ICA decompositions for each group. In practice, whenever such situations arise, the common and tempting solution is to simply pool the data across groups and run an ICA algorithm on all data whereby the group structure is ignored.

Here, we present a novel, methodologically sound framework that extends the ordinary ICA, respects the group structure and accounts for group-wise confounding. More precisely, we consider a model of the form

$$X_i = A \cdot (S_i + H_i),$$

where  $A$  remains fixed across different groups,  $S_i$  is a vector of independent source signals and  $H_i$  is a vector of confounding variables with fixed covariance within each group. Based on this extension to ordinary ICA, we construct a method and an easy to implement algorithm to extract one common set of sources that are stable across the different groups and can be used for across-group analyses.

## 1.1. Contributions and relation to existing work

ICA is extremely well-studied with a tremendous amount of research related to various types of extensions and relaxations of the ordinary ICA model. In light of this, it is important to understand where our proposed procedure is positioned and why it is an interesting and useful extension. ICA research can be categorized and systematized according to different criteria. The first is related to the type of assumptions under which the ICA model is identifiable: roughly speaking, there is one group of non-Gaussianity- and one of non-stationarity-based methods. Many methods for model identification under the respective assumptions have been developed (e.g., [Cardoso and Souloumiac, 1993](#); [Bell and Sejnowski, 1995](#); [Belouchrani et al., 1997](#); [Hyvärinen, 1999](#); [Pham and Cardoso, 2001](#)). Our `groupICA` falls in the category of non-stationarity-based methods such as the methods by [Matsuoka et al. \(1995\)](#); [Pham and Cardoso \(2001\)](#); [Hyvärinen \(2001\)](#) and likewise relies on whitening covariance matrices across time. A second criterion to categorize ICA research is to separate different extensions of the underlying model. Some of the most prominent extensions are noisy ICA (e.g., [Moulines et al., 1997](#)), overcomplete ICA (e.g., [Lewicki and Sejnowski, 1997](#)), independent subspace analysis ([Hyvärinen and Hoyer, 2000](#)) and topographic ICA ([Hyvärinen et al., 2001](#)). The `groupICA` procedure has connections to several of these extensions which we discuss in more detail in Section 3.2.

Apart from this more general positioning within the ICA literature, our methodology is aimed at a more specific setting: data with a group structure. While such applications frequently arise in practice (e.g., multi-subject EEG/fMRI studies), theoretical results are still scarce. ICA does not naturally generalize to a method suitable to extract components that are robust across groups. In an attempt to overcome this limitation, groups are often simply pooled together and an ICA is performed on the concatenated data, which is the default behavior in toolboxes like the widely used `eeglab` for EEG analyses ([Delorme and Makeig, 2004](#)). Pooling routines like this, however, ignore any group structure and are prone to distortions resulting from variations across groups, as also illustrated in our simulations in Section 4.3. Often one can remedy some of these problems using clever normalization or denoising techniques to preprocess the data. Indeed, most existing approaches for grouped data involve some level of engineering and are hence no longer based on explicit statistical models. They are rather a combination of several preprocessing, normalization, concatenation and analysis steps (cf. [Guo and Pagnoni \(2008\)](#); [Calhoun et al. \(2009\)](#) for an overview). One frequently used example is the method introduced by [Calhoun et al. \(2001\)](#), which first whitens and projects the grouped data by PCA (thus performing a normalization within groups) followed by pooling the data and running an ordinary ICA algorithm on the concatenated data. Another technique is tensorial ICA due to [Beckmann and Smith \(2005\)](#) who add an extra dimension to their data to incorporate the groups and propose a three dimensional ICA procedure which learns a joint model of related decompositions for each group.

One strength of our methodology is that it explicitly fixes a statistical model that is sensible for data with group structure and can be estimated efficiently, while being backed by provable identification results. Furthermore, providing an explicit model with all required assumptions enables a constructive discussion about the appropriateness of such modeling decisions in specific application scenarios. The model itself is based on a notion of invariance that is robust against confounding structures from groups: this idea is also related to invariance ideas in causality ([Haavelmo, 1944](#); [Peters et al., 2016](#)), see also the relation to causality which we discuss in Section 3.3.

We believe that `groupICA` is a valuable contribution to the ICA literature in the following ways:

- We introduce a methodologically sound framework which extends ordinary ICA to settings with grouped data.
- We prove identifiability of the unmixing matrix under mild assumptions.
- We provide an easy to implement estimation procedure.
- We illustrate the usefulness, robustness, applicability, and limitations of our newly introduced `groupICA` algorithm as well as characterize the advantage of `groupICA` over pooled ICA: the source separation by `groupICA` is more stable across groups since it explicitly accounts for group-wise confounding instead of pooling the data across groups.
- We propose a way to rank the recovered components according to their stability.
- We provide an open-source scikit-learn compatible ready-to-use Python implementation available as `groupICA` from the Python Package Index repository as well as R and Matlab implementations and an intuitive audible example which is available at <https://sweichwald.de/groupICA>.

## 2. Methodology

For the model description, let  $S_i = (S_i^1, \dots, S_i^d)^\top \in \mathbb{R}^{d \times 1}$  and  $H_i = (H_i^1, \dots, H_i^d)^\top \in \mathbb{R}^{d \times 1}$  for  $i \in \{1, \dots, n\}$  be two independent vector-valued sequences of random variables. The components  $S_i^1, \dots, S_i^d$  are assumed to be mutually independent at all points  $i$  while there are no further constraints on  $H$ . Moreover, let  $A \in \mathbb{R}^{d \times d}$  be an invertible matrix. The  $d$ -dimensional data process  $(X_i)_{i \in \{1, \dots, n\}}$  is generated by the following model

$$X_i = A \cdot (S_i + H_i), \quad \text{for all } i \in \{1, \dots, n\}. \quad (2.1)$$

That is,  $X$  is a linear combination of source signals  $S$  and confounding variables  $H$ . In this model, both  $S$  and  $H$  are unobserved and one aims at recovering the pure source signals  $S$ . Without additional assumptions, however, this is impossible as there is no way to distinguish between the confounding  $H$  and the source signals  $S$  and even with additional assumptions it remains a difficult task (see Section 3.2 for an overview of related existing ICA models). Here, we aim at solving the problem of only recovering the confounded source signals  $\tilde{S}_i = S_i + H_i$ .

Throughout this paper we denote by  $\mathbf{X} = (X_1, \dots, X_n) \in \mathbb{R}^{d \times n}$  the observed data matrix and similarly by  $\mathbf{S}$  and  $\mathbf{H}$  the corresponding (unobserved) source and confounding data matrices. For a finite data sample generated by this model we hence have

$$\mathbf{X} = A \cdot (\mathbf{S} + \mathbf{H}).$$

In order to distinguish between the confounding  $H$  and the source signals  $S$  we need to assume that the two processes are sufficiently different. This can be achieved by assuming the existence of a group structure such that the covariance of the confounding  $H$  remains invariant within a group and only changes across groups.

**Assumption 1 (group-wise confounding)**

There exists a collection of  $m$  disjoint groups  $\mathcal{G} = \{g_1, \dots, g_m\}$  with  $g_k \subseteq \{1, \dots, n\}$  and  $\cup_{k=1}^m g_k = \{1, \dots, n\}$ , such that for all  $g \in \mathcal{G}$  and for all  $i \in g$  it holds that

$$\text{Cov}(H_i) \equiv \Sigma_{H,g}.$$

Under this assumption and given that the source signals change enough within groups, the mixing matrix  $A$  is identifiable, see Section 2.1. Similar to existing ICA methods, our approach to estimate the mixing matrix  $A$  is based on shifts in second moments of the source signals  $S$  (e.g., Matsuoka et al., 1995; Pham and Cardoso, 2001; Hyvärinen, 2001). In contrast to these existing methods, we also adjust for the confounding  $H$ .

For  $V = A^{-1}$  and using (2.1) it holds for all  $i \in \{1, \dots, n\}$  that

$$V \text{Cov}(X_i) V^\top = \text{Cov}(S_i) + \text{Cov}(H_i).$$

Since the source signal components  $S_i^j$  are mutually independent, the covariance matrix  $\text{Cov}(S_i)$  is diagonal. Moreover, due to Assumption 1 the covariance matrix of the confounding  $H$  is constant within each group. This implies for all groups  $g \in \mathcal{G}$  and for all  $k, l \in g$  that

$$V (\text{Cov}(X_k) - \text{Cov}(X_l)) V^\top = \text{Cov}(S_k) - \text{Cov}(S_l) \quad (2.2)$$

is a diagonal matrix. We can therefore identify  $V$  by simultaneously diagonalizing differences of covariance matrices as in (2.2). This process of finding a matrix  $V$  that simultaneously diagonalizes a set of matrices is known as joint matrix diagonalization and has been studied extensively (e.g., Ziehe et al., 2004; Tichavsky and Yeredor, 2009). In Section 2.2, we show how to construct an estimator for  $V$  based on approximate joint matrix diagonalization.

**2.1. Identifiability and theoretical properties**

Identifiability requires that the source signals  $S$  change sufficiently strong within groups. We prove identifiability under the following mild assumption that characterizes “how much” the source signals change.

**Assumption 2 (non-stationary, independently changing second moments)**

For each pair of components  $p, q \in \{1, \dots, d\}$  we require the existence of three (not necessarily unique) groups  $g_1, g_2, g_3 \in \mathcal{G}$  and three corresponding pairs  $l_1, k_1 \in g_1$ ,  $l_2, k_2 \in g_2$  and  $l_3, k_3 \in g_3$  such that the two vectors

$$\begin{pmatrix} \text{Var}(S_{l_1}^p) - \text{Var}(S_{k_1}^p) \\ \text{Var}(S_{l_2}^p) - \text{Var}(S_{k_2}^p) \\ \text{Var}(S_{l_3}^p) - \text{Var}(S_{k_3}^p) \end{pmatrix} \quad \text{and} \quad \begin{pmatrix} \text{Var}(S_{l_1}^q) - \text{Var}(S_{k_1}^q) \\ \text{Var}(S_{l_2}^q) - \text{Var}(S_{k_2}^q) \\ \text{Var}(S_{l_3}^q) - \text{Var}(S_{k_3}^q) \end{pmatrix}$$

are not collinear and not equal to zero.

This assumption can be interpreted as requiring the individual components to have pairwise independently changing variance processes. Based on this assumption, it is possible to show that the mixing matrix  $A$  is uniquely identifiable.

**Theorem 2.1 (identifiability of the mixing matrix)**

Assume the data process  $(X_i)_{i \in \{1, \dots, n\}}$  satisfies the model in (2.1) and additionally Assumption 1 and Assumption 2 are satisfied. Then  $A$  is unique up to permutation and rescaling of its columns.

**Proof** A proof is given in Appendix A. □

## 2.2. Estimation

In order to estimate  $V$  from a finite observed sample  $\mathbf{X} \in \mathbb{R}^{d \times n}$  we first partition each group into subgroups. Then, we compute the empirical covariance matrices on each subgroup. Finally, using an approximate joint matrix diagonalization technique, we estimate a matrix that simultaneously diagonalizes the differences of these empirical covariance matrices.

To make this more precise, for each group  $g \in \mathcal{G}$  we first construct a partition  $\mathcal{P}_g$  consisting of subsets of  $g$  such that each  $e \in \mathcal{P}_g$  satisfies that  $e \subseteq g$  and  $\cup_{e \in \mathcal{P}_g} e = g$ . This partition  $\mathcal{P}_g$  should be granular enough to capture the changes in the signals described in Assumption 2. We propose partitioning each group based on a grid such that the separation between grid points is large enough for a reasonable estimation of the covariance matrix and at the same time small enough to capture the variations in the signals. In our experiments, we observed robustness with respect to the exact choice and that only extreme values should be avoided (cf. Section 4.4.1).

Next, for each group  $g \in \mathcal{G}$  and each pair  $e, f \in \mathcal{P}_g$ , we define the matrix

$$M_{e,f}^g := \widehat{\text{Cov}}(\mathbf{X}_e) - \widehat{\text{Cov}}(\mathbf{X}_f),$$

where  $\widehat{\text{Cov}}(\cdot)$  denotes the empirical covariance matrix and  $\mathbf{X}_e$  is the data matrix restricted to the columns corresponding to the subgroup  $e$ . From (2.2) it follows that  $VM_{e,f}^gV^\top$  should be approximately diagonal. We are therefore interested in finding an invertible matrix  $V$  which approximately jointly diagonalizes all the matrices in the set

$$\mathcal{M}^{\text{all}} := \{M_{e,f}^g \mid g \in \mathcal{G} \text{ and } e, f \in \mathcal{P}_g\}. \quad (2.3)$$

The number of matrices in this set grows quadratically in the number of partitions. In practice, this can lead to large numbers of matrices to be diagonalized. Another option that reduces the computational load is to compare each partition to its complement, which leads to the following collection of matrices

$$\mathcal{M}^{\text{comp}} := \{M_{e,\bar{e}}^g \mid g \in \mathcal{G} \text{ and } e \in \mathcal{P}_g \text{ (with } \bar{e} := g \setminus e)\}. \quad (2.4)$$

The task of jointly diagonalizing a set of matrices is a well-studied topic in the literature and is referred to as approximate joint matrix diagonalization. Many solutions have been proposed for different assumptions made on the matrices to be diagonalized. In this paper we use the `uwedge` algorithm<sup>1</sup> introduced by Tichavsky and Yeredor (2009). The basic idea behind `uwedge` is to find a minimizer of a proxy for the following loss function

$$\ell(V) = \sum_{M \in \mathcal{M}^*} \left( \sum_{k \neq l} [VMV^\top]_{k,l}^2 \right),$$

over the set of invertible matrices, where  $\mathcal{M}^* \in \{\mathcal{M}^{\text{all}}, \mathcal{M}^{\text{comp}}\}$  (cf. (2.3) and (2.4)).

The full estimation procedure based on the set  $\mathcal{M}^{\text{comp}}$  defined in (2.4) is made explicit in the pseudo code in Algorithm 1 (where `ApproximateJointDiagonalizer` stands for a general approximate joint diagonalizer, e.g., `uwedge`).

---

<sup>1</sup>As a byproduct of our work, we are able to provide a new stable open-source Python/R/Matlab implementation of the `uwedge` algorithm which is also included in our respective `groupICA` packages.

---

**Algorithm 1: groupICA**

---

**input** : data matrix  $\mathbf{X}$   
          group index  $\mathcal{G}$   
          group-wise partition  $(\mathcal{P}_g)_{g \in \mathcal{G}}$  (user selected)  
initialize empty list  $\mathcal{M}$   
**for**  $g \in \mathcal{G}$  **do**  
    **for**  $e \in \mathcal{P}_g$  **do**  
      append  $\widehat{\text{Cov}}(\mathbf{X}_e) - \widehat{\text{Cov}}(\mathbf{X}_{\bar{e}})$  to list  $\mathcal{M}$   
    **end**  
**end**  
 $\widehat{V} \leftarrow \text{ApproximateJointDiagonalizer}(\mathcal{M})$   
 $\widehat{\mathbf{S}} \leftarrow \widehat{V}\mathbf{X}$   
**output**: unmixing matrix  $\widehat{V}$   
          sources  $\widehat{\mathbf{S}}$

---

### 3. Stability, related ICA models and a causal perspective

A key aspect of our model is that, given the existence of a group structure, it aims to improve the stability of the unmixing across these groups. Here, we refer to the following notion of stability: if we estimate unmixing matrices on two different groups we want the two unmixing matrices to be the same in the sense that both extract the same set of sources. Given that the data generating process is truly given by our confounded ICA model in (2.1) a standard ICA method is not able to estimate the correct unmixing  $V = A^{-1}$ . Hence, it extracts different mixtures of sources and confounding across different groups and is not stable in the aforementioned sense. This is illustrated by the “America’s Got Talent Duet Problem” (cf. Example 3.1), an extension and alteration of the classical “cocktail party problem”.

#### Example 3.1 (America’s Got Talent Duet Problem)

*Consider the problem of evaluating two singers at a duet audition individually. This requires to somehow listen to the two voices separately, while the singers perform simultaneously. There are two sound sources in the audition room (the two singers) and additionally several noise sources which corrupt the recordings at the two microphones (or the jury member’s two ears). A schematic of such a setting is illustrated in Figure 1. The additional noise comes from an audience and two open windows. One can assume that this noise satisfies our Assumption 1 on a single group. The sound stemming from the audience can be seen as an average of many little sounds, hence remaining approximately constant over time. Also typical sounds from an open window often satisfy this assumption, for example think of sound from a river or a busy road. Our methodology, however, also allows for more complicated settings in which the noise shifts at known points in times, for example if someone opens or closes a window or starts mowing the lawn outside. In such cases we use the known time blocks of constant noise as groups and apply groupICA on this grouped data. An example with artificial sound data related to this setting is available at <https://sweichwald.de/groupICA>, where we show that groupICA is able to recover useful sound signals, while ICA on pooled data fails to unmix the two singers.*

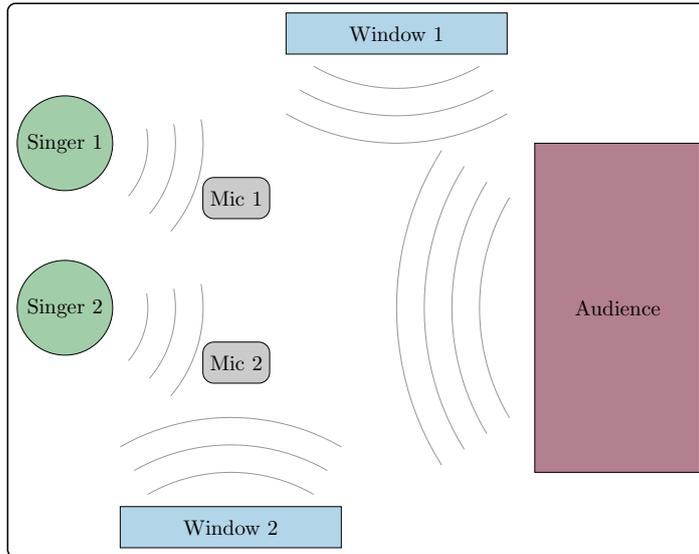


Figure 1. Schematic of the “America’s Got Talent Duet Problem” described in Example 3.1. The sound from the windows and audience is taken to be confounding noise which has fixed covariance structure over given time blocks. The challenge is to recover the sound signals from the individual singers given the recordings of the two microphones.

### 3.1. Ranking recovered sources by stability

The ordering of the independent components in ICA is not identifiable since there is no natural ordering between them. For practical purposes, however, it is often useful to rank the components in some meaningful way. One common option, is to proceed similar to PCA and rank components by their explained variance, i.e., the amount of variance in the observed data that can be explained by each component. Such a ranking can be misleading, especially when the observations are corrupted by confounding noise, as it only takes into account the strength of the signal. We propose a different type of ranking which builds upon the ICA or groupICA model and ranks the recovered components by their stability across observations.

To this end, let  $A = (a_1, \dots, a_d) \in \mathbb{R}^{d \times d}$  be the mixing and  $V = (v_1, \dots, v_d)^\top \in \mathbb{R}^{d \times d}$  the corresponding unmixing matrix (i.e.,  $V = A^{-1}$ ,  $a_i$  are columns of  $A$  and  $v_i$  are rows of  $V$ ). Then it holds that,

$$\begin{aligned}
 \text{Cov}(X_i)v_j^\top &= A[\text{Cov}(S_i) + \text{Cov}(H_i)]A^\top v_j^\top \\
 &= A[\text{Cov}(S_i) + \text{Cov}(H_i)]e_j^\top \\
 &= A \text{Cov}(S_i)e_j + A \text{Cov}(H_i)e_j^\top
 \end{aligned} \tag{3.1}$$

Under our group-wise confounding assumption (Assumption 1) this implies that within all groups  $g \in \mathcal{G}$ , it holds for all  $l, k \in g$  that

$$(\text{Cov}(X_l) - \text{Cov}(X_k))v_j^\top = a_j \left( \text{Var}(S_l^j) - \text{Var}(S_k^j) \right). \tag{3.2}$$

This reflects the contribution of the  $j$ -th recovered source  $S^j$  (in terms of differences) to the variance of each component of the original multivariate data  $X$ . While in the population case

the equality in (3.2) is satisfied exactly this is no longer the case when empirically estimating the mixing (and unmixing) matrix. Consider for example two subsets  $e, f \in g$  for some group  $g \in \mathcal{G}$ , then using the notation from Section 2.2 and denoting by  $\hat{v}_j$  and  $\hat{a}_j$  the estimated counterparts of  $v_j$  and  $a_j$ , respectively, it holds that

$$\begin{aligned} M_{e,f}^g \hat{v}_j^\top &= [\widehat{\text{Cov}}(\mathbf{X}_e) - \widehat{\text{Cov}}(\mathbf{X}_f)] \hat{v}_j^\top \\ &= \widehat{A} [\widehat{\text{Cov}}(\widehat{S}_e) - \widehat{\text{Cov}}(\widehat{S}_f)] \widehat{A}^\top \hat{v}_j^\top \\ &= \widehat{A} \underbrace{[\widehat{\text{Cov}}(\widehat{S}_e) - \widehat{\text{Cov}}(\widehat{S}_f)]}_{\approx \text{Id}(\text{Var}(S_e) - \text{Var}(S_f))} e_j^\top, \end{aligned} \quad (3.3)$$

where the approximation only holds if our **groupICA** procedure is able to correctly unmix the  $j$ -th source. The extent to which this approximation holds for different components across various subsets  $e, f$  gives a natural measure of how stable each component is recovered. In particular, it allows us to construct an intuitive ranking of recovered components by comparing how well  $\hat{a}_j$  correlates with the terms  $M_{e,f}^g \hat{v}_j^\top$ .

More precisely, for each recovered source we compute its stability by

$$p^j := \frac{1}{|\mathcal{M}^*|} \sum_{M \in \mathcal{M}^*} \left| \widehat{\text{Corr}}(M \hat{v}_j^\top, \hat{a}_j) \right|$$

where  $\widehat{\text{Corr}}(\cdot, \cdot)$  denotes the empirical correlation coefficient and  $\mathcal{M}^* \in \{\mathcal{M}^{\text{all}}, \mathcal{M}^{\text{comp}}\}$  (cf. (2.3) and (2.4)). Based on these stabilities  $p^1, \dots, p^d$  we then rank the components from most stable to least stable. Interestingly, this ranking has a direct connection to stability of topographic maps in EEG experiments, which we explain in more detail in Section 4.4.3.

Such a stability ranking is of course by no means restricted to the **groupICA** model. In fact, a similar idea can be used in the ordinary ICA model where, due to the lack of confounding, it holds that

$$\text{Cov}(X_i) v_j^\top = a_j \text{Var}(S_i^j), \quad (3.4)$$

which can similarly be utilized for a corresponding stability ranking of the components.

### 3.2. Relation to existing ICA methods

It is informative to understand our **groupICA** model (2.1) in relation to three commonly used ICA models in the literature: ordinary, overcomplete and noisy ICA.

The ordinary ICA model assumes that the observed process  $X$  is a linear mixture of independent source signals  $S$  without a confounding term  $H$ . Identifiability of the source signals  $S$  is then guaranteed by assumptions on  $S$  as for example, non-Gaussianity or specific time structures. One problem with the ordinary ICA model is that it requires at least as many observed signals as there are source signals. In practice, this assumption is often violated since the underlying data generating process consists of many more independent variables than are observed (or measured).

Overcomplete ICA (e.g., [Lewicki and Sejnowski, 1997](#)) extends the ordinary ICA model to these settings. The key difficulty in this extension is that the mixing and unmixing matrices become non-identifiable, as one is essentially interested in solving an underdetermined system of linear equations. Therefore, existing solutions to this problem generally add additional assumptions, e.g., assuming that at most as many source signals are active at the same time

as observed signals exists (see Lee et al., 1999). By considering the confounding terms  $H$  in our model (2.1) as additional source signals our model can be seen as an overcomplete ICA model, i.e., more source signals than observed signals. Hence, our model allows to solve the overcomplete ICA model by separating out the true (or interesting) signals from the group-wise constant sources that are viewed as confounding in our model. The drawback being that the obtained sources are still “contaminated” by the confounding.

A further related ICA model is known as noisy ICA (e.g., Moulines et al., 1997) in which the data generating process is assumed to be an ordinary ICA model with additive noise. Similarly, as in the overcomplete ICA model this leads to identifiability issues, which in this setting are generally resolved by assuming that the additive noise is Gaussian, hence allowing a separation of the true (non-Gaussian) sources  $S$  from the noise. Once again this is closely related to our groupICA model in the sense that our confounding term  $H$  is also additive noise. The difference is that we assume that the noise remains constant across groups while circumventing distributional assumptions on the sources (non-Gaussian) and noise (Gaussian) as is necessary in noisy ICA.

### 3.3. Causal interpretation

There is a strong connection between ICA and the problem of identifying structural causal models (SCMs) introduced by Pearl (2009). Shimizu et al. (2006) were the first to make use of this connection by using ICA to infer causal structures. To make this more precise consider the following linear SCM

$$X = B \cdot X + \varepsilon. \quad (3.5)$$

Assuming that the matrix  $\text{Id} - B$  is invertible, we can rewrite this equation as

$$X = (\text{Id} - B)^{-1} \varepsilon,$$

which corresponds to the underlying linear ICA model with mixing matrix  $A = (\text{Id} - B)^{-1}$  and source signals  $S = \varepsilon$ . Instead of taking the noise term  $\varepsilon$  as independent noise sources one can also consider  $\varepsilon = S + H$ . Then, the linear SCM in (3.5) describes a causal model between the  $X$  variables containing hidden confounding. This is illustrated in Figure 2, which depicts a 3 variable SCM with feedback loops and confounding. Since the sources have inherently unidentifiable scale and permutation it is, in general, not straight forward to recover the causal effect matrix  $B$  from the mixing matrix  $A$ . This is only possible under additional restrictions on the causal structure, for example by restricting to acyclic graphs as in Shimizu et al. (2006) or by restricting the types of allowed feedback loops as in Hoyer et al. (2008) and Rothenhäusler et al. (2015).

## 4. Experiments

In this section, we analyze empirical properties of our procedure. To this end, we first introduce two performance measures, which enable us to empirically assess the quality of the recovered components in terms of stability across groups. We illustrate the performance of groupICA as compared to a pooled version of ICA on simulated data with and without confounding. We also compare on real data and outline potential benefits of using our method when analyzing multi-subject EEG data.

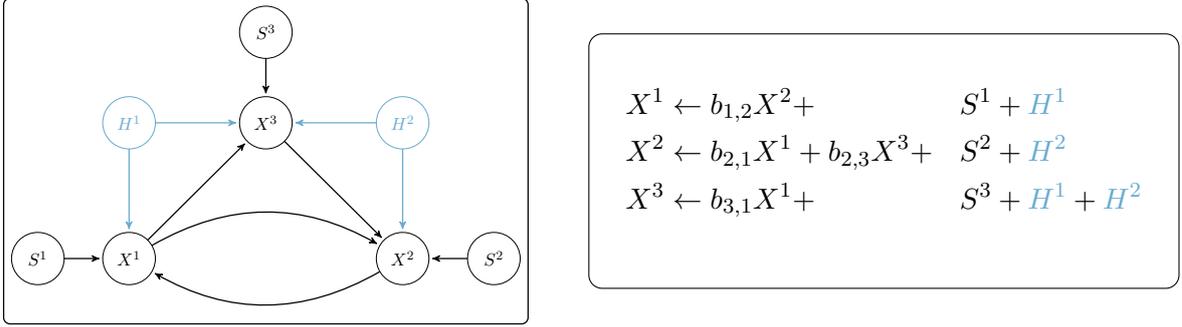


Figure 2. Illustration of an SCM with (including colored nodes  $H^1$ ,  $H^2$ ) and without (excluding colored nodes) confounding.

#### 4.1. Assessing the quality of recovered sources

Assessing the quality of the recovered sources in an ICA setting is an inherently difficult task, as is typical for unsupervised learning procedures. In particular, the unidentifiable scale and ordering of the sources as well as the unclear choice of a performance measure render this task difficult. Provided that ground truth is known, one way to measure how close the estimated components  $\widehat{V} = (\widehat{v}_1, \dots, \widehat{v}_d)^\top$  are to the true ones  $V = (v_1, \dots, v_d)^\top$  is by the following order- and scale-invariant correlation accuracy (CA) score

$$\text{CA}(V, \widehat{V}) := \sum_{k=1}^d \max_{i,j \in I_k \times J_k} |\widehat{\text{Corr}}(v_i, \widehat{v}_j)|,$$

where  $I_1 \times J_1 = \{1, \dots, d\}^2$  and  $I_{k+1} \times J_{k+1} = I_k \times J_k \setminus \{\arg \max_{i,j \in I_k \times J_k} |\widehat{\text{Corr}}(v_i, \widehat{v}_j)|\}$  is recursively defined as the set of all row-index pairs that have previously been unmatched. This leads to a greedy matching of row pairs, ensuring that eventually each row of  $V$  is matched to exactly one row of  $\widehat{V}$ . In short, the CA score measures the (empirical) correlation between rows of the true and the estimated unmixing matrix, while the order indeterminacy is resolved by a greedy pairwise matching of the rows.

Obviously, we require a different performance measure for our real data experiments where the true unmixing matrix is unknown. Here, we check whether the desired independence (without constant confounding) is achieved by computing the following covariance instability score (CIS) matrix which measures the instability of the covariance structure of the unmixed sources  $\widehat{\mathbf{S}}$  and is defined for a set of groups  $\mathcal{G}$  and a partition  $(\mathcal{P}_g)_{g \in \mathcal{G}}$  (see Section 2.2) by

$$\text{CIS}(\widehat{\mathbf{S}}, \mathcal{G}, (\mathcal{P}_g)_{g \in \mathcal{G}}) := \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} \frac{2}{|\mathcal{P}_g|(|\mathcal{P}_g| - 1)} \sum_{e \in \mathcal{P}_g} \left( \frac{\widehat{\text{Cov}}(\widehat{\mathbf{S}}_e) - \widehat{\text{Cov}}(\widehat{\mathbf{S}}_{\bar{e}})}{\widehat{\sigma}_{\widehat{\mathbf{S}}_g} \cdot \widehat{\sigma}_{\widehat{\mathbf{S}}_g}^\top} \right)^2,$$

where  $\widehat{\sigma}_{\widehat{\mathbf{S}}} \in \mathbb{R}^{d \times 1}$  is the empirical standard deviation of  $\widehat{\mathbf{S}}$  and the fraction is taken element-wise. The CIS matrix is approximately diagonal whenever  $\widehat{\mathbf{S}}$  can be written as the sum of independent source signals  $\mathbf{S}$  and confounding  $\mathbf{H}$  with fixed covariance. This is condensed into one scalar that reflects how stable the sources' covariance structure is by averaging the

off-diagonals of the CIS matrix

$$\text{MCIS}(\widehat{\mathbf{S}}, \mathcal{G}, (\mathcal{P}_g)_{g \in \mathcal{G}}) := \frac{1}{d(d-1)} \sum_{\substack{i,j=1 \\ i \neq j}}^d \left[ \text{CIS}(\widehat{\mathbf{S}}, \mathcal{G}, (\mathcal{P}_g)_{g \in \mathcal{G}}) \right]_{i,j}.$$

The mean covariance instability score (MCIS) can be understood as a measure of independence of the signal process  $S$  after subtracting the constant confounding term. High values of MCIS imply large dependences (beyond constant confounding) between the signals and low values imply small dependences. In particular, the measure remains valid in settings without confounding in which case small values imply uncorrelated sources across observations. Depending on the approximate joint matrix diagonalizer used in the `groupICA` estimation, the MCIS is directly or indirectly minimized. We therefore expect that this score is small on the training data whenever the joint diagonalization was successful.

## 4.2. Competing methods

In all of our numerical experiments, we apply `groupICA` as outlined in Algorithm 1, where we partition each group based on equally spaced grids, run a fixed number of  $10 \cdot 10^3$  iterations of the uwedge approximate joint diagonalizer and compare it to two other methods. The first, is a pooled version of the popular FastICA algorithm due to Hyvärinen (1999), where we simply neglect the group structure. As implementation we use the FastICA implementation from the scikit-learn Python library due to Pedregosa et al. (2011) and take the defaults for all parameters. The comparison to this method aims to illustrate that it is indeed useful to explicitly account for existing group-structures. For the remainder of the paper we will denote this method by `pooledICA`. The second method we compare to is a random projection of the sources, where the unmixing matrix is simply sampled with iid standard normal entries. The idea of this comparison is to give a baseline of the unmixing problem. Moreover, it gives some quantification to the variance of our assessment measures which we discussed in Section 4.1. In order to illustrate the variance in this method, we generally sample 100 random projections and show the results for each of them.

## 4.3. Simulations

In this section, we investigate empirical properties of `groupICA` in simulated and hence well-controlled scenarios. First off, we show that we can recover the correct mixing matrix given that the data is generated according to our model (2.1) and Assumptions 1 and 2 hold, while `pooledICA` necessarily falls short in this setting (cf. Section 4.3.1). Moreover, in Section 4.3.2 we show that even in the absence of any confounding (i.e., when the data follows the ordinary ICA model and  $H \equiv 0$  in our model) we remain competitive with ordinary ICA. All of our simulations are based on block-wise shifted variance data, which we describe in Data Set 1.

### 4.3.1. Dependence on confounding strength

For this simulation experiment, we sample data according to Data Set 1 and choose to simulate  $n = 100 \cdot 10^3$  (dimension  $d = 22$ ) samples from  $m = 20$  groups where each group contains  $n/m = 5 \cdot 10^3$  observations. Within each group, we then select a random partition consisting of 20 (i.e.,  $|\mathcal{P}_g| = 20$ ) subsets while ensuring that these have the same size on average. We fix the

### Block-wise shifted variance data

For our simulations, we select  $m$  equally sized groups  $\mathcal{G} := \{g_1, \dots, g_m\}$  on the data points  $\{1, \dots, n\}$  and for each group  $g \in \mathcal{G}$  construct a partition  $\mathcal{P}_g$ . Then, we sample a model of the form

$$X_i = A \cdot (S_i + C \cdot H_i),$$

where the values on the right-hand side are sampled as follows:

- $A, C \in \mathbb{R}^{d \times d}$  are sampled with iid entries from  $\mathcal{N}(0, 1)$  and  $\mathcal{N}(0, \frac{1}{d})$ , respectively.
- For each  $g \in \mathcal{G}$  the variables  $H_i \in \mathbb{R}^m$  are sampled from  $\mathcal{N}(0, \sigma_g^2 \text{Id}_m)$ , where the  $\sigma_g^2$  are sampled iid from  $\text{Unif}(0.1, b_1)$ .
- For each  $g \in \mathcal{G}$  and  $e \in \mathcal{P}_g$  the variables  $S_i \in \mathbb{R}^d$  are sampled from  $\mathcal{N}(0, \eta_e^2 \text{Id}_d)$ , where the  $\eta_e^2$  are sampled iid from  $\text{Unif}(0.1, b_2)$ .

The parameters  $b_1$  and  $b_2$  are selected in such a way that the expected signal strength  $c_1 := \mathbb{E}(|\eta_e^2 - \eta_f^2|)$  and confounding strength  $c_2 := \mathbb{E}(|\sigma_g^2 - \sigma_h^2|)$  are as dictated by the respective experiment. Due to the uniform distribution this reduces to

$$b_1 = 3(c_1 + 0.1) \quad \text{and} \quad b_2 = 3(c_2 + 0.1).$$

Data Set 1. Description of the synthetic data sets used in the simulations in Section 4.3.

signal strength to  $c_1 = 1$  and consider the behavior of **groupICA** (applied to half of the groups with an equally spaced grid of 10 partitions per group) for different confounding strengths  $c_2 = \{0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4\}$ . The results for 1000 repetitions are shown in Figure 3.

The results indicate that in terms of the CA score **pooledICA** becomes worse as the confounding strength increases and systematically estimates an incorrect unmixing matrix. While **groupICA** shows an increased variance of the CA score as the confounding strength increases, it is less severely affected in terms of bias; the increase in variance is expected due to the decreasing signal to noise ratio. In terms of MCIS the behavior is analogous; with increasing confounding strength the **pooledICA** unmixing estimation is systematically biased resulting in high MCIS scores both out-of-sample and even in-sample.

#### 4.3.2. Efficiency in absence of group confounding

For this simulation experiment, we sample data according to Data Set 1 and we choose to simulate  $n = 20 \cdot 10^3$  (dimension  $d = 22$ ) samples from  $m = 10$  groups where each group contains  $n/m = 10^3$  observations. Within each group, we then select a random partition consisting of 10 (i.e.,  $|\mathcal{P}_g| = 10$ ) subsets while ensuring that these have the same size on average. This time, to illustrate performance in the absence of confounding, we fix the confounding strengths  $c_2 = 0$  and consider the behavior of **groupICA** (applied to half of the groups with an equally spaced grid of 4 partitions per group) for different signal strengths  $c_1 = \{0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4\}$ . The results for 1000 repetitions are shown in Figure 4.

The results indicate that overall **groupICA** performs competitive in the confounding-free case. In particular, there is no drastic negative hit on the performance of **groupICA** as compared to **pooledICA** in settings where the data follows the ordinary ICA model. In sum-

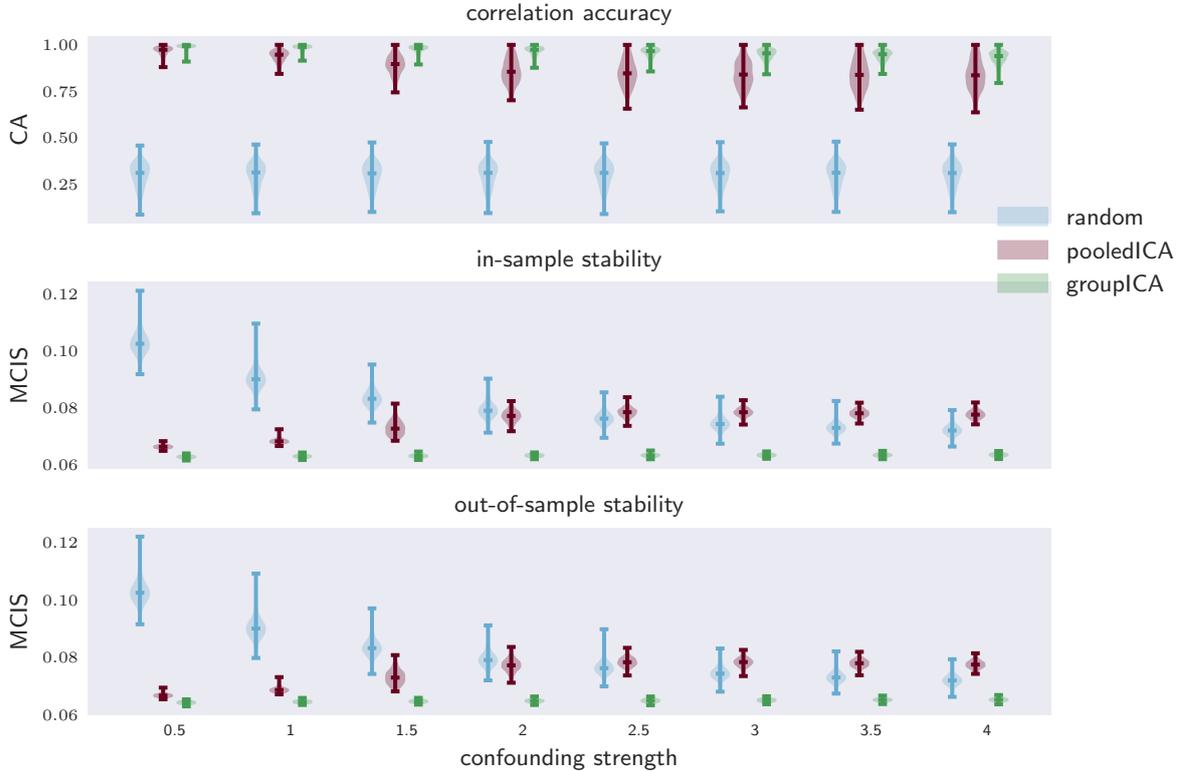


Figure 3. Results of the simulation experiment described in Section 4.3.1. Plot shows performance measures (CA: large implies close to truth; MCIS: small implies stable) for fixed signal strength and various confounding strengths. The difference between pooledICA and groupICA is more prominent for higher confounding strengths where the estimates of pooledICA are increasingly different from the true unmixing matrix and the sources become increasingly unstable.

mary, this means that groupICA performs well on a larger model class consisting of both the group-wise confounded as well as the confounding-free case. More precisely, an advantage over pooledICA is gained in confounded cases (as shown in Section 4.3.1) while there is no (strong) disadvantage in unconfounded cases. This suggests that the ordinary ICA algorithms are to be preferred if the data is known to exactly follow the ordinary ICA model, otherwise, if in doubt about the presence of hidden confounding, one may be better off using groupICA and allowing for a richer model class.

#### 4.4. EEG experiments

To illustrate the applicability of our method to real data, we apply it to two publicly available EEG data sets, CovertAttention and BCICompIV2a as described in Data Set 2 and Data Set 3, respectively. For both data sets, we compare the recovered sources of groupICA with those recovered by pooledICA and the sources received from a random projection. Since ground truth is unknown we report comparisons in the following sections that are based on the following three criteria:

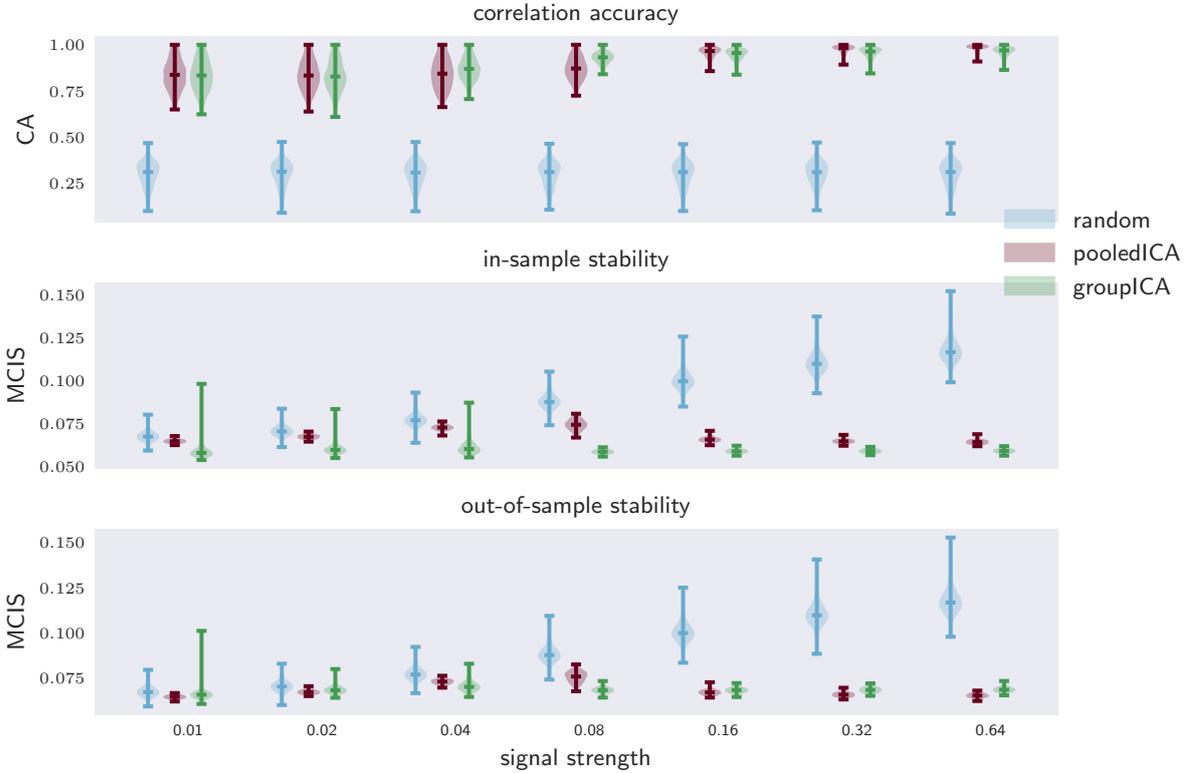


Figure 4. Results of the simulation experiment described in Section 4.3.2. Plot shows performance measures (CA: large implies close to truth; MCIS: small implies stability) for data generated without confounding and for various signal strengths. These results are reassuring, as they indicate that when applied to data that follows the ordinary ICA model, `groupICA` still performs competitive to `pooledICA` although it allows for a richer model class.

### stability and independence

We use MCIS (cf. Section 4.1) to assess the stability and independence of the recovered sources both in- and out-of-sample.

### classification accuracy

For both data sets there is label information available that associates certain time windows of the EEG recordings with the task the subjects were performing at that time. Based on the recovered sources, we build a classification pipeline relying on feature extraction and classification techniques that are common in the field (Lotte et al., 2018). The achieved classification accuracy serves as a proxy of how informative and suitable for this common feature and classification pipeline the extracted signals are.

### topographies

For a qualitative assessment, we inspect the topographic maps of the extracted sources, as well as the corresponding power spectra and a raw time series chunk. This is used to illustrate that the sources recovered by `groupICA` do not appear random or implausible for EEG recordings and are qualitatively similar to what is expected from ordinary ICA.

### CovertAttention data

This data set is due to [Treder et al. \(2011\)](#) and consists of EEG recordings of 8 subjects performing multiple trials of covertly shifting visual attention to one out of 6 cued directions. The data set contains recordings of

- 8 subjects,
- for each subject there exist 6 runs with 100 trials,
- each recording consists of 60 EEG channels recorded at 1000 Hz sampling frequency, while we work with the publicly available data that is downsampled to 200 Hz.

Since visual inspection of the data revealed data segments with huge artifacts and details about how the publicly available data was preprocessed was unavailable to us, we removed outliers and high-pass filtered the data at 0.5 Hz. In particular, along each dimension we set those values to the median along its dimension that deviate more than 10 times the median absolute distance from this median. We further preprocess the data by re-referencing to common average reference (car) and projecting onto the orthogonal complement of the null component. For our unmixing estimations, we use the entire data, i.e., including intertrial breaks.

For classification experiments (cf. Section 4.4.2) we use, in line with [Treder et al. \(2011\)](#), the 8–12 Hz bandpass-filtered data during the 500–2000 ms window of each trial, and use the log-variance as bandpower feature ([Lotte et al., 2018](#)).

Data Set 2. Description of the CovertAttention data, a publicly available EEG data set of subjects performing a covert attention shift task.

#### 4.4.1. Stability and independence

We aim to probe stability not only in-sample but also verify the expected increase in stability when applying the unmixing matrix to data of new unseen subjects, i.e., to new groups of samples with different confounding specific to that subject. In order to analyze stability in terms of the MCIS both in- and out-of-sample and for different amounts of training samples, we proceed by repeatedly splitting the data into a training and a test data set. More precisely, we construct all possible splits into training and test subjects for any given number of training subjects. For each pair of training and test set, we fit an unmixing matrix using `groupICA`, `pooledICA` and `random` as described in Section 4.2. We then compute the MCIS on the training and test data for each method separately and collect the results of each training-test split with the same number of subjects used for training.

Results obtained on the CovertAttention data set (with equally spaced partitions of  $\approx 15$  seconds length) are given in Figure 5 and the results for the BCICompIV2a data set (with equally spaced partitions of  $\approx 9$  seconds length) are shown in Appendix B.1, Figure 8. For both data sets the results are qualitatively similar and support the claim that the unmixing obtained by `groupICA` is more stable when transferred to new unseen subjects. While for `pooledICA` the instability on held-out subjects does not follow a clear decreasing trend with increasing number of training subjects, `groupICA` can successfully make use of additional training subjects to learn a more stable unmixing matrix.

### BCICompIV2a data

This data set is due to [Tangermann et al. \(2012, Section 5\)](#) and consists of EEG recordings of 9 subjects performing multiple trials of 4 different motor imagery tasks. The data set contains recordings of

- 9 subjects, each recorded on 2 different days,
- for each subject and day there exist 6 runs with 48 trials,
- each recording consists of 22 EEG channels recorded at 250 Hz sampling frequency,
- and is bandpass filtered between 0.5 and 100 Hz and is 50 Hz notch filtered.

For our analysis we only use the trial-data, i.e., the concatenated segments of seconds 3–6 of each trial (corresponding to the motor imagery part of the trials ([Tangermann et al., 2012](#))). We further preprocess the data by re-referencing to common average reference (car) and projecting onto the orthogonal complement of the null component.

As features for classification experiments (cf. [Section 4.4.2](#)) on this data set we use band-power in the 8–30 Hz band as measured by the log-variance of the 8–30 Hz bandpass-filtered trial data ([Lotte et al., 2018](#)).

Data Set 3. Description of the BCICompIV2a data that consists of a publicly available motor imagery data set.

The results for random projections provide a reference for the variance in MCIS which varies for different numbers of training subjects that determine the number of training-test splits. Furthermore, they hint at a general problem when assessing blind source separation results. As expected, random projections perform quite well in terms of stability when transferred to new subjects (out-of-sample) since they avoid any bias towards training data whatsoever. In other aspects, however, these sources may be clearly inferior, though stable, which we demonstrate in the following section.

To further probe the behavior for different partition grid sizes, we repeat the experiment for an equal training-test subject split (i.e., 4 training/4 test subjects for CovertAttention and 4 training/5 test sets for BCICompIV2a) for varying partition grid sizes. The upshot is, the higher the number of partitions the greater the stability gain of `groupICA` over `pooledICA`. This suggests that a small partition grid size that still allows for a reasonable covariance estimation is preferable. The detailed results are shown in [Appendix B.2](#), [Figure 9](#) and [Figure 10](#).

#### 4.4.2. Classification based on recovered sources

While the results in the previous section indicate that `groupICA` can recover more stable sources from EEG signals than `pooledICA`, they also reveal that random projections are stable when assessed with the MCIS. This implies that in scenarios with an unknown ground truth the stability of the recovered sources cannot serve as the sole determining criterion for assessing the quality of recovered sources. In addition to asking whether the recovered sources are stable, we hence also need to investigate whether the sources extracted by `groupICA` are

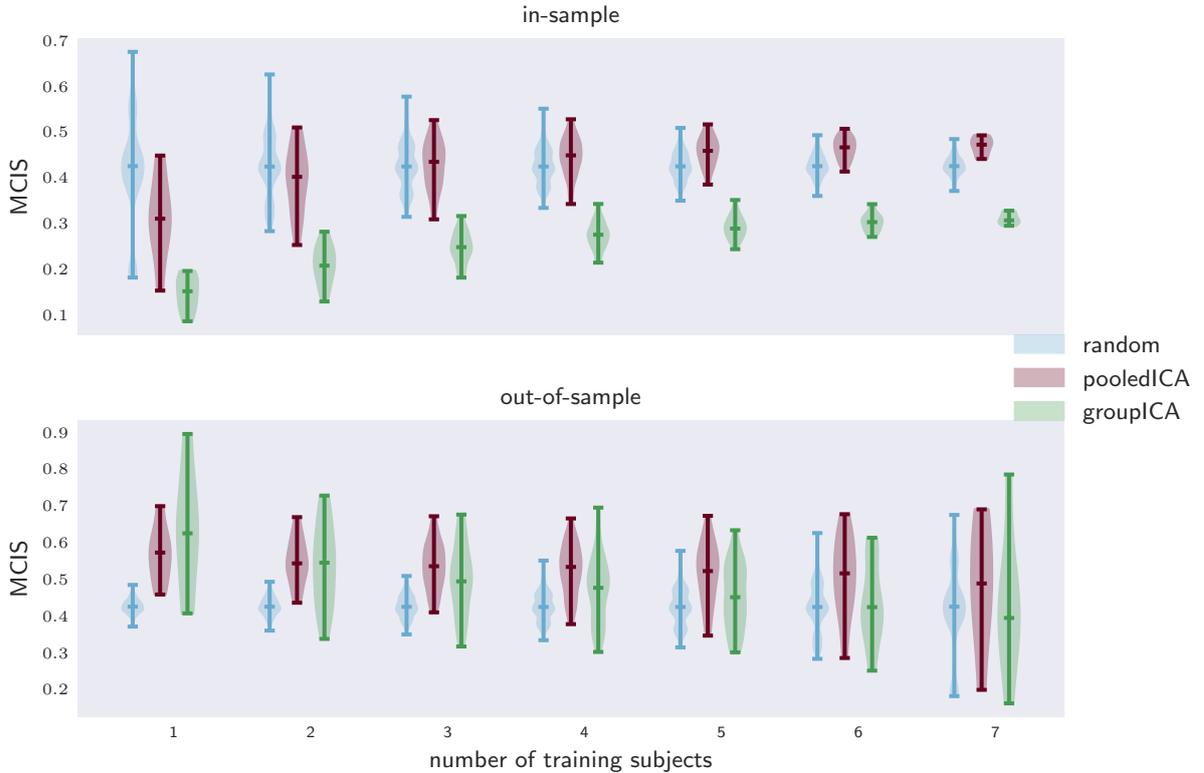


Figure 5. Experimental results for comparing the stability of sources (MCIS: small implies stable) trained on different numbers of training subjects by either `groupICA` or `pooledICA` (cf. Section 4.4.1), here on the CovertAttention Data Set 2, demonstrating that `groupICA` in contrast to `pooledICA` can successfully incorporate more training subjects to learn more stable unmixing matrices when applied to new unseen subjects. Random projections are expected to perform quite well out-of-sample as they avoid any bias towards the training data. They are shown here only as a point of reference since they are clearly not an option for real data analysis.

in fact reasonable or meaningful (as, most likely, opposed to the randomly projected sources). In the “America’s Got Talent Duet Problem” (cf. Example 3.1) this means that each of the recovered sources should only contain the voice of one singer. For EEG data, this assessment is not quite as easy. Here, we approach this problem from two angles: (a) in this section we show that the recovered sources are informative and suitable for common EEG classification pipelines, (b) in Section 4.4.3 we qualitatively assess the extracted sources based on their power spectra and topographic maps.

In both data sets there are labeled trials, i.e., segments of data during which the subject covertly shifts attention to one of six cues (cf. Data Set 2) or performs one of four motor imagery tasks (cf. Data Set 3). Based on these, one can try to predict the trial label given the trial EEG data. To mimic a situation where the sources are transferred from other subjects, we assess the informativeness of the extracted sources in a leave-one-subject-out fashion as follows. We use the unmixing matrix estimated on data from all but one subject and apply it to the held-out subject to obtain the extracted source signals. On top of these, we compute

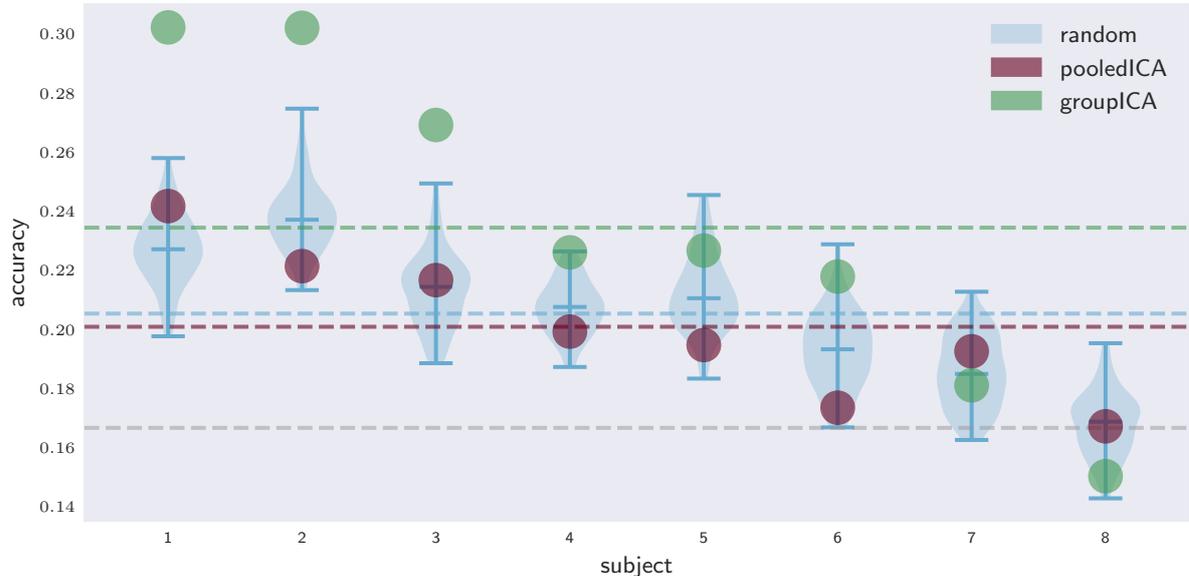


Figure 6. Cross-validated classification accuracy achieved for each subject when using a linear support vector classifier and the components obtained on the other remaining subjects by `groupICA`, `pooledICA`, or 100 random projections (cf. Section 4.4.2), here on the CovertAttention Data Set 2, demonstrating that components extracted by `groupICA` contain label information. The green/red/blue/gray dashed line corresponds to `groupICA` mean, `pooledICA` mean, random projection mean, and random guessing level, respectively. For an average of about 310 to be classified trials per subject an accuracy above 0.203 is significant for non-trivial classification at a 95% confidence level.

bandpower features for each trial (as described in Data Set 2 and 3) and use a linear support vector classifier in a one-vs-rest scheme for multiclass classification of the task labels from these features.

The results are reported in Figure 6 and Appendix B.3, Figure 11 which show for each subject the cross-validated classification accuracy<sup>2</sup> when using the unmixing obtained on the remaining subjects by either `groupICA`, `pooledICA` or 100 random unmixing matrices. The results on both data sets support the claim that the sources recovered by `groupICA` are not only stable but in addition also capture meaningful aspects of the data that enable classification accuracies that cannot be achieved based on random projections of the EEG signals or when using `pooledICA`.

It is worth noting that these classification results depend heavily on the employed classification pipeline following the source separation. Here, our goal is only to show that `groupICA` does indeed separate the data into informative sources. In practice, and when only classification accuracy matters, one might also consider using a label-informed source separation (Dähne et al., 2014), employ common spatial patterns (Koles et al., 1990) or use decoding techniques based on Riemannian geometry (Barachant et al., 2012). For completeness, in Appendix B.3 we show accuracies when using a perceptron classifier or a linear discriminant

<sup>2</sup>We average the accuracies obtained on 200 random stratified 9:1 splits into training and test data.

analysis classifier with shrinkage instead of the linear support vector classifier.

### 4.4.3. Topographic maps

The components that `groupICA` extracts from EEG signals are stable (cf. Section 4.4.1) and meaningful in the sense that they contain information that enables classification of trial labels, which is a common task in EEG studies (cf. Section 4.4.2). In this section, we complement the assessment of the recovered sources by demonstrating that the results obtained by `groupICA` lead to topographies, activation maps, power spectra and raw time series that are similar to what is commonly obtained during routine ICA analyses of EEG data when the plausibility and nature of ICA components is to be judged.

Topographies are common in the EEG literature to depict the relative projection strength of extracted sources to the scalp sensors. More precisely, the column-vector  $a_j$  of  $A = V^{-1}$  that specifies the mixing of the  $j$ -th source component is visualized as follows. A sketched top view of the head is overlaid with a heatmap where the value at each electrodes' position is given by the corresponding entry in  $a_j$ . These topographies are indicative of the nature of the extracted sources, for example the dipolarity of source topographies is a criterion invoked to identify cortical sources (Delorme et al., 2012) or the topographies reveal that the source mainly picks up changes in the electromagnetic field induced by eye movements. Another way to visualize an extracted source is an activation map, which is obtained by depicting the vector  $\widehat{\text{Cov}}(X)v_j^\top$  (where  $v_j$  is  $j$ -th row of unmixing matrix  $V$ ) and shows for each electrode how the signal observed at that electrode covaries with the signal extracted by  $v_j$  (Haufe et al., 2014). Besides inspecting the raw time series data, another criterion invoked to separate cortical from muscular components is the log power spectrum. For example, a monotonic increase in spectral power starting at around 20 Hz is understood to indicate muscular activity (Goncharova et al., 2003) and peaks in typical EEG frequency ranges are used to identify brain-related components.<sup>3</sup>

In Figure 7, we depict the aforementioned criteria for three exemplary components (2<sup>nd</sup>, 6<sup>th</sup> and 51<sup>st</sup> based on our stability ranking) extracted by `groupICA` on the CovertAttention Data Set 2. The idea is to demonstrate that `groupICA` components are qualitatively similar in these aspects to components extracted by the commonly employed ICA. Therefore, we choose to display one example of an ocular component (2<sup>nd</sup> where the topography is indicative of eye movement), a cortical component (6<sup>th</sup> where the dipolar topography, the typical frequency peak at around 8–12 Hz, and the amplitude modulation visible in the raw time series are indicative of the cortical nature), and an artifactual component (51<sup>st</sup> where the irregular topography and the high frequency components indicate an artifact). For comparison, we additionally show for each component the topographies of the components extracted by `pooledICA` or the random projection by matching the recovered source which most strongly correlates with the one extracted by `groupICA`. Once again, this illustrates, that, while random components are stable, they lead to irregular topographies. This is not the case for the components extracted by `groupICA`, the topographies of which closely resemble the results one would obtain from a commonly employed ICA analysis on EEG data.

---

<sup>3</sup>These are commonly employed criteria which are also advised in the eeglab tutorial (Delorme and Makeig, 2004, [https://sccn.ucsd.edu/wiki/Chapter\\_09:\\_Decomposing\\_Data\\_Using\\_ICA](https://sccn.ucsd.edu/wiki/Chapter_09:_Decomposing_Data_Using_ICA)) and the neurophysiological biomarker toolbox wiki (Hardstone et al., 2012, [https://www.nbtwiki.net/doku.php?id=tutorial:how\\_to\\_use\\_ica\\_to\\_remove\\_artifacts](https://www.nbtwiki.net/doku.php?id=tutorial:how_to_use_ica_to_remove_artifacts)).

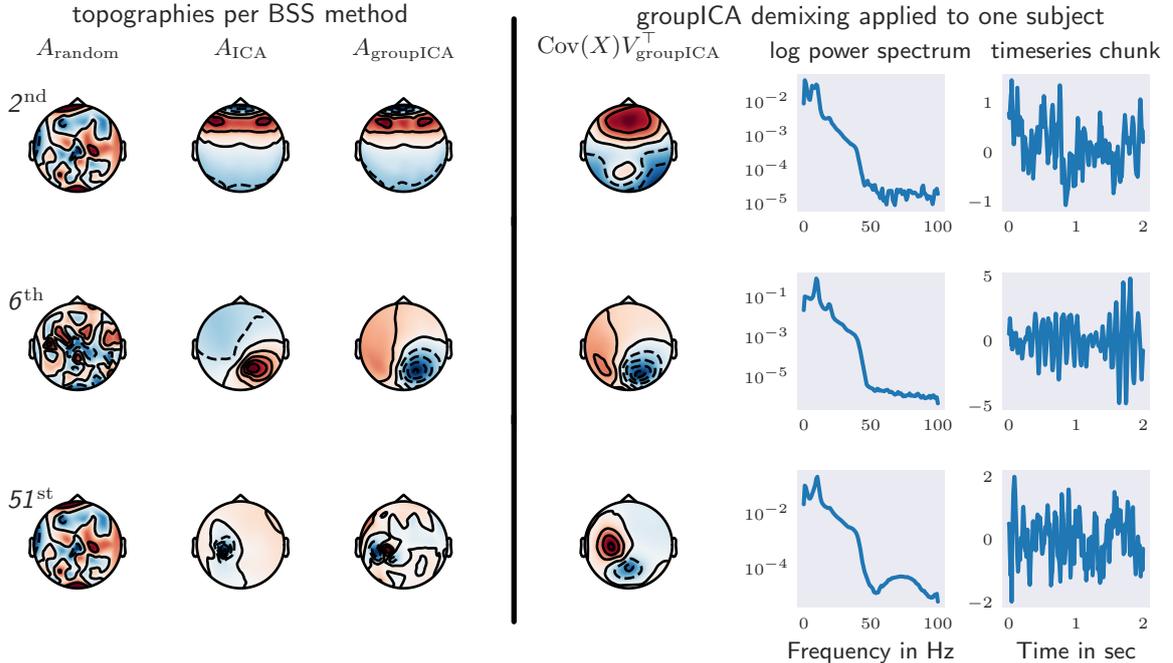


Figure 7. Visualization of exemplary EEG components recovered on the CovertAttention Data Set 2. On the left the topographies of three components are shown where the mixing matrix is the inverse of a random unmixing matrix ( $A_{\text{random}}$ ), the unmixing matrix obtained by pooledICA ( $A_{\text{ICA}}$ ) and that of **groupICA** ( $A_{\text{groupICA}}$ ). On the right we depict, for a randomly chosen subject, the activation maps, the log power spectra, and randomly chosen chunks of the raw time series data corresponding to the respective **groupICA** components. Components extracted by **groupICA** are qualitatively similar to those of the commonly employed pooledICA; see Section 4.4.3 for details.

## 5. Conclusion

In this paper, we construct a method for recovering independent sources corrupted by group-wise confounding. It extends ordinary ICA to an easily interpretable model which we believe is relevant for many practical problems as is demonstrated in Section 4.4 for EEG data. Based on this model, we give explicit assumptions under which one can expect the sources to be identifiable in the population case (cf. Section 2.1). Moreover, we introduce a straightforward method for estimating the sources based on the well-understood concept of approximate joint matrix diagonalization. As illustrated in the simulations in Section 4.3, this estimation procedure performs competitive even for data from an ordinary ICA model, while additionally being able to adjust for group-wise confounding. Finally, we show that the **groupICA** model indeed performs reasonably on typical EEG data sets and leads to improvements in comparison to naive pooling approaches, while at the same time preserving an easy interpretation of the recovered sources.

## Acknowledgements

The authors thank Nicolai Meinshausen, Jonas Peters, and Vinay Jayaram for helpful discussions.

## References

- Back, A. and Weigend, A. (1997). A first application of independent component analysis to extracting structure from stock returns. *International Journal of Neural Systems*, 8(4):473–484.
- Barachant, A., Bonnet, S., Congedo, M., and Jutten, C. (2012). Multiclass Brain-Computer Interface Classification by Riemannian Geometry. *IEEE Transactions on Biomedical Engineering*, 59(4):920–928.
- Beckmann, C. and Smith, S. (2005). Tensorial extensions of independent component analysis for multisubject fMRI analysis. *NeuroImage*, 25(1):294–311.
- Bell, A. and Sejnowski, T. (1995). An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129–1159.
- Belouchrani, A., Abed-Meraim, K., Cardoso, J.-F., and Moulines, E. (1997). A blind source separation technique using second-order statistics. *IEEE Transactions on Signal Processing*, 45(2):434–444.
- Calhoun, V., Adalı, T., Hansen, L., Larsen, J., and Pekar, J. (2003). ICA of functional MRI data: an overview. In *Proceedings of the International Workshop on Independent Component Analysis and Blind Signal Separation*, pages 281–288.
- Calhoun, V., Adalı, T., Pearlson, G., and Pekar, J. (2001). A method for making group inferences from functional MRI data using independent component analysis. *Human Brain Mapping*, 14(3):140–151.
- Calhoun, V., Liu, J., and Adalı, T. (2009). A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data. *NeuroImage*, 45(1):S163–S172.
- Cardoso, J.-F. (1989). Blind Identification Of Independent Components With Higher-order Statistics. In *Workshop on Higher-Order Spectral Analysis*, pages 157–162.
- Cardoso, J.-F. and Souloumiac, A. (1993). Blind beamforming for non-Gaussian signals. *IEE Proceedings F - Radar and Signal Processing*, 140(6):362–370.
- Comon, P. (1994). Independent component analysis, A new concept? *Signal Processing*, 36(3):287–314.
- Dähne, S., Meinecke, F., Haufe, S., Höhne, J., Tangermann, M., Müller, K.-R., and Nikulin, V. (2014). SPoC: a novel framework for relating the amplitude of neuronal oscillations to behaviorally relevant parameters. *NeuroImage*, 86:111–122.
- Delorme, A. and Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *Journal of Neuroscience Methods*, 134(1):9–21.

- Delorme, A., Palmer, J., Onton, J., Oostenveld, R., and Makeig, S. (2012). Independent EEG Sources Are Dipolar. *PLOS One*, 7:1–14.
- Delorme, A., Sejnowski, T., and Makeig, S. (2007). Enhanced detection of artifacts in EEG data using higher-order statistics and independent component analysis. *NeuroImage*, 34(4):1443–1449.
- Ghahremani, D., Makeig, S., Jung, T., Bell, A., and Sejnowski, T. (1996). Independent Component Analysis of Simulated EEG Using a Three-Shell Spherical Head Model. Technical report, Naval Health Research Center San Diego CA.
- Goncharova, I., McFarland, D., Vaughan, T., and Wolpaw, J. (2003). EMG contamination of EEG: spectral and topographical characteristics. *Clinical Neurophysiology*, 114(9):1580–1593.
- Guo, Y. and Pagnoni, G. (2008). A unified framework for group independent component analysis for multi-subject fMRI data. *NeuroImage*, 42(3):1078–1093.
- Haavelmo, T. (1944). The Probability Approach in Econometrics. *Econometrica*, 12:iii–115.
- Hardstone, R., Poil, S.-S., Schiavone, G., Jansen, R., Nikulin, V., Mansvelder, H., and Linkenkaer-Hansen, K. (2012). Detrended fluctuation analysis: a scale-free view on neuronal oscillations. *Frontiers in Physiology*, 3:450.
- Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.-D., Blankertz, B., and Bießmann, F. (2014). On the interpretation of weight vectors of linear models in multivariate neuroimaging. *NeuroImage*, 87:96–110.
- Hoyer, P., Shimizu, S., Kerminen, A., and Palviainen, M. (2008). Estimation of causal effects using linear non-Gaussian causal models with hidden variables. *International Journal of Approximate Reasoning*, 49(2):362–378.
- Hyvärinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE Transactions on Neural Networks*, 10(3):626–634.
- Hyvärinen, A. (2001). Blind source separation by nonstationarity of variance: a cumulant-based approach. *IEEE Transactions on Neural Networks*, 12(6):1471–1474.
- Hyvärinen, A. and Hoyer, P. (2000). Emergence of Phase- and Shift-Invariant Features by Decomposition of Natural Images into Independent Feature Subspaces. *Neural Computation*, 12(7):1705–1720.
- Hyvärinen, A., Hoyer, P., and Inki, M. (2001). Topographic Independent Component Analysis. *Neural Computation*, 13(7):1527–1558.
- Hyvärinen, A., Karhunen, J., and Oja, E. (2001). *Independent Component Analysis*. John Wiley & Sons, Inc.
- Jung, T.-P., Makeig, S., Humphries, C., Lee, T.-W., Mckeown, M., Iragui, V., and Sejnowski, T. (2000). Removing electroencephalographic artifacts by blind source separation. *Psychophysiology*, 37(2):163–178.

- Kleinstenber, M. and Shen, H. (2013). Uniqueness Analysis of Non-Unitary Matrix Joint Diagonalization. *IEEE Transactions on Signal Processing*, 61(7):1786–1796.
- Koles, Z. J., Lazar, M. S., and Zhou, S. Z. (1990). Spatial patterns underlying population differences in the background EEG. *Brain Topography*, 2(4):275–284.
- Lee, T.-W., Lewicki, M., Girolami, M., and Sejnowski, T. (1999). Blind source separation of more sources than mixtures using overcomplete representations. *IEEE Signal Processing Letters*, 6(4):87–90.
- Lewicki, M. and Sejnowski, T. (1997). Learning nonlinear overcomplete representations for efficient coding. In *Advances in Neural Information Processing Systems (NIPS 10)*, pages 556–562.
- Lotte, F., Bougrain, L., Cichocki, A., Clerc, M., Congedo, M., Rakotomamonjy, A., and Yger, F. (2018). A review of classification algorithms for EEG-based brain-computer interfaces: a 10 year update. *Journal of Neural Engineering*, 15(3):031005.
- Makeig, S., Bell, A., Jung, T.-P., and Sejnowski, T. (1996). Independent Component Analysis of Electroencephalographic Data. In *Advances in Neural Information Processing Systems (NIPS 8)*, pages 145–151.
- Makeig, S., Jung, T.-P., Bell, A., Ghahremani, D., and Sejnowski, T. (1997). Blind separation of auditory event-related brain responses into independent components. *Proceedings of the National Academy of Sciences*, 94(20):10979–10984.
- Makeig, S., Westerfield, M., Jung, T.-P., Enghoff, S., Townsend, J., Courchesne, E., and Sejnowski, T. (2002). Dynamic Brain Sources of Visual Evoked Responses. *Science*, 295(5555):690–694.
- Matsuoka, K., Ohoya, M., and Kawamoto, M. (1995). A neural net for blind separation of nonstationary signals. *Neural Networks*, 8(3):411–419.
- McKeown, M., Jung, T.-P., Makeig, S., Brown, G., Kindermann, S., Lee, T.-W., and Sejnowski, T. (1998a). Spatially independent activity patterns in functional MRI data during the Stroop color-naming task. *Proceedings of the National Academy of Sciences*, 95(3):803–810.
- McKeown, M., Makeig, S., Brown, G., Jung, T.-P., Kindermann, S., Bell, A., and Sejnowski, T.-J. (1998b). Analysis of fMRI data by blind separation into independent spatial components. *Human Brain Mapping*, 6(3):160–188.
- Moulines, E., Cardoso, J.-F., and Gassiat, E. (1997). Maximum likelihood for blind separation and deconvolution of noisy signals using mixture models. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 5, pages 3617–3620. IEEE.
- Nunez, P. and Srinivasan, R. (2006). *Electric Fields of the Brain: The Neurophysics of EEG*. Oxford University Press.
- Pearl, J. (2009). *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York, USA, 2nd edition.

- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Peters, J., Bühlmann, P., and Meinshausen, N. (2016). Causal inference using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 78(5):947–1012.
- Pham, D.-T. and Cardoso, J.-F. (2001). Blind separation of instantaneous mixtures of non-stationary sources. *IEEE Transactions on Signal Processing*, 49(9):1837–1848.
- Rothenhäusler, D., Heinze, C., Peters, J., and Meinshausen, N. (2015). BACKSHIFT: Learning causal cyclic graphs from unknown shift interventions. In *Advances in Neural Information Processing Systems (NIPS 28)*, pages 1513–1521.
- Shimizu, S., Hoyer, P., Hyvärinen, A., and Kerminen, A. (2006). A Linear Non-Gaussian Acyclic model for Causal Discovery. *Journal of Machine Learning Research*, 7(10):2003–2030.
- Tangermann, M., Müller, K.-R., Aertsen, A., Birbaumer, N., Braun, C., Brunner, C., Leeb, R., Mehring, C., Miller, K., Müller-Putz, G., Nolte, G., Pfurtscheller, G., Preissl, H., Schalk, G., Schlögl, A., Vidaurre, C., Waldert, S., and Blankertz, B. (2012). Review of the BCI Competition IV. *Frontiers in Neuroscience*, 6:55.
- Tichavsky, P. and Yeredor, A. (2009). Fast Approximate Joint Diagonalization Incorporating Weight Matrices. *IEEE Transactions on Signal Processing*, 57(3):878–891.
- Treder, M., Bahramisharif, A., Schmidt, N., Van Gerven, M., and Blankertz, B. (2011). Brain-computer interfacing using modulations of alpha activity induced by covert shifts of attention. *Journal of NeuroEngineering and Rehabilitation*, 8(1):24.
- Zhukov, L., Weinstein, D., and Johnson, C. (2000). Independent component analysis for EEG source localization. *IEEE Engineering in Medicine and Biology Magazine*, 19(3):87–96.
- Ziehe, A., Laskov, P., Nolte, G., and Müller, K.-R. (2004). A Fast Algorithm for Joint Diagonalization with Non-orthogonal Transformations and its Application to Blind Source Separation. *Journal of Machine Learning Research*, 5(7):777–800.

## A. Supplementary proof

**Proof** The theorem is proven by the correct invocation of [Kleinsteuber and Shen \(2013, Theorem 1\)](#). We first define the unmixing matrix  $V = A^{-1}$  and introduce the set of matrices

$$\mathcal{D} := \{V(\text{Cov}(X_k) - \text{Cov}(X_l))V^\top \mid g \in \mathcal{G} \text{ and } k, l \in g\}.$$

Due to the underlying group ICA model and Assumption 1 all matrices in the set  $\mathcal{D}$  are diagonal (cf. (2.2)). Moreover, for  $g \in \mathcal{G}$  and  $k, l \in g$  it holds that

$$\begin{aligned} V(\text{Cov}(X_k) - \text{Cov}(X_l))V^\top &= \text{Cov}(S_k) - \text{Cov}(S_l) \\ &= \text{diag}(\text{Var}(S_k^1) - \text{Var}(S_l^1), \dots, \text{Var}(S_k^d) - \text{Var}(S_l^d)). \end{aligned}$$

Using notation as in [Kleinsteuber and Shen \(2013\)](#) we define for all  $j \in \{1, \dots, d\}$  the vectors

$$\mathbf{z}_j = \left( \left( \text{Var}(S_k^j) - \text{Var}(S_l^j) \right)_{k, l \in g} \right)_{g \in \mathcal{G}}.$$

Then, Assumption 2 implies for all distinct pairs  $p, q \in \{1, \dots, d\}$  that

$$|\widehat{\text{Corr}}(\mathbf{z}_p, \mathbf{z}_q)| = \frac{|\mathbf{z}_p \cdot \mathbf{z}_q|}{\|\mathbf{z}_p\| \|\mathbf{z}_q\|} < 1.$$

Hence, it holds that  $\rho(\mathcal{D}) < 1$ , where  $\rho$  is as defined in [Kleinsteuber and Shen \(2013\)](#). Hence, we can invoke [Kleinsteuber and Shen \(2013, Theorem 1\)](#) to conclude that any matrix  $M \in \mathbb{R}^{d \times d}$  which satisfies that  $MDM^\top$  is diagonal for all  $D \in \mathcal{D}$  is equal to the identity matrix up to scaling and permutation of its columns.

Finally, assume there exists another (invertible) mixing matrix  $\tilde{A}$  such that for all  $g \in \mathcal{G}$  and all  $k, l \in g$  it holds that

$$\tilde{A}^{-1}(\text{Cov}(X_k) - \text{Cov}(X_l))(\tilde{A}^{-1})^\top = \text{Cov}(S_k) - \text{Cov}(S_l).$$

Then, it also holds that

$$(V\tilde{A})(\text{Cov}(S_k) - \text{Cov}(S_l))(V\tilde{A})^\top = V(\text{Cov}(X_k) - \text{Cov}(X_l))V^\top,$$

which is diagonal. By the above reasoning it follows that  $V\tilde{A}$  is equal to the identity matrix up to permutation and rescaling of its columns. Moreover, this implies that  $\tilde{A}$  is equal to  $A$  up to scaling and permutation of its columns. This completes the proof of [Theorem 2.1](#).  $\square$

## B. Complementary material showing results of EEG experiments

### B.1. Stability and independence

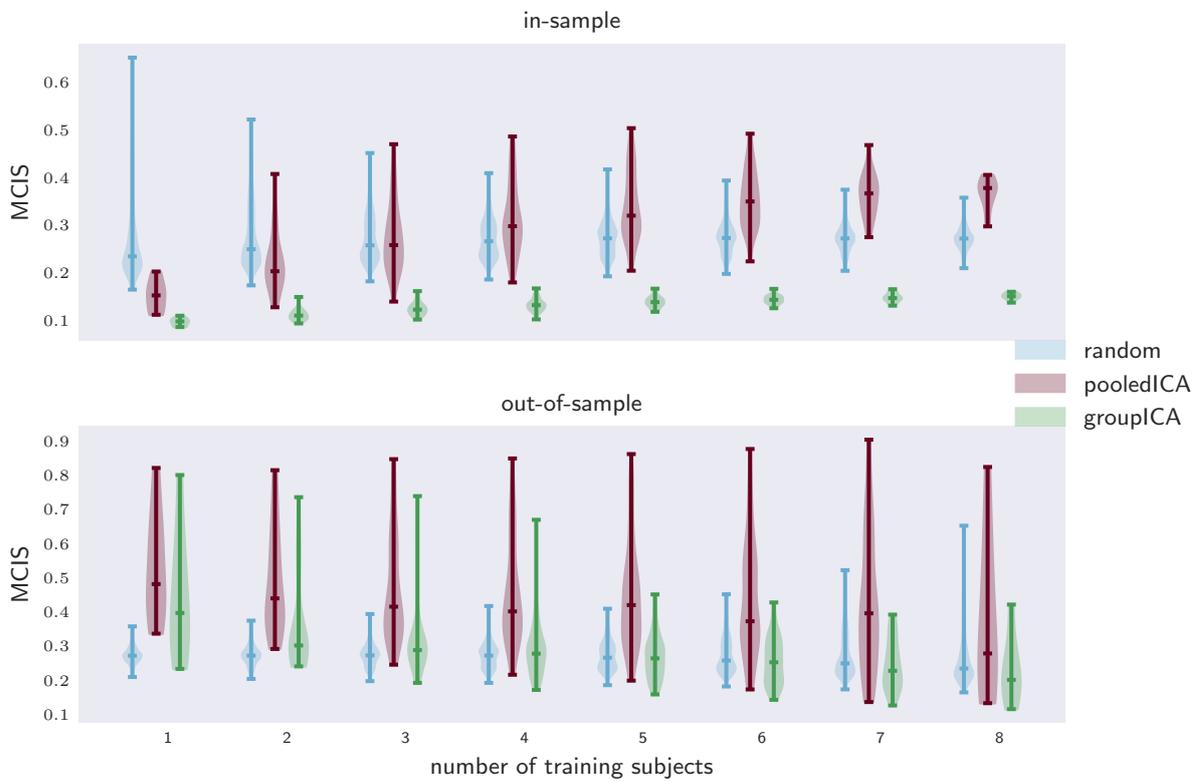


Figure 8. Experimental results for comparing the stability of sources (MCIS: small implies stable) trained on different numbers of training subjects by either `groupICA` or `pooledICA` (cf. Section 4.4.1), here on the BCICompIV2a Data Set 3, demonstrating that `groupICA` can successfully incorporate more training subjects to learn more stable unmixing matrices when applied to new unseen subjects.

## B.2. Dependence on partition grid size

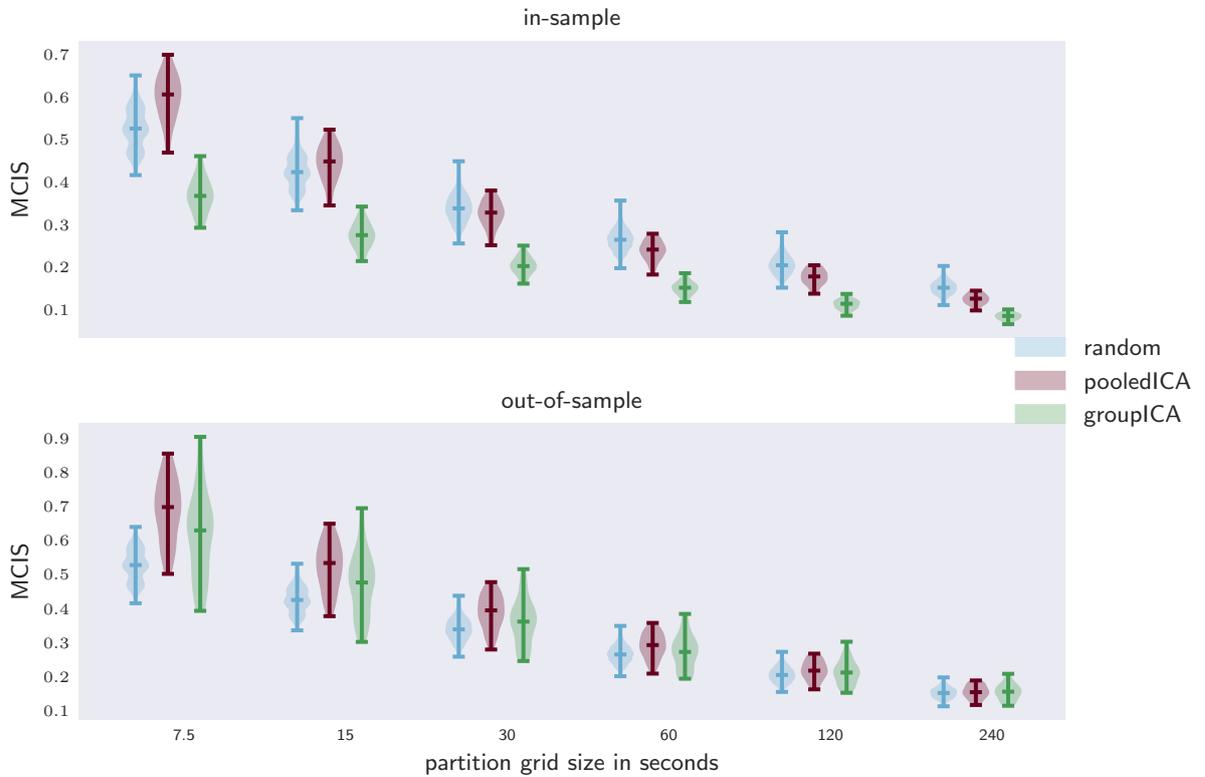


Figure 9. Experimental results for comparing the stability of sources (MCIS: small implies stable) trained for varying partition grid sizes by either groupICA or pooledICA (cf. Section 4.4.1), here on the CovertAttention Data Set 2, demonstrating a greater gain in stability of groupICA over pooledICA for smaller partition grid sizes.

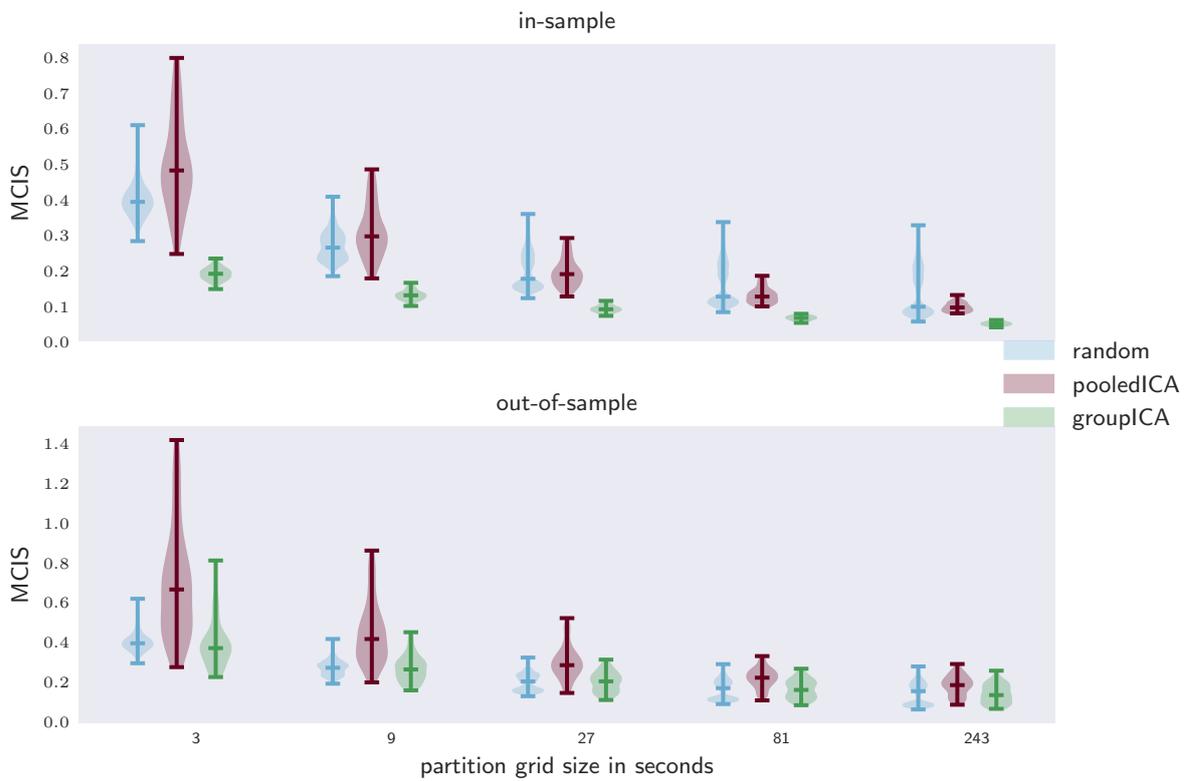


Figure 10. Experimental results for comparing the stability of sources (MCIS: small implies stable) trained for varying partition grid sizes by either `groupICA` or `pooledICA` (cf. Section 4.4.1), here on the BCICompIV2a Data Set 3, demonstrating a greater gain in stability of `groupICA` over `pooledICA` for smaller partition grid sizes.

### B.3. Classification

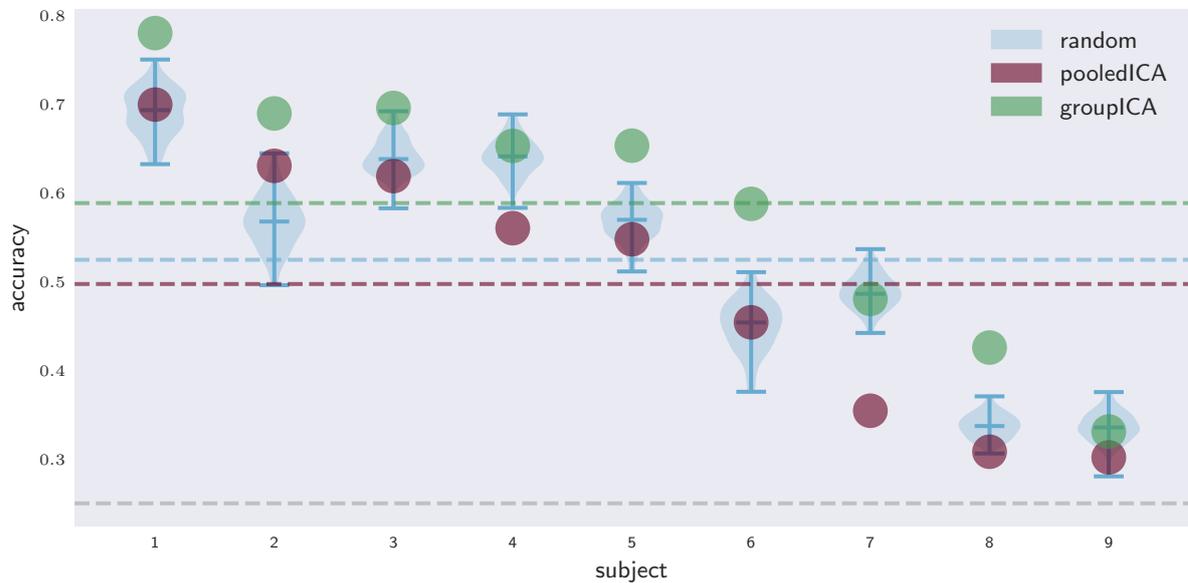


Figure 11. Cross-validated classification accuracy achieved for each subject when using a linear support vector classifier and the components obtained on the other remaining subjects by `groupICA`, `pooledICA`, or 100 random projections (cf. Section 4.4.2), here on the BCICompIV2a Data Set 3, demonstrating that components extracted by `groupICA` contain label information. The green/red/blue/gray dashed line corresponds to `groupICA` mean, `pooledICA` mean, random projection mean, and random guessing level, respectively. For the 576 to be classified trials per subject an accuracy above 0.280 is significant for non-trivial classification at a 95% confidence level.

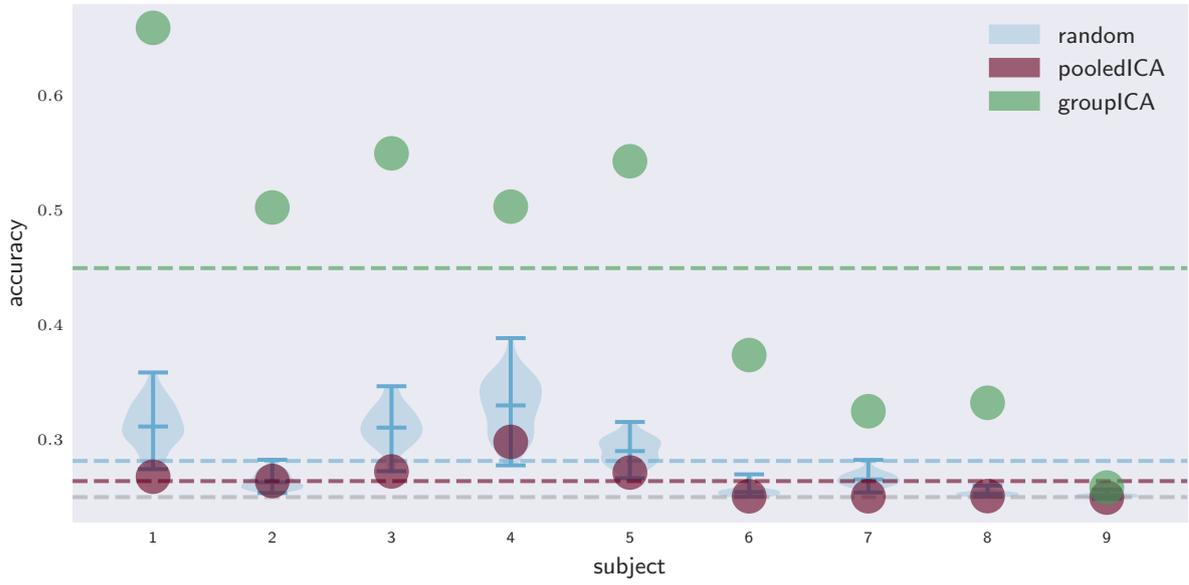


Figure 12. Results as in Figure 11 (on BCICompIV2a data set) using a perceptron classifier instead of a linear support vector classifier.

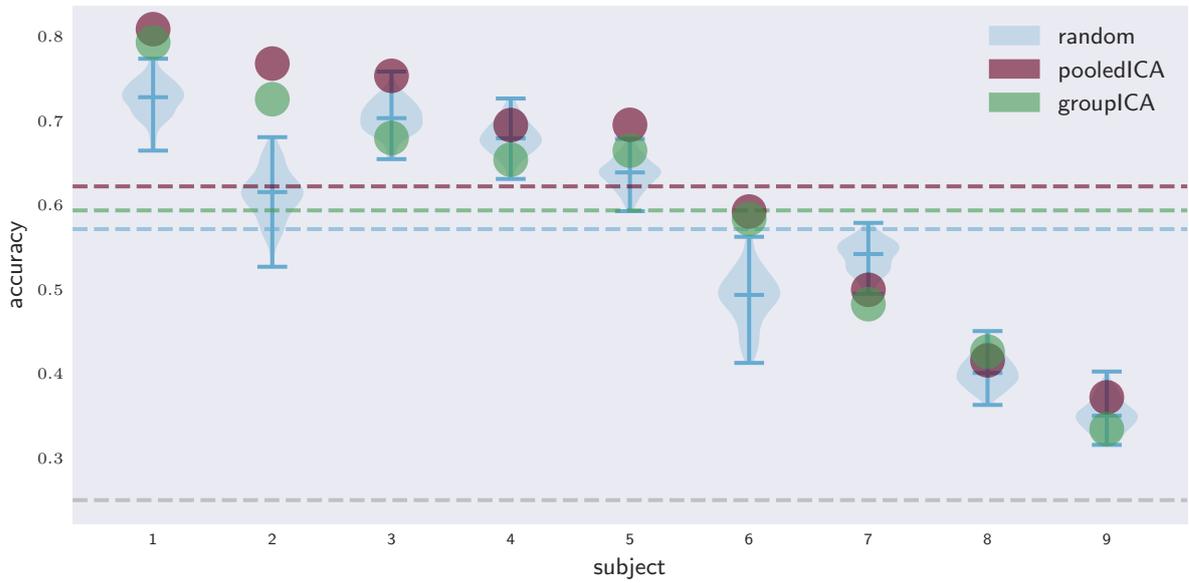


Figure 13. Results as in Figure 11 (on BCICompIV2a data set) using a shrinkage linear discriminant analysis instead of a linear support vector classifier.

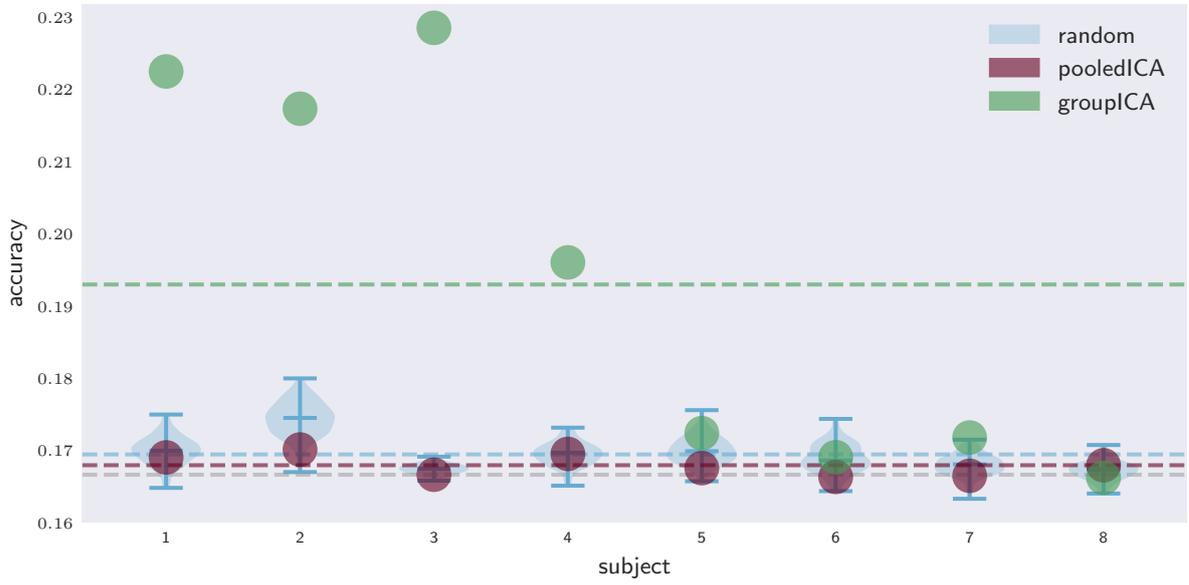


Figure 14. Results as in Figure 6 (on CovertAttention data set) using a perceptron classifier instead of a linear support vector classifier.

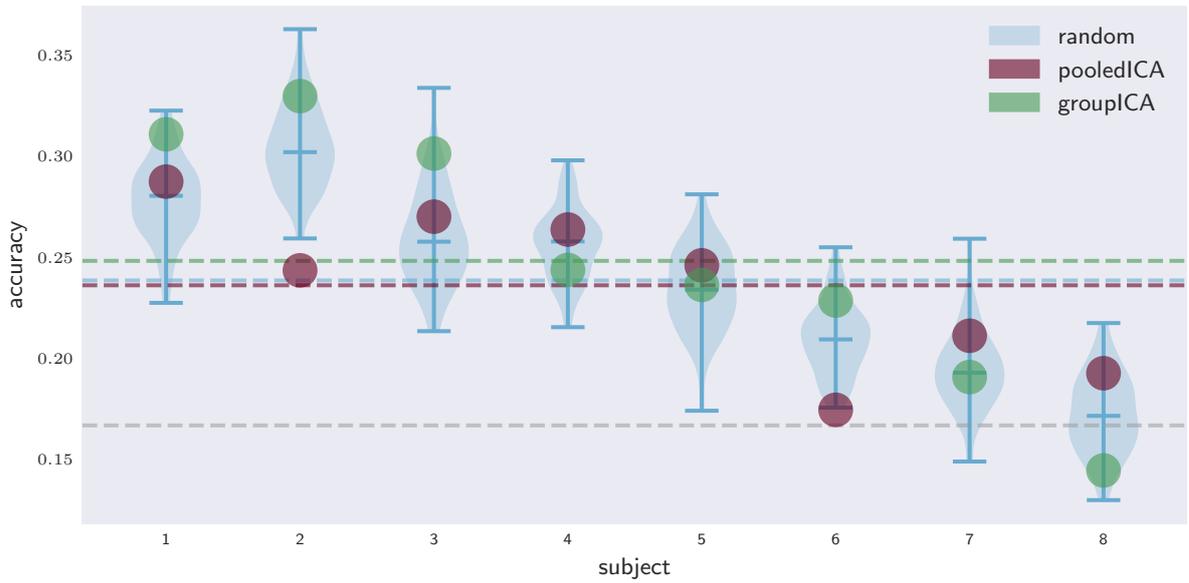


Figure 15. Results as in Figure 6 (on CovertAttention data set) using a shrinkage linear discriminant analysis instead of a linear support vector classifier.