

ASSISTED DICTIONARY LEARNING FOR FMRI DATA ANALYSIS

Manuel Morante Moreno^{1,2}, Yannis Kopsinis^{2,3}, Eleftherios Kofidis^{4,2}, Christos Chatzichristos^{1,2},
Sergios Theodoridis^{1,2,5}

¹ Dept. of Informatics and Telecommunications, University of Athens (Greece), morante@cti.gr

² Computer Technology Institute & Press “Diophantus” (CTI), Patras (Greece), chatzichris@cti.gr

³ LIBRA MLI Ltd, Edinburgh (UK), kopsinis@ieee.org

⁴ Dept. of Statistics and Insurance Science, University of Piraeus, Piraeus (Greece), kofidis@unipi.gr

⁵ IAASARS, National Observatory of Athens, GR-15236, Penteli (Greece), stheodor@di.uoa.gr

ABSTRACT

Extracting information from functional magnetic resonance images (fMRI) has been a major area of research for more than two decades. The goal of this work is to present a new method for the analysis of fMRI data sets, that is capable to incorporate a priori available information, via an efficient optimization framework. Tests on synthetic data sets demonstrate significant performance gains over existing methods of this kind.

Index Terms— fMRI Data Analysis, Dictionary Learning, Blind Source Separation

1. INTRODUCTION

Functional magnetic resonance imaging (fMRI) is a powerful non-invasive technique suitable for providing important information concerning the brain activity. Studying the different areas in the brain that correspond to important tasks such as vision, perception, recognition, etc., constitutes a major open area of research, demanding robust and high precision techniques for the analysis of fMRI data [1], [2], [3], [4].

Such data are generated as a sequence of 3D images of the brain, which are successively acquired along time. Each one of these images is formed by the concatenation of elementary cubes, called *voxels*. Accordingly, the signal measured at each voxel reflects the degree of activity in a certain brain spot. Each 3D image is unfolded into a large row vector, $\mathbf{x} = [x_1, x_2, \dots, x_N] \in \mathbb{R}^N$, where N is the total number of voxels. Then, all such data vectors are concatenated together to form the data matrix, $\mathbf{X} \in \mathbb{R}^{T \times N}$, where T is the total number of successive acquisition time instants.

In the brain, a number of different functions/processes take place simultaneously; the obtained measurements consist of a mixture of various activation signals referred to as *sources*. The aim of fMRI data analysis is to unmix those sources in order to reveal both their activation patterns as well as the corresponding activated brain areas.

From a mathematical point of view, the source unmixing task can be described as a problem of factorization of the data matrix, i.e.,

$$\mathbf{X} \approx \mathbf{D}\mathbf{S}, \quad (1)$$

where $\mathbf{D} \in \mathbb{R}^{T \times K}$ is a matrix, whose columns represent the activation patterns or time courses (also called atoms), associated with each one of the sources, $\mathbf{S} \in \mathbb{R}^{K \times N}$ is the matrix whose rows model the brain areas activated by the corresponding sources, and K is the number of sources, whose value is set by the user. The rows of the matrix \mathbf{S} are usually referred to as spatial maps.

fMRI essentially measures the changes in the level of oxygen in blood caused by the neural activity, which yields an indirect measure of the latter. More specifically, the observed/measured signal results from the convolution of the true activations with the, so called, Hemodynamic Response Function (HRF). The HRF varies across different persons as well as across different brain areas of the same person [5].

A widely used tool in fMRI analysis is the General Linear Model (GLM), which relies on the assumed form of the HRF in order to construct the matrix \mathbf{D} in Eq. (1). In particular, the specific design of each experiment allows to make a guess of the true time instances, where the activations are expected to appear. Adopting a functional form for HRF and convolving it with the expected activation sequence, the time course corresponding to the specific task can be estimated and considered known. Hereafter, such estimated time courses are referred to as *task-related* time courses.

Alternatively, one might use a blind source separation (BSS) approach, which can simultaneously estimate \mathbf{D} and \mathbf{S} without having a resort to any assumptions regarding HRF. To this end, different assumptions concerning either statistical or structural properties of the involved matrices are adopted. Namely, Independent Component Analysis (ICA) [6], [7], [8], which has been widely used in the fMRI unmixing problem, assumes independence among the sources, whereas Dictionary Learning (DL)-based techniques [9], which have been gaining more attention recently, exploit the fact that the matrix \mathbf{S} is expected to be sparse. This is true, since the brain can be considered as a *sparse system*; each task/function produces an activation pattern which appears localized in specific regions [10].

Recently, a method called *Supervised Dictionary Learning* (SDL) [11], which allows the incorporation of information related to the HRF in a BSS framework, was presented, leading to enhanced results. However, both GLM and SDL suffer from the same shortcoming; that is, a sufficiently accurate assumption about the functional form of the HRF needs to be made.

The research leading to these results has received funding from the European Union's H2020 Framework Programme (H2020-MSCA-ITN-2014) under grant agreement n 642685 MacSeNet.

In this paper, a new DL method is proposed, which, although it exploits the benefits of incorporating some *a priori* knowledge concerning the HRF, allows for substantial tolerance against inaccurate choices of its respective form.

2. ASSISTED DL FOR fMRI DATA ANALYSIS

2.1. Supervised Dictionary Learning

The starting point in the formulation of the SDL lies in the splitting of the main dictionary in two parts:

$$\mathbf{D} = [\mathbf{\Delta}, \mathbf{D}_F] \in \mathbb{R}^{T \times K}, \quad (2)$$

where the first part, $\mathbf{\Delta} \in \mathbb{R}^{T \times M}$, is constrained to contain the imposed task-related time courses and is considered fixed. In contrast, the second part, $\mathbf{D}_F \in \mathbb{R}^{T \times (K-M)}$, is the variable one to be estimated via DL optimizing arguments.

The result is still a DL scheme but it incorporates a number of specific time courses; it turns out that the reported results lead to an enhanced performance, compared to those obtained via a standard DL technique. Nevertheless, this approach still inherits the same major drawback associated with GLM. That is, the constrained dictionary atoms (columns of the $\mathbf{\Delta}$ matrix) only help if the *a priori* imposed information is sufficiently accurate. If the imposed time courses are shifted or miss-modelled, their contribution can have a detrimental effect, leading to wrong results.

2.2. Atom-Assisted DL

In this paper, an alternative approach is presented, which provides a more relaxed way of incorporating the *a priori* adopted forms of the time courses. The main idea is to consider that the atoms of the constrained part are not necessarily equal to the *a priori* selected ones. Instead, a looser constraint is employed, embedded in the optimization process. Thus, the strong *equality* requirement is relaxed by a looser *similarity* distance-measuring norm constraint.

Thus, if part of the *a priori* information is not accurate enough, since the constrained atoms are not considered fixed any more, the method is free to readjust them, with respect to the information that resides in the data, in an optimal way. It turns out that such an approach robustifies the procedure against the major drawback associated with the HRF-based methods: the miss-modelling.

The starting point is, again, to split the dictionary:

$$\mathbf{D} = [\mathbf{D}_C, \mathbf{D}_F] \in \mathbb{R}^{T \times K}. \quad (3)$$

In contrast to the SDL approach, the constrained part, $\mathbf{D}_C \in \mathbb{R}^{T \times M}$, is not considered fixed any more, instead, it can vary in line with the constrained optimization cost.

The optimization task, adopted here, is formulated as:

$$(\hat{\mathbf{D}}, \hat{\mathbf{S}}) = \underset{\mathbf{D}, \mathbf{S}}{\operatorname{argmin}} \|\mathbf{X} - \mathbf{D}\mathbf{S}\|_F^2 + \lambda \|\mathbf{S}\|_{1,1} \text{ s.t. } \mathbf{D} \in \mathfrak{D} \quad (4)$$

where $\|\mathbf{S}\|_{1,1} = \sum_{i=1}^K \sum_{j=1}^N |s_{ij}|$ is the sparsity-promoting term over the coefficient matrix, $\|\cdot\|_F$ stands for the Frobenius

norm and \mathfrak{D} is an admissible set of dictionaries. In this case, \mathfrak{D} comprises the dictionaries sharing the following property:

$$\mathfrak{D} = \left\{ \mathbf{D} \in \mathbb{R}^{T \times K} : \begin{array}{l} \|\mathbf{d}_i - \boldsymbol{\delta}_i\|_2^2 \leq c_\delta, \forall i \in [1, M] \subset \mathbb{N} \\ \|\mathbf{d}_i\|_2^2 \leq c_d, \forall i \in [M+1, K] \subset \mathbb{N} \end{array} \right\}, \quad (5)$$

where \mathbb{N} is the set of natural numbers, $\|\cdot\|_2$ denotes the Euclidean norm, \mathbf{d}_i is the i^{th} column of the dictionary \mathbf{D} and $\boldsymbol{\delta}_i$ is the i^{th} *a priori* selected task-related time course. The constant c_δ is a user-defined parameter which controls the *degree of similarity* between the constrained atoms and the imposed time courses. Note that the first term of the objective function above is invariant to scale transformations [12], and hence, in order to prevent degenerate solutions, the remaining dictionary atoms are constrained to have a bounded norm no larger than a prefixed parameter c_d .

2.3. Optimization Method

In order to solve the previous optimization task, the Majorization-Minimization (MM) principle [13] is adopted. No doubt, any other relevant optimization method can be mobilized, and its most appropriate choice is currently under study. Although the adoption of the MM method does not necessarily involve a Lagrangian relaxation, this approach is followed here for simplicity. Thus, the equivalent optimization task, via the corresponding Lagrangian formulation of the minimization problem in Eq. (4), is formulated as

$$(\hat{\mathbf{D}}, \hat{\mathbf{S}}) = \underset{\mathbf{D}, \mathbf{S}}{\operatorname{argmin}} \phi_{\lambda, \gamma}(\mathbf{D}, \mathbf{S}), \quad \text{where} \quad (6)$$

$$\phi_{\lambda, \gamma}(\mathbf{D}, \mathbf{S}) = \|\mathbf{X} - \mathbf{D}\mathbf{S}\|_F^2 + \lambda \|\mathbf{S}\|_{1,1} + \mathcal{P}_\gamma(\mathbf{D}). \quad (7)$$

$\mathcal{P}_\gamma(\mathbf{D})$ depends on the dictionary and is defined as

$$\begin{aligned} \mathcal{P}_\gamma(\mathbf{D}) = & \sum_{i=1}^M \gamma_i \left[(\mathbf{d}_i - \boldsymbol{\delta}_i)^T (\mathbf{d}_i - \boldsymbol{\delta}_i) - c_\delta \right] + \\ & + \sum_{i=M+1}^K \gamma_i (\mathbf{d}_i^T \mathbf{d}_i - c_d), \end{aligned} \quad (8)$$

where the introduced parameters, γ_i , $i = 1, 2, \dots, K$ correspond to the associated K Lagrangian multipliers.

Eq. (8) can be compactly expressed as:

$$\mathcal{P}_\Gamma(\mathbf{D}) = \operatorname{tr} \left[\mathbf{\Gamma} (\mathbf{D} - \mathbf{\Delta}\mathbf{M})^T (\mathbf{D} - \mathbf{\Delta}\mathbf{M}) - \mathbf{C} \right], \quad (9)$$

where $\mathbf{M} \in \mathbb{R}^{M \times K}$ is a rectangular diagonal matrix with diagonal entries equal to one, $\mathbf{\Gamma} = \operatorname{diag}(\gamma_1, \gamma_2, \dots, \gamma_K)$ and \mathbf{C} is the diagonal matrix with the corresponding parameters c_δ and c_d on its diagonal. Accordingly, the cost function (7) is rewritten as:

$$\phi_{\lambda, \Gamma}(\mathbf{D}, \mathbf{S}) = \|\mathbf{X} - \mathbf{D}\mathbf{S}\|_F^2 + \lambda \|\mathbf{S}\|_{1,1} + \mathcal{P}_\Gamma(\mathbf{D}). \quad (10)$$

2.4. The Algorithm

The optimization with respect to \mathbf{D} and \mathbf{S} is a challenging one and is greatly simplified by adopting a two-step alternating

minimization iterative procedure. In particular, starting from arbitrary estimates, $\mathbf{D}_{(0)}$ and $\mathbf{S}_{(0)}$, the algorithm comprises the following steps:

$$\text{Step I} \quad \min_{\mathbf{S}} \phi_{\lambda, \Gamma}(\mathbf{D}, \mathbf{S}) \quad \text{given } \mathbf{D}, \quad (11)$$

$$\text{Step II} \quad \min_{\mathbf{D}} \phi_{\lambda, \Gamma}(\mathbf{D}, \mathbf{S}) \quad \text{given } \mathbf{S}. \quad (12)$$

Following the MM scheme, for each step, the objective function is replaced by a surrogate one, which majorizes it and is easier to be iteratively minimized compared to the original one. The surrogate function is not unique, but it has to satisfy specific conditions, e.g., [13].

2.4.1. Step I: Coefficient Update

At the t^{th} step of the alternating minimization of Eq. (11), the objective function is minimized with respect to \mathbf{S} keeping \mathbf{D} fixed at its currently available estimate, $\mathbf{D} = \mathbf{D}_{(t)}$. This minimization is also achieved in an iterative way and through the introduction of a surrogate function. Starting the iterations from the currently available estimate, $\mathbf{S}^{[0]} = \mathbf{S}_{(t)}$, the estimate, $\mathbf{S}^{[n]}$, at the n^{th} iteration, is obtained in terms of the previous estimate, $\mathbf{S}^{[n-1]}$, by minimizing the following surrogate function [13],

$$\psi_{\lambda}(\mathbf{S}, \mathbf{S}^{[n-1]}) = \phi_{\lambda, \Gamma}(\mathbf{D}, \mathbf{S}) + \pi_S(\mathbf{S}, \mathbf{S}^{[n-1]}), \quad (13)$$

where

$\pi_S(\mathbf{S}, \mathbf{S}^{[n-1]}) := c_S \|\mathbf{S} - \mathbf{S}^{[n-1]}\|_F^2 - \|\mathbf{D}\mathbf{S} - \mathbf{D}\mathbf{S}^{[n-1]}\|_F^2$ and $c_S > \|\mathbf{D}^T \mathbf{D}\|_2$ is a constant with $\|\cdot\|_2$ now standing for the spectral norm. Thus, two different iterations run in a nested form; for each iteration with respect to (t) , there is an (inner) iteration with respect to $[n]$.

Let $\mathbf{A} := \frac{1}{c_S} (\mathbf{D}^T \mathbf{X} + (c_S \mathbf{I}_K - \mathbf{D}^T \mathbf{D}) \mathbf{S}^{[n-1]})$. It can be shown that the optimum value of the surrogate function above is found by shrinking the elements in \mathbf{A} , that is,

$$\mathbf{S}^{[n]} = \mathcal{S}_{\lambda}(\mathbf{A}), \quad \text{where} \quad (14)$$

$$\mathcal{S}_{\lambda}(\mathbf{A}) : s_{ij} = \begin{cases} a_{ij} - \frac{\lambda}{2} \text{sign}(a_{ij}) & \text{if } |a_{ij}| > \frac{\lambda}{2} \\ 0 & \text{otherwise} \end{cases}. \quad (15)$$

The iterative update continues until a stopping criterion is satisfied. The pseudocode for this coefficient update is presented in *Algorithm 1*.

2.4.2. Step II: Dictionary Update

In the second step of the alternating minimization, the objective function is minimized with respect to \mathbf{D} , keeping \mathbf{S} fixed at its currently available estimate, $\mathbf{S} = \mathbf{S}_{(t+1)}$. A majorization rationale is also used for this step as well.

To this end, an appropriate surrogate function is introduced given by

$$\psi_{\Gamma}(\mathbf{D}, \mathbf{R}) = \phi_{\lambda, \Gamma}(\mathbf{D}, \mathbf{S}) + \pi_D(\mathbf{D}, \mathbf{R}), \quad (16)$$

where

$$\pi_D(\mathbf{D}, \mathbf{R}) = c_D \|\mathbf{D} - \mathbf{R}\|_F^2 - \|\mathbf{D}\mathbf{S} - \mathbf{D}\mathbf{R}\|_F^2,$$

Algorithm 1 - Step I (Coefficient Update)

```

1: Initialization:  $c_S > \|\mathbf{D}^T \mathbf{D}\|_2$ ,  $\mathbf{S}^{[0]} = \mathbf{S}_{(t)}$ ,  $n = 0$ 
2: repeat
3:    $n = n + 1$ 
4:    $\mathbf{A} = \frac{1}{c_S} (\mathbf{D}^T \mathbf{X} + (c_S \mathbf{I}_K - \mathbf{D}^T \mathbf{D}) \mathbf{S}^{[n-1]})$ 
5:    $\mathbf{S}^{[n]} = \mathcal{S}_{\lambda}(\mathbf{A})$ 
6: until stop criterion is satisfied*
7: output:  $\mathbf{S}_{(t+1)} = \mathbf{S}^{[n]}$ 

```

Algorithm 2 - Step II (Dictionary Update)

```

1: Initialization:  $c_D > \|\mathbf{S}^T \mathbf{S}\|_2$ ,  $\mathbf{D}^{[0]} = \mathbf{D}_{(t)}$ ,  $n = 0$ 
2: repeat
3:    $n = n + 1$ 
4:    $\mathbf{B} = \frac{1}{c_D} (\mathbf{X}\mathbf{S}^T + \mathbf{R} (c_D \mathbf{I}_K - \mathbf{S}\mathbf{S}^T))$ 
5:    $\mathbf{D}^{[n]} = \mathcal{U}(\mathbf{B})$ 
6: until stop criterion is satisfied*
7: output:  $\mathbf{D}_{(t+1)} = \mathbf{D}^{[n]}$ 

```

* In this paper, a fixed number of iterations is used.

$c_D > \|\mathbf{S}^T \mathbf{S}\|_2$ is a constant and $\mathbf{R} = \mathbf{D}^{[n-1]}$ is the estimate of the dictionary of the previous step.

Minimizing Eq. (16) with respect to \mathbf{D} takes place also iteratively, starting from $\mathbf{D}^{[0]} = \mathbf{D}_{(t)}$. The optimum value of the surrogate function is found at the point of zero gradient:

$$\nabla_{\mathbf{D}} \psi_{\Gamma} = -2\mathbf{X}\mathbf{S}^T + 2(\mathbf{D} - \Delta\mathbf{M})\Gamma + 2c_D(\mathbf{D} - \mathbf{R}) + 2\mathbf{R}\mathbf{S}\mathbf{S}^T. \quad (17)$$

Setting the derivative above equal to zero, solving for \mathbf{D} and setting γ_i to values that satisfy the Karush-Kuhn-Tucker (KKT) conditions (details are omitted due to lack of space), a two-step procedure for the dictionary update results, following arguments similar to those in [13]. Namely, an intermediate quantity \mathbf{B} is first defined and computed as $\mathbf{B} := \frac{1}{c_D} (\mathbf{X}\mathbf{S}^T + \mathbf{R} (c_D \mathbf{I} - \mathbf{S}\mathbf{S}^T))$. Thus, the estimates of the update dictionary atom can be summarized using the operator $\mathbf{D}^{[n]} = \mathcal{U}(\mathbf{B})$ where $\mathcal{U}(\mathbf{B}) : \mathbf{u}_i \rightarrow \mathbf{d}_i^{[n]}$ is given by:

$$\mathbf{d}_i^{[n]} = \begin{cases} \text{for } i \in [1, M], \quad \begin{cases} \mathbf{b}_i & \text{if } \|\mathbf{b}_i - \delta_i\|_2^2 \leq c_{\delta} \\ \frac{c_{\delta}(\mathbf{b}_i - \delta_i)}{\|\mathbf{b}_i - \delta_i\|_2^2} + \delta_i & \text{otherwise} \end{cases} \\ \text{for } i \in [M+1, K], \quad \begin{cases} \mathbf{b}_i & \text{if } \|\mathbf{b}_i\|_2^2 \leq c_d \\ \frac{c_d}{\|\mathbf{b}_i\|_2^2} \mathbf{b}_i & \text{otherwise} \end{cases} \end{cases}, \quad (18)$$

It is not difficult to show that after this update, the KKT conditions for all dictionary columns hold. Moreover, the set of all admissible dictionaries, \mathcal{D} , constitutes a convex non-empty set. It can be shown that this fact guarantees that the proposed algorithm converges for random initialization. Due to the space limitations imposed by a conference paper, details are omitted. The pseudo-code for this dictionary update is presented in *Algorithm 2*. Furthermore, the MATLAB code for this method can be freely downloaded from <https://github.com/MorCTI/Assisted-DL.git>.

3. PERFORMANCE EVALUATION

The aim of this section is twofold. First, to demonstrate the advantages of incorporating external information about task-

related time courses. Second, to compare the sensitivity of the proposed scheme with that of SDL, in cases where the imposed time courses deviate from the true ones.

The data set used is synthetic and generated with SimTB¹ [19]. In order to make the data more realistic, the sources 3, 4, 5, 7, 8 of the data set in [14], which represent machine artifacts, are also added. The data set used can be downloaded from (<https://github.com/MorCTI/Assisted-DL.git>). In Fig. 1, as an example, 6 among the 20 sources used in total are depicted. Note that the current performance evaluation cannot be realized based on real fMRI data, since in such a case the ground truth is not known.

With respect to the SDL and atom-assisted DL methods, the larger the number of time courses which are imposed as constraints in the algorithm, better is the performance observed due to the fact that a larger amount of information is provided. Therefore, in order to make things harder, in the evaluation tests that follow, only one task-related time course is considered. Moreover, two different miss-modeling cases are examined. In the first one, the task-related time course is a time-shifted version of the true one. The result is shown in Fig. 2a, where the solid and the dashed lines correspond to the atom-assisted DL and the SDL, respectively. The horizontal axis represents the time shifting of the imposed task-related time course in relation to the true one, expressed in seconds. The vertical axis shows $1 - R^2$, with R being the correlation coefficient between the estimated and the true source. It is apparent that the proposed scheme offers enhanced robustness allowing time discrepancies up to 4 seconds (2 seconds in each direction) without any performance degradation. If some performance loss is allowed, 6 seconds of time shift are well tolerated.

In the second miss-modeling scenario, shown in Fig. 2b, the imposed time course results from the convolution of the experimental task event with an HRF which is different from the true one. For the construction of the different tested HRFs, the canonical HRF model is adopted [20]. In order to perform this study, the free parameters of the canonical HRF model are gradually modified leading to HRFs with a successively narrower shape compared to the true HRF. In particular, the horizontal axis shows the squared correlation coefficient, R_{HRF}^2 , between the true HRF and the modified HRF of the corresponding imposed time course. Again, the vertical axis shows $1 - R^2$, with R being the correlation coefficient between the estimated and the true source. Once again, it is

¹SimTB simulator is a free MATLAB toolbox available for download in (<http://mialab.mrn.org/software>), which has been lately adopted in a number of fMRI data analysis studies, e.g. [15],[16],[17],[18].

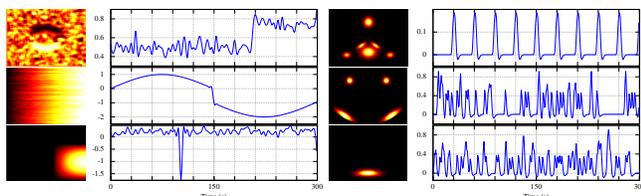


Fig. 1. Selection of different simulated sources. In the first column, Gaussian, subgaussian and supergaussian artefacts are plotted from the artificial data set in [14]. In the second column, three other simulated yet realistic brain sources are shown. The first one corresponds to the source of interest.

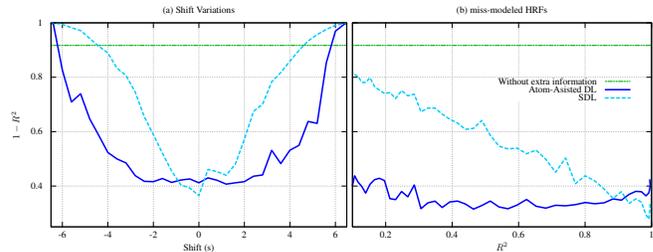


Fig. 2. Squared correlation coefficient between the estimated source and the true one for the two miss-modelling experiments.

observed that the proposed method is insensitive to large deviations between the provided information and the true one.

In both cases, the dot-dashed curve corresponds to the fully blind approach, i.e., when no information regarding the task-related time course is provided. Obviously, the fully blind approach fails to estimate the signal of interest. This happens since in the experimental setup it has been provisioned that the signal of interest a) exhibit significant space overlap with artefact sources and other physical sources and b) have overall energy not higher than of its neighbouring sources. This design generates a hard but realistic experimental scenario, in which other conventional blind source separation methods, such as ICA [7] or k-SVD [9] fail to recover the source of interest. The latter was confirmed with various simulation studies, which will be presented elsewhere due to space limitations.

Note that both in the current and in the next experiment, all curves result from the ensemble average of 20 independent runs. Besides, the majorization optimization approach was also used in the SDL case for the dictionary learning task substituting the online DL optimization, [21], used in the original paper. Note that the SDL is just a particular case of ADL, with $c_\delta = 0$. In any case, it was verified via extensive simulations that the two optimization approaches resulted in similar performance. In the atom-assisted DL case, the parameters, c_δ , c_d , λ were set equal to 0.3, 1, 200, respectively. Moreover, 50 iterations were performed for Algorithms 1 and 2. Finally, 1000 alternating minimization iterations were used in all cases.

4. CONCLUSIONS

In this paper, a new source separation approach for fMRI data analysis is proposed. The method allows for the incorporation of task-related a priori information which leads to vast performance improvements compared to conventional fully blind approaches. Moreover, the proposed method exhibits enhanced robustness against miss-modelling of the imposed extra information.

5. REFERENCES

- [1] Lee Y. et al., “Sparse SPM: Group sparse-dictionary learning in SPM framework for resting-state functional connectivity MRI analysis,” *NeuroImage*, vol. 125, pp. 1032–1045, 2016.
- [2] A. Protopapas et al., “Evaluating cognitive models of visual word recognition using fMRI: Effects of lexical and sublexical variables,” *NeuroImage*, vol. 128, pp. 328–341, 2016.
- [3] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*, Academic Press, 2015.
- [4] M. G. Bright and K. Murphy, “Is fMRI ‘noise’ really a noise? Resting state nuisance regressors remove variance with network structure,” *NeuroImage*, vol. 114, pp. 158–169, 2015.
- [5] G. K. Aguirre, E. Zarahn, and M. D’Esposito, “The variability of human, BOLD hemodynamic responses,” *NeuroImage*, vol. 8, no. 4, pp. 360–369, 1998.
- [6] M. J. McKeown et al., “Analysis of fMRI data by blind separation into independent spatial components,” *Human Brain Mapping*, vol. 6, pp. 160–188, June 1998.
- [7] V. D. Calhoun and T. Adalı, “Unmixing fMRI with independent component analysis,” *IEEE Engineering in Medicine and Biology Magazine*, vol. 25, no. 2, pp. 79–90, 2006.
- [8] L. Griffanti et al., “ICA-based artefact removal and accelerated fMRI acquisition for improved resting state network imaging,” *NeuroImage*, vol. 95, pp. 232–247, 2014.
- [9] Y. Kopsinis, H. Georgiou, and S. Theodoridis, “fMRI unmixing via properly adjusted dictionary learning,” in *22nd European Signal Processing Conference (EU-SIPCO)*, Lisbon, Portugal, Sep. 2014.
- [10] P. A. Valdés-Sosa et al., “Estimating brain functional connectivity with sparse multivariate autoregression,” *Philosophical Transactions of the Royal Society B: Biological Sciences*, vol. 360, no. 1457, pp. 969–981, 2005.
- [11] J. Lu et al., “Supervised dictionary learning for inferring concurrent brain networks,” *IEEE Trans. Medical Imaging*, vol. 34, no. 10, pp. 2036–2045, Oct. 2015.
- [12] Genevera I. Allen, “Sparse and functional principal components analysis,” *arXiv*, vol. 1, pp. 1–21, 2013.
- [13] M. Yaghoobi, T. Blumensath, and M. E. Davies, “Dictionary learning for sparse approximations with the majorization method,” *IEEE Trans. Signal Processing*, vol. 57, no. 6, pp. 2178–2191, Jun. 2009.
- [14] N. Correa, T. Adalı, and Y. Li, “Comparison of blind source separation algorithms for fMRI using a new matlab toolbox: GIFT,” in *30th Int. Conf. Acoustics, Speech, and Signal Processing (ICASSP-2005)*, Philadelphia, PA. IEEE, Mar. 2005, vol. 5, pp. v–401.
- [15] E. A. Allen et al., “Capturing inter-subject variability with group independent component analysis of fMRI data: A simulation study,” *NeuroImage*, vol. 59, no. 4, pp. 4141–4159, 2012.
- [16] S. Ma et al., “Automatic identification of functional clusters in fMRI data using spatial dependence,” *IEEE Trans. Biomedical Engineering*, vol. 58, no. 12, pp. 3406–3417, Dec. 2011.
- [17] S. Ma et al., “Capturing group variability using IVA: A simulation study and graph-theoretical analysis,” in *38th Int. Conf. Acoustics, Speech, and Signal Processing (ICASSP-2013)*, Vancouver, BC, Canada, May. 2013.
- [18] R. Ge et al., “A two-step super-gaussian independent component analysis approach for fMRI data,” *NeuroImage*, vol. 118, pp. 344–358, 2015.
- [19] E. B. Erhardt et al., “SimTB, a simulation toolbox for fMRI data under a model of spatiotemporal separability,” *NeuroImage*, vol. 59, no. 4, pp. 4160–4167, 2012.
- [20] D. A. Handwerker, J. M. Ollinger, and M. D’Esposito, “Variation of BOLD hemodynamic responses across subjects and brain regions and their effects on statistical analyses,” *NeuroImage*, vol. 21, no. 4, pp. 1639–1651, 2004.
- [21] J. Mairal et al., “Online learning for matrix factorization and sparse coding,” *Journal of Machine Learning Research*, vol. 11, pp. 19–60, Jan. 2010.