



Unification of sparse Bayesian learning algorithms for electromagnetic brain imaging with the majorization minimization framework

Ali Hashemi^{a,b,c,d}, Chang Cai^{e,f}, Gitta Kutyniok^{g,h}, Klaus-Robert Müller^{b,i,j,k,*},
Srikantan S. Nagarajan^{e,*}, Stefan Haufe^{a,c,l,m,*}

^a Uncertainty, Inverse Modeling and Machine Learning Group, Technische Universität Berlin, Germany

^b Machine Learning Group, Technische Universität Berlin, Germany

^c Berlin Center for Advanced Neuroimaging (BCAN), Charité – Universitätsmedizin Berlin, Germany

^d Institut für Mathematik, Technische Universität Berlin, Germany

^e Department of Radiology and Biomedical Imaging, University of California, San Francisco, CA, USA

^f National Engineering Research Center for E-Learning, Central China Normal University, China

^g Mathematisches Institut, Ludwig-Maximilians-Universität München, Germany

^h Department of Physics and Technology, University of Tromsø, Norway

ⁱ BIFOLD – Berlin Institute for the Foundations of Learning and Data, Berlin, Germany

^j Department of Brain and Cognitive Engineering, Korea University, Seoul, South Korea

^k Max Planck Institute for Informatics, Saarbrücken, Germany

^l Mathematical Modelling and Data Analysis Department, Physikalisch-Technische Bundesanstalt Braunschweig und Berlin, Germany

^m Bernstein Center for Computational Neuroscience, Berlin, Germany

ARTICLE INFO

Keywords:

Electro-/magnetoencephalography
Brain source imaging
Type I/II Bayesian learning
Non-convex
Majorization-Minimization
Noise learning
Hyperparameter learning

ABSTRACT

Methods for electro- or magnetoencephalography (EEG/MEG) based brain source imaging (BSI) using sparse Bayesian learning (SBL) have been demonstrated to achieve excellent performance in situations with low numbers of distinct active sources, such as event-related designs. This paper extends the theory and practice of SBL in three important ways. First, we reformulate three existing SBL algorithms under the *majorization-minimization* (MM) framework. This unification perspective not only provides a useful theoretical framework for comparing different algorithms in terms of their convergence behavior, but also provides a principled recipe for constructing novel algorithms with specific properties by designing appropriate bounds of the Bayesian marginal likelihood function. Second, building on the MM principle, we propose a novel method called *LowSNR-BSI* that achieves favorable source reconstruction performance in low signal-to-noise-ratio (SNR) settings. Third, precise knowledge of the noise level is a crucial requirement for accurate source reconstruction. Here we present a novel principled technique to accurately learn the noise variance from the data either jointly within the source reconstruction procedure or using one of two proposed cross-validation strategies. Empirically, we could show that the monotonous convergence behavior predicted from MM theory is confirmed in numerical experiments. Using simulations, we further demonstrate the advantage of LowSNR-BSI over conventional SBL in low-SNR regimes, and the advantage of learned noise levels over estimates derived from baseline data. To demonstrate the usefulness of our novel approach, we show neurophysiologically plausible source reconstructions on averaged auditory evoked potential data.

1. Introduction

Electro- and Magnetoencephalography (EEG/MEG) are non-invasive techniques for measuring brain electrical activity with high temporal resolution. As such, both have become indispensable tools in basic neuroscience and clinical neurology. The downside of both techniques, however, is that their sensors are located far away from the neural generators

of the measured brain electrical activity. EEG/MEG measurements are therefore characterized by low spatial resolution and highly overlapping contributions of multiple brain sources in each sensor. The mathematical model of the EEG/MEG sensing procedure can be described by the linear *forward model*

$$\mathbf{Y} = \mathbf{LX} + \mathbf{E}, \quad (1)$$

* Corresponding authors.

E-mail addresses: alhashemi.ee@gmail.com (A. Hashemi), klaus-robot.mueller@tu-berlin.de (K.-R. Müller), Srikantan.Nagarajan@ucsf.edu (S.S. Nagarajan), haufe@tu-berlin.de (S. Haufe).

<https://doi.org/10.1016/j.neuroimage.2021.118309>.

Received 27 March 2021; Received in revised form 17 May 2021; Accepted 23 June 2021

Available online 26 June 2021.

1053-8119/© 2021 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>)

which maps the electrical activity of the brain sources, \mathbf{X} , to the sensor measurements, \mathbf{Y} . The *measurement matrix* $\mathbf{Y} \in \mathbb{R}^{M \times T}$ captures the activity of M sensors attached at different parts of the scalp at T time instants, $\mathbf{y}(t) \in \mathbb{R}^{M \times 1}, t = 1, \dots, T$, while the *source matrix*, $\mathbf{X} \in \mathbb{R}^{N \times T}$, consists of the unknown activity of N brain sources located in the cortical gray matter at the same time instants, $\mathbf{x}(t) \in \mathbb{R}^{N \times 1}, t = 1, \dots, T$. The matrix $\mathbf{E} = [\mathbf{e}(1), \dots, \mathbf{e}(T)] \in \mathbb{R}^{M \times T}$ represents T time instances of independent and identically distributed (i.i.d.) zero mean white Gaussian noise with variance σ^2 , $\mathbf{e}(t) \in \mathbb{R}^{M \times 1} \sim \mathcal{N}(0, \sigma^2 \mathbf{I}_M), t = 1, \dots, T$, which is assumed to be independent of the source activations. The linear forward mapping from \mathbf{X} to \mathbf{Y} is given by the *lead field matrix* $\mathbf{L} \in \mathbb{R}^{M \times N}$, which is here assumed to be known. In practice, \mathbf{L} can be computed using discretization methods such as the Finite Element Method (FEM) for a given head geometry and known electrical conductivities using the quasi-static approximation of Maxwell's equations (Baillet et al., 2001; Gramfort, 2009; Hämäläinen et al., 1993; Huang et al., 2016).

The goal of brain source imaging (BSI) is to infer the underlying brain activity \mathbf{X} from the EEG/MEG measurement \mathbf{Y} given the lead field matrix \mathbf{L} . Unfortunately, this inverse problem is highly ill-posed as the number of sensors is typically much smaller than the number of locations of potential brain sources. Thus, a unique solution cannot be found without introducing further mathematical constraints or penalties, which are often referred to as regularizers. In addition, the leadfield matrix is typically highly ill-conditioned even for small numbers of sensors, introducing numerical instabilities in the inverse estimates.

Interestingly, regularization can also be interpreted in a Bayesian framework, where the regularizer introduces prior knowledge or assumptions about the nature of the true sources into the estimation (Calvetti and Somersalo, 2018; Stuart, 2010). A common assumption is that the number of active brain sources during the execution of a specific mental task is small, i.e., that the spatial distribution of the brain activity is sparse. This assumption can be encoded in various ways. Classical approaches (Matsuura and Okabe, 1995) employ super-Gaussian prior distributions to identify solutions in which most of brain regions are inactive.

In these approaches Maximum-a-Posteriori (MAP) estimation, also termed *Type-I learning*, is used. Later work (Wipf et al., 2010) has shown that hierarchical Bayesian models achieve better reconstructions of sparse brain signals by employing a separate Gaussian prior for each brain location. The variances at each location are treated as unknown (hyper-) parameters, which are estimated jointly with the source activity. This approach is called Sparse Bayesian Learning (SBL), Type-II Maximum-Likelihood (Type-II ML) estimation or simply *Type-II learning* (Mika et al., 2001; Tipping, 2001; Wipf and Rao, 2004).

Type-II learning generally leads to non-convex objective functions, which are non-trivial to optimize. A number of iterative algorithms have been proposed (Mika et al., 2001; Tipping, 2001; Wipf and Nagarajan, 2009; 2010; Wipf et al., 2010; 2011), which, due to employing distinct parameter update rules, differ in their convergence guarantees, rates and overall computational complexity. Being derived using vastly different mathematical concepts such as fixed point theory and expectation-maximization (EM), it has, however, so far been difficult to explain the observed commonalities and differences, advantages and disadvantages of Type-II methods in absence of a common theoretical framework, even if the properties of individual algorithms have been extensively studied (Wipf and Nagarajan, 2009).

The primary contribution of this paper is to introduce *Majorization-Minimization (MM)* (Hunter and Lange, 2004; Sun et al., 2017, and references therein) as a flexible algorithmic framework within which different SBL approaches can be theoretically analyzed. Briefly, MM is a family of iterative algorithms to optimize general non-linear objective functions. In a minimization setting, MM replaces the original cost function in each iteration by an upper bound, or majorization function, whose minimum is usually easy to find. The objective value at the minimum is then used to construct the bound for the following iteration, and the procedure is repeated until a local minimum of the objective

is reached. Notably, MM algorithms are popular in many disciplines in which Type-II learning problems arise, such as, e.g., telecommunications (Haghighatshoar and Caire, 2017; Khalilarai et al., 2020; Oguz-Ekim et al., 2011; Prasad et al., 2015; Shen et al., 2019) and finance (Benidis et al., 2018; Feng et al., 2016). The concept of MM is, however, rarely explicitly referenced in EEG/MEG brain source imaging, even though it has been used implicitly (Bekhti et al., 2018; Hashemi and Haufe, 2018). We demonstrate here that three popular SBL variants, denoted as *EM*, *MacKay*, and *convex-bounding based* SBL, can be cast as majorization-minimization methods employing different types of upper bounds on the marginal likelihood. This view as variants of MM helps explain, among other things, the guaranteed convergence of these algorithms to a local minimum. The characteristics of the chosen bounds determine the reconstruction performance and convergence rates of the resulting algorithms. The MM framework additionally offers a principled way of constructing new SBL algorithms for specific purposes by designing appropriate bounds.

Therefore, a second contribution of this paper is the development of a new SBL algorithm, called LowSNR-BSI, that is especially suitable for low signal-to-noise ratio (SNR) regimes. Real-world applications of EEG/MEG brain source imaging are often characterized by low SNR, where the power of unwanted noise sources can be comparable to the power of the signal of interest. This holds in particular for the reconstruction of ongoing as well as induced (non-phase-locked) oscillatory activity, where no averaging can be performed prior to source reconstruction. Current SBL algorithms may suffer from reduced performance in such low-SNR regimes (Cai et al., 2021; Khanna and Murthy, 2017a; Owen et al., 2012). To overcome this limitation, we propose a novel MM algorithm for EEG/MEG source imaging, which employs a bound on the SBL cost function that is particularly suitable for low-SNR regimes.

As a third contribution, this paper discusses principled ways to estimate the sensor noise variance σ^2 , which is assumed to be known in the first part of the paper. Determining the goodness-of-fit of the optimal model, the value of this variable exerts a strong impact on the overall reconstruction (Habermehl et al., 2014). Technically being another model hyperparameter, the noise variance is, however, rarely estimated as part of the model fitting. Instead, it is often determined prior to the model fitting from a baseline recording. This approach can, however, lead to suboptimal results in practice or be even inapplicable, e.g., when resting state data are analyzed. Here we present a number of alternatives to estimate the noise variance in Type-I and Type-II brain source imaging approaches. Building on work by (Wipf and Rao, 2007), we derive an analytic update rule, which enables the adaptive estimation of the noise variance within various SBL schemes. Moreover, we propose two novel cross-validation (CV) schemes from the machine learning field to determine the noise variance parameter.

We conduct extensive ground-truth simulations in which we compare LowSNR-BSI with popular source reconstruction schemes including existing SBL variants, and in which we systematically study the impact of different strategies to estimate the noise level σ^2 from the data.

The outline of the paper is as follows: In Section 2, a comprehensive review of Type-II BSI methods is presented. In Section 3, we unify the Type-II methods described in Section 2 within the MM framework, and in Section 4, we derive LowSNR-BSI algorithm within the same framework. Section 5 introduces numerous principled ways for estimating the sensor noise variance. Simulation studies, real data analysis, and discussions are presented in Sections 6, 7, and 8, respectively. Finally, Section 9 concludes the paper.

2. Bayesian learning

The ill-posed nature of the EEG/MEG inverse problem can be overcome by assuming a prior distribution $p(\mathbf{X})$ for the source activity. The posterior distribution of the sources after observing the data \mathbf{Y} , $p(\mathbf{X}|\mathbf{Y})$,

is given by Bayes' rule:

$$p(\mathbf{X}|\mathbf{Y}) = \frac{p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})}{\int p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})d\mathbf{X}}, \quad (2)$$

where the conditional probability $p(\mathbf{Y}|\mathbf{X})$ in the numerator denotes the *likelihood*, while the term in the denominator, $\int p(\mathbf{Y}|\mathbf{X})p(\mathbf{X})d\mathbf{X} = p(\mathbf{Y})$, is referred to as *model evidence* or *marginal likelihood*. However, note that the posterior is often not analytically tractable, as evaluating the integral in the model evidence is intractable for many choices of prior distributions and likelihoods.

Remark 1. Priors fulfill the same practical purpose as regularizers even if they are motivated from a different perspective, e.g., the Bayesian formalism, in this paper. Besides, we regard the Bayesian perspective as a helpful technical vehicle to inspire and generate flexible priors for shaping more plausible solutions.

2.1. Type-I Bayesian learning

As the model evidence in Eq. (2) only acts as a scalar normalization for the posterior, its evaluation can be avoided if one is only interested in the most probable source configuration \mathbf{X} rather than the full posterior distribution. This point estimate is known as the maximum-a-posteriori (MAP) estimate:

$$\mathbf{X}^{\text{MAP}} := \arg \max_{\mathbf{X}} \underbrace{p(\mathbf{Y}|\mathbf{X})}_{\text{likelihood}} \underbrace{p(\mathbf{X})}_{\text{prior}}. \quad (3)$$

Assuming i.i.d. Gaussian sensor noise, the likelihood reads:

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{t=1}^T p(\mathbf{y}(t)|\mathbf{x}(t)) = \prod_{t=1}^T \mathcal{N}(\mathbf{L}\mathbf{x}(t), \sigma^2\mathbf{I}), \quad (4)$$

and the resulting MAP estimate (3) is given by

$$\begin{aligned} \mathbf{X}^{\text{MAP}} &:= \arg \max_{\mathbf{X}} \left[\prod_{t=1}^T \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\mathbf{x}(t)\|_2^2\right) \right] p(\mathbf{X}) \\ &= \arg \min_{\mathbf{X}} \left[\frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L}\mathbf{x}(t)\|_2^2 \right] + \sigma^2 \mathcal{R}^1(\mathbf{X}) = \arg \min_{\mathbf{X}} \mathcal{L}^1(\mathbf{X}), \end{aligned} \quad (5)$$

where $\mathcal{R}^1(\mathbf{X}) = \log(p(\mathbf{X}))$ and $\mathcal{L}^1(\mathbf{X})$ denotes the Bayesian Type-I learning (MAP) objective function.

Note that this expression can be interpreted as a trade-off between two optimization goals, where the first (log-likelihood) term in (5) penalizes model errors using a quadratic loss function and the second (log-prior) term penalizes deviations of the solution from the assumed spatial or temporal properties of the brain sources encoded in $\mathcal{R}^1(\mathbf{X})$. The trade-off between these two optimization goals is defined by the ratio of the noise variance σ^2 and the variance of the prior distribution. As the latter is hardly known in practice, a *regularization parameter* $\lambda \propto \sigma^2$ subsuming both variables is introduced, which can be tuned to adjust the relative importance of both penalties in the optimization.

Several existing algorithms are characterized by different choices of a prior. For instance, choosing a Gaussian prior distribution leads to the classical minimum-norm estimate (Hämäläinen and Ilmoniemi, 1994; Pascual-Marqui, 2007; Pascual-Marqui et al., 1994), which also goes by the names ℓ_2^2 -norm (or Tikhonov) regularization and “ridge regression” in the statistics and machine learning literature. The choice of a Laplace prior leads to the minimum-current estimate (Matsuura and Okabe, 1995), which is also known as ℓ_1 -norm regularization or “LASSO” regression. Besides, hierarchical Bayesian priors with automatic depth weighting have been used to infer brain activity from EEG/MEG data (Calvetti et al., 2019). More complex priors have been also used to incorporate anatomical information of the sources (Dale and Sereno, 1993; Pascual-Marqui et al., 2002; Trujillo-Barreto et al., 2004) or to encode assumptions on the spatial, temporal and/or spectral structure of the sources. Respective methods include FOCUSS (Gorodnitsky et al., 1995),

S-FLEX (Haufe et al., 2008; 2011), MxNE (Gramfort et al., 2012), ir-MxNE (Strohmeier et al., 2016), TF-MxNE (Gramfort et al., 2013), irTF-MxNE (Strohmeier et al., 2015), and STOUT (Castaño-Candamil et al., 2015), which all enforce sparsity in different domains such as Gabor frames or cortical patches through appropriate norm constraints.

2.2. Type-II Bayesian learning

While in the MAP approach the prior distribution is fixed, it is sometimes desirable to consider entire families of distributions $p(\mathbf{X}|\gamma)$ parameterized by a set of hyper-parameters γ . These hyper-parameters can be learned from the data along with the model parameters using a hierarchical empirical Bayesian approach (Mika et al., 2001; Tipping, 2001; Wipf and Rao, 2004). In this maximum-likelihood Type-II (ML-II, or simply Type-II) approach, γ is estimated through the maximum-likelihood principle:

$$\gamma^{\text{II}} := \arg \max_{\gamma} p(\mathbf{Y}|\gamma) = \arg \max_{\gamma} \int p(\mathbf{Y}|\mathbf{X}, \gamma) p(\mathbf{X}|\gamma) d\mathbf{X}. \quad (6)$$

Computation of the conditional density $p(\mathbf{Y}|\gamma)$ is formally achieved by integrating over all possible source distributions \mathbf{X} for any given choice of γ . The maximizer of Eq. (6) then determines a data-driven prior distribution $p(\mathbf{X}|\gamma^{\text{II}})$. Plugged into the MAP estimation framework Eq. (3), this gives rise to the Type-II source estimate \mathbf{X}^{II} .

As the conditional density $p(\mathbf{Y}|\gamma)$ for a given γ is identical to the model evidence in Eq. (2), this approach also goes by the name evidence maximization (Wipf and Rao, 2007; Wipf et al., 2011). Concrete instantiations of this approach have further been introduced under the names *sparse Bayesian learning* (SBL) (Tipping, 2001) or *automatic relevance determination* (ARD) (Tipping, 2000), *kernel Fisher discriminant* (KFD) (Mika et al., 2001), *variational Bayes* (VB) (Seeger and Wipf, 2010; Wipf and Nagarajan, 2009) and iteratively-reweighted MAP estimation (Gorodnitsky et al., 1995; Wipf and Nagarajan, 2010). Interested readers are referred to (Wu et al., 2016) for a comprehensive survey on Bayesian machine learning techniques for EEG/MEG signals. To distinguish all these Type-II variants from classical ML and MAP approaches not involving hyperparameter learning, the latter are also referred to as Type-I approaches.

Remark 2. The marginal likelihood formulation in Type-II Bayesian learning, Eq. (6), enables estimation of flexible priors with many parameters from data. This stands in contrast to the use of classical cross-validation techniques to learn hyperparameters of regularizers, which works for very few parameters only (in most cases only a single scalar regularization constant).

2.3. Sparse Bayesian learning and Champagne

A Type-II estimation framework with particular relevance for EEG/MEG source imaging is SBL. In this framework, the N modeled brain sources are assumed to follow independent univariate Gaussian distributions with zero mean and distinct unknown variances γ_n : $x_n(t) \sim \mathcal{N}(0, \gamma_n)$, $n = 1, \dots, N$. In the SBL solution, the majority of variances is zero, thus effectively inducing spatial sparsity of the corresponding source activities. Such sparse solutions are physiologically plausible in task-based analyses, where only a fraction of the brain's macroscopic structures is expected to be consistently engaged. This consideration has led (Wipf and Rao, 2004) to propose the *Champagne* algorithm for brain source imaging, which is rooted in the concept of SBL. Compared to Type-I approaches achieving sparsity through ℓ_1 -norm minimization, Champagne has shown significant performance improvement with respect to EEG/MEG source localization (Owen et al., 2012; Wipf et al., 2010).

Just as most existing approaches, Champagne makes the simplifying assumption of statistical independence between time samples. This leads

to the following expression for the distribution of the sources:

$$p(\mathbf{X}|\boldsymbol{\gamma}) = \prod_{t=1}^T p(\mathbf{x}(t)|\boldsymbol{\gamma}) = \prod_{t=1}^T \mathcal{N}(\mathbf{0}, \boldsymbol{\Gamma}), \quad (7)$$

where $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$ and $\boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma})$. Note that, in task-based analyses, the noise variance σ^2 can be estimated from a baseline (resting state) recording. In the first part of this paper, it is, therefore, assumed to be known.

Remark 3. Electrophysiological data are known to possess a complex intrinsic autocorrelation structure. Here, we consider priors that make the simplifying assumption of independence between time samples, which is consistent with most existing works in the field (Gramfort et al., 2012; Hämäläinen and Ilmoniemi, 1994; Haufe et al., 2008; Mat-suura and Okabe, 1995; Pascual-Marqui et al., 1994). Importantly, using such simplifying priors generally does not prevent the resulting inverse solutions to have time structure. Nevertheless, priors modeling the known properties of the latent variables more accurately might lead to better reconstructions especially in low-sample regimes. Preliminary work shows that priors modeling temporal structure with autoregressive models can indeed improve the reconstruction of autocorrelated source (Hashemi and Haufe, 2018).

The parameters of the SBL model are the unknown sources as well as their variances. As computation of the integral in Eq. (6) is infeasible, Champagne considers an approximation, where the variances $\gamma_n, n = 1, \dots, N$, are optimized based on the current estimates of the sources in an alternating iterative process. Given an initial estimate of the variances, the posterior distribution of the sources is a Gaussian of the form (Wipf et al., 2010), (Sekihara and Nagarajan, 2015, Chapter 4)

$$p(\mathbf{X}|\mathbf{Y}, \boldsymbol{\gamma}) = \prod_{t=1}^T \mathcal{N}(\bar{\mathbf{x}}(t), \boldsymbol{\Sigma}_{\mathbf{x}}), \quad \text{where} \quad (8)$$

$$\bar{\mathbf{x}}(t) = \boldsymbol{\Gamma} \mathbf{L}^\top (\boldsymbol{\Sigma}_{\mathbf{y}})^{-1} \mathbf{y}(t) \quad (9)$$

$$\boldsymbol{\Sigma}_{\mathbf{x}} = \boldsymbol{\Gamma} - \boldsymbol{\Gamma} \mathbf{L}^\top (\boldsymbol{\Sigma}_{\mathbf{y}})^{-1} \mathbf{L} \boldsymbol{\Gamma} \quad (10)$$

$$\boldsymbol{\Sigma}_{\mathbf{y}} = \sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top. \quad (11)$$

The estimated posterior parameters $\bar{\mathbf{x}}(t)$ and $\boldsymbol{\Sigma}_{\mathbf{x}}$ are then in turn used to update the estimate of the variances $\gamma_n, n = 1, \dots, N$ as the minimizer of the negative log of the marginal likelihood $p(\mathbf{Y}|\boldsymbol{\gamma})$, which is given by Wipf et al. (2010):

$$\mathcal{L}^{\text{II}}(\boldsymbol{\gamma}) = -\log p(\mathbf{Y}|\boldsymbol{\gamma}) = \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} \mathbf{y}(t) + \log |\boldsymbol{\Sigma}_{\mathbf{y}}|, \quad (12)$$

where $|\cdot|$ denotes the determinant of a matrix. This process is repeated until convergence. Given the final solution of the hyperparameter $\boldsymbol{\gamma}^{\text{II}}$, the point estimate \mathbf{x}^{II} of the source activity is obtained from the posterior mean of the estimated source distribution: $\mathbf{x}^{\text{II}}(t) = \bar{\mathbf{x}}(t)$. Note that given the definition of the empirical sample covariance matrix as $\mathbf{C}_{\mathbf{y}} = \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t) \mathbf{y}(t)^\top$, the term $\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_{\mathbf{y}}^{-1} \mathbf{y}(t)$ in Eq. (12) can be rewritten as $\text{tr}(\mathbf{C}_{\mathbf{y}} \boldsymbol{\Sigma}_{\mathbf{y}}^{-1})$, so that Eq. (12) becomes (Wipf et al., 2010, Section II)

$$\mathcal{L}^{\text{II}}(\boldsymbol{\gamma}) = \text{tr}(\mathbf{C}_{\mathbf{y}} \boldsymbol{\Sigma}_{\mathbf{y}}^{-1}) + \log |\boldsymbol{\Sigma}_{\mathbf{y}}|. \quad (13)$$

Note that, in this form, the loss function Eq. (13) bears an interesting similarity to the *log-determinant (log-det) Bregman divergence* in information geometry (James and Stein, 1992). This perspective on Type-II loss function enables a common viewpoint for Type-I and Type-II methods.

By invoking mathematical tools based on Legendre-Fenchel duality theory, the cost function Eq. (12) can be formulated equivalently as another cost function, $\mathcal{L}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma})$, whose optimizers, $\{\boldsymbol{\gamma}^*, \mathbf{X}^*\}$, are derived

by performing a joint minimization over \mathbf{X} and $\boldsymbol{\gamma}$ (Wipf et al., 2011, see also Section II-B), Bauschke and Combettes (2017); Rockafellar (1970):

$$\boldsymbol{\gamma}^*, \mathbf{X}^* = \arg \min_{\boldsymbol{\gamma} \geq 0, \mathbf{X} \geq 0} \mathcal{L}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma}), \quad \text{where}$$

$$\begin{aligned} \mathcal{L}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma}) &= \frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L} \mathbf{x}(t)\|_2^2 + \sigma^2 \mathcal{R}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma}) \\ \mathcal{R}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma}) &= \frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{x_n(t)^2}{\gamma_n} + \log |\boldsymbol{\Sigma}_{\mathbf{y}}|, \end{aligned} \quad (14)$$

where $\mathcal{R}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma})$ denotes a regularizer that depends on the data, $\mathbf{x}(t)$, and where $x_n(t)$ denotes the activity of source n at time instant t . Then, as each source $x_n(t)$ is also a function of γ_n according to Eq. (9), the term $\frac{x_n(t)^2}{\gamma_n}$ goes to zero when $\gamma_n \rightarrow 0$.

Remark 4. In contrast to standard MAP estimation, the effective priors obtained within our hierarchical Bayesian framework, e.g., $\mathcal{R}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma})$ in Eq. (14), are not fixed. They depend on parameters that can be tuned and learned from the data; thus, Type-II priors have the ability and flexibility to capture the actual properties of the observed real data.

We will use the formulation in Eq. (14) to derive alternative optimization schemes for Champagne in Sections 2.3.2 and 2.3.3.

2.3.1. EM Champagne

As the cost function Eq. (12) is non-convex in $\boldsymbol{\gamma}$, the quality of the obtained solution depends substantially on the properties of the employed numerical optimization algorithm. Crucially, algorithms might not only differ with respect to their convergence properties but may also lead to different solutions representing distinct local minima of Eq. (12). The first algorithm for minimizing Eq. (12) has been introduced by Wipf and Nagarajan (2009) and is an application of the expectation-maximization (EM) formalism (Dempster et al., 1977). As can be shown, Eqs. (9)–(11) correspond to the expectation (E) step of the EM algorithms with respect to the posterior distribution $p(\mathbf{X}|\mathbf{Y}, \boldsymbol{\gamma})$. The maximization (M) step of the EM formalism with respect to $\boldsymbol{\gamma}$ then leads to the update rule

$$\gamma_n^{k+1} := [\boldsymbol{\Sigma}_{\mathbf{x}}^k]_{n,n} + \frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2 \quad \text{for } n = 1, \dots, N. \quad (15)$$

Final estimates of both parameters are obtained by iterating the updates (9)–(11) and (15) until convergence. The resulting algorithm is known as the EM variant of the Champagne algorithm (Sekihara and Nagarajan, 2015, Chapter 4) (Wipf and Nagarajan, 2009) in the field of brain source imaging.

2.3.2. Convex-bounding based Champagne

As the EM algorithm outlined above has been shown to have slow convergence speed, alternative minimization strategies have been proposed. Two such variants, a convex-approximation based approach and the so-called MacKay update, have been proposed in Wipf and Nagarajan (2009) and further practically investigated in Owen et al. (2012). Considering that the log-determinant in Eq. (14) is concave, the convex-bounding based variant of Champagne constructs a linear upper bound based on the concave conjugate of $\log |\sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top|$, defined as $w^*(\mathbf{z})$,

$$\log |\sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top| = \log |\sigma^2 \mathbf{I} + \mathbf{L} \text{diag}(\boldsymbol{\gamma}) \mathbf{L}^\top| = \min_{\mathbf{z} > 0} \mathbf{z}^\top \boldsymbol{\gamma} - w^*(\mathbf{z}). \quad (16)$$

With this upper bound, and for a fixed value of $\boldsymbol{\gamma}$, the auxiliary variable \mathbf{z} can be derived as the tangent hyperplane of the $\log |\boldsymbol{\Sigma}_{\mathbf{y}}|$:

$$\mathbf{z} = \nabla_{\boldsymbol{\gamma}} \log |\sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top|.$$

Note that the concave conjugate is obtained as a result of applying Legendre-Fenchel duality theory (see, e.g., Bauschke and Combettes (2017); Rockafellar (1970)) on the concave function $\log |\sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top|$ as follows: $w^*(\mathbf{z}) = \inf_{\boldsymbol{\gamma} > 0} [\boldsymbol{\gamma}^\top \mathbf{z} - w(\boldsymbol{\gamma})]$, where $w(\boldsymbol{\gamma}) = \log |\sigma^2 \mathbf{I} + \mathbf{L} \boldsymbol{\Gamma} \mathbf{L}^\top|$ denotes our target concave function.

By inserting Eq. (16) instead of $\log|\Sigma_y|$ into Eq. (14), the non-convex penalty function Eq. (14) is replaced by the convex function

$$\mathcal{R}_{\text{conv}}^{\text{II-x}}(\mathbf{X}, \boldsymbol{\gamma}) = \min_{\boldsymbol{\gamma} \geq 0, \mathbf{z} > 0} \left[\frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{\bar{x}_n(t)^2}{\gamma_n} \right] + \mathbf{z}^T \boldsymbol{\gamma} - w^*(\mathbf{z})$$

in each step of the optimization. The final estimates of \mathbf{X} , $\boldsymbol{\gamma}$ and \mathbf{z} are obtained by iterating between following update rules until convergence:

$$\mathbf{z}_n^k = \mathbf{L}_n^T (\Sigma_y^k)^{-1} \mathbf{L}_n, n = 1, \dots, N \quad (17)$$

$$\bar{\mathbf{x}}^k(t) = \mathbf{L} \mathbf{L}^T (\Sigma_y^k)^{-1} \mathbf{y}(t) \quad (18)$$

$$\gamma_n^{k+1} = \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2}{z_n^k}}, n = 1, \dots, N. \quad (19)$$

Here, \mathbf{L}_n in (17) denotes the n -th column of the lead field matrix.

2.3.3. MacKay update for Champagne

The MacKay update proposed in (Wipf and Nagarajan, 2009, Section III.A-2) can be derived in a similar fashion as the convex-bounding based update using different auxiliary functions and variables. By defining new variables $\kappa_n := \log(\gamma_n)$ for $n = 1, \dots, N$, the non-convex term $\log|\sigma^2 \mathbf{I} + \mathbf{L} \text{diag}(\boldsymbol{\gamma}) \mathbf{L}^T|$ in Eq. (16) can be written as:

$$\begin{aligned} \log|\sigma^2 \mathbf{I} + \mathbf{L} \text{diag}(\boldsymbol{\gamma}) \mathbf{L}^T| &= \log \left| \sigma^2 \mathbf{I} + \sum_{n=1}^N \gamma_n \mathbf{L}_n^T \mathbf{L}_n \right| \\ &= \log \left| \sigma^2 \mathbf{I} + \sum_{n=1}^N \exp(\kappa_n) \mathbf{L}_n^T \mathbf{L}_n \right|. \end{aligned}$$

Then, one can introduce another surrogate function (Wipf and Nagarajan, 2009, Appendix-B)

$$\log \left| \sigma^2 \mathbf{I} + \sum_{n=1}^N \exp(\kappa_n) \mathbf{L}_n^T \mathbf{L}_n \right| = \max_{\mathbf{z} > 0} \mathbf{z}^T \log(\boldsymbol{\gamma}) - h^*(\mathbf{z}) \quad (20)$$

for the $\log|\sigma^2 \mathbf{I} + \mathbf{L} \mathbf{L}^T|$, where $h^*(\mathbf{z})$ denotes the convex conjugate of $\log|\sigma^2 \mathbf{I} + \sum_{n=1}^N \exp(\kappa_n) \mathbf{L}_n^T \mathbf{L}_n|$ in contrast to the concave conjugate counterpart, $w^*(\mathbf{z})$ used in Eq. (16). Substituting (20) into Eq. (14) leads to a so-called *min-max optimization program* for optimizing the non-convex penalty function $\mathcal{R}^{\text{II-x}}(\mathbf{X})$, which alternates between minimizations over $\boldsymbol{\gamma}$ and maximizations of the bound in (20):

$$\mathcal{R}_{\text{conv}}^{\text{MacKay}}(\mathbf{X}, \boldsymbol{\gamma}) = \min_{\boldsymbol{\gamma} \geq 0} \max_{\mathbf{z} > 0} \left[\frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{\bar{x}_n(t)^2}{\gamma_n} \right] + \mathbf{z}^T \log(\boldsymbol{\gamma}) - h^*(\mathbf{z}). \quad (21)$$

Let γ_n^k denote the value of γ_n in the k -th iteration. Inserting $\boldsymbol{\gamma}^k$ into Eq. (21) and minimizing with respect to $\boldsymbol{\gamma}_n^k$ requires that the derivatives

$$\frac{\partial}{\partial \gamma_n^k} \left[\frac{1}{T} \sum_{t=1}^T \frac{\bar{x}_n^k(t)^2}{\gamma_n^k} + \mathbf{z}^T \log(\boldsymbol{\gamma}_n^k) - h^*(\mathbf{z}) \right] = 0,$$

for $n = 1, \dots, N$, vanish. The resulting function is then maximized with respect to \mathbf{z} (Wipf and Nagarajan, 2009, Appendix-B), which leads to the so-called MacKay update for optimizing Eq. (14) (Wipf and Nagarajan, 2009, Section A-2):

$$\begin{aligned} \gamma_n^{k+1} &:= \left[\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2 \right] \left(\gamma_n^k \mathbf{L}_n^T (\Sigma_y^k)^{-1} \mathbf{L}_n \right)^{-1} \\ &= \left[\frac{1}{T} \sum_{t=1}^T (\gamma_n^k) \beta_n^k(t)^2 \right] \left(\gamma_n^k \mathbf{L}_n^T (\Sigma_y^k)^{-1} \mathbf{L}_n \right)^{-1} \\ &= \gamma_n^k \left[\frac{1}{T} \sum_{t=1}^T \beta_n^k(t)^2 \right] \left(\mathbf{L}_n^T (\Sigma_y^k)^{-1} \mathbf{L}_n \right)^{-1}, \text{ for } n = 1, \dots, N, \quad (22) \end{aligned}$$

where $\beta_n^k(t)$ is defined as follows: $\beta_n^k(t) := \mathbf{L}_n^T (\Sigma_y^k)^{-1} \mathbf{y}(t)$ for $n = 1, \dots, N$.

3. Unification of sparse Bayesian learning algorithms with the majorization-minimization (MM) framework

In this section, we first briefly review theoretical concepts behind the MM algorithmic framework (Hunter and Lange, 2004; Jacobson and Fessler, 2007; Razaviyayn et al., 2013; Wu et al., 2010). Then, we formally characterize Champagne variants as MM algorithms by suggesting upper bounds on the cost function Eq. (14) that, when employed within the MM framework, yield the same update rules as the original algorithms. The first three rows of Table 1 list the update rules and mathematical formalism used in this section.

3.1. Majorization-Minimization

Majorization-minimization is a promising strategy for solving general non-linear optimization programs. Compared to other popular optimization paradigms such as (quasi)-Newton methods, MM algorithms enjoy guaranteed convergence to a stationary point (Sun et al., 2017). The MM class covers a broad range of common optimization algorithms such as *proximal methods* and *convex-concave procedures (CCCP)* (Sun et al., 2017, Section IV), Lipp and Boyd (2016); Yuille and Rangarajan (2003). While such algorithms have been applied in various contexts, such as non-negative matrix factorization (Févotte, 2011) and massive MIMO systems for wireless communication (Haghighatshoar and Caire, 2017; Khalilsarai et al., 2020), their advantages have so far rarely been made explicit in the context of brain source imaging (Bekhti et al., 2018; Hashemi and Haufe, 2018; Luessi et al., 2013).

We define an original optimization problem with the objective of minimizing a continuous function $f(\mathbf{u})$ within a closed convex set $\mathcal{U} \subset \mathbb{R}^n$:

$$\min_{\mathbf{u}} f(\mathbf{u}) \quad \text{subject to } \mathbf{u} \in \mathcal{U}. \quad (23)$$

Then, the idea of MM can be summarized as follows. First, construct a continuous *surrogate function* $g(\mathbf{u}|\mathbf{u}^k)$ that upper-bounds, or *majorizes*, the original function $f(\mathbf{u})$ and coincides with $f(\mathbf{u})$ at a given point \mathbf{u}^k :

$$\begin{aligned} \text{[A1]} \quad & g(\mathbf{u}^k|\mathbf{u}^k) = f(\mathbf{u}^k) \quad \forall \mathbf{u}^k \in \mathcal{U} \\ \text{[A2]} \quad & g(\mathbf{u}|\mathbf{u}^k) \geq f(\mathbf{u}) \quad \forall \mathbf{u}, \mathbf{u}^k \in \mathcal{U}. \end{aligned}$$

Second, starting from an initial value \mathbf{u}^0 , generate a sequence of feasible points $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^k, \mathbf{u}^{k+1}$ as solutions of a series of successive simple optimization problems, where

$$\text{[A3]} \quad \mathbf{u}^{k+1} := \arg \min_{\mathbf{u} \in \mathcal{U}} g(\mathbf{u}|\mathbf{u}^k).$$

Note that the performance of MM algorithms heavily depends on the choice of a suitable surrogate function, which should, on one hand, faithfully reflect the behavior of the original non-convex function Eq. (23) while, on the other hand, be easy to minimize.

Definition 1. Any algorithm fulfilling conditions [A1]–[A3] is called a Majorization Minimization (MM) algorithm.

Corollary 1. An MM algorithm has a *descending trend property*, whereby the value of the cost function f decreases in each iteration: $f(\mathbf{u}^{k+1}) \leq f(\mathbf{u}^k)$.

Proof. The proof is included in Appendix B. \square

While Corollary 1 guarantees a descending trend, convergence requires additional assumptions on particular properties of f and g (Jacobson and Fessler, 2007; Razaviyayn et al., 2013). For the smooth functions considered in this paper, we require that the derivatives of the original and surrogate functions coincide at \mathbf{u}^k :

$$\text{[A4]} \quad \nabla g(\mathbf{u}^k|\mathbf{u}^k) = \nabla f(\mathbf{u}^k) \quad \forall \mathbf{u}^k \in \mathcal{U}.$$

Then, the following, stronger, theorem holds.

Theorem 1. For an MM algorithm that additionally satisfies [A4], every limit point of the sequence of minimizers generated through [A3] is a stationary point of the original optimization problem Eq. (23).

Table 1

This table summarizes the update rules presented in Section 2.3 and their corresponding MM upper-bounds that will be utilized in Sections 3 and 4.

| | Update Rule | Mathematical Formalism Used | Inequality Used |
|-----------------|---|---------------------------------------|---------------------|
| EM | $\gamma_n^{k+1} := [\Sigma_n^k]_{n,n} + \frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^2$ | Expectation-Maximization Formalism | Jensen's Inequality |
| Convex Bounding | $\gamma_n^{k+1} := \sqrt{\left[\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^2 \right] \left(\mathbf{L}_n^\top (\Sigma_y^k)^{-1} \mathbf{L}_n \right)^{-1/2}}$ | Concave Conjugate | Taylor Expansion |
| MacKay | $\gamma_n^{k+1} := \left[\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^2 \right] \left(\gamma_n^k \mathbf{L}_n^\top (\Sigma_y^k)^{-1} \mathbf{L}_n \right)^{-1}$ | Change of Variable + Convex Conjugate | Taylor Expansion |
| LowSNR-BSI | $\gamma_n^{k+1} := \sqrt{\left[\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^2 \right] \left(\mathbf{L}_n^\top \mathbf{L}_n \right)^{-1/2}}$ | MM Principle in Low-SNR Setting | Taylor Expansion |

Proof. A detailed proof can be found in (Razaviyayn et al., 2013, Theorem 1). \square

Note that since we are working with smooth functions, conditions [A1]–[A4] are sufficient to prove convergence to a stationary point according to Theorem 1 (see Hunter and Lange (2004); Razaviyayn et al. (2013); Wu et al. (2010) and Dempster et al. (1977); Wu (1983)) for proofs of the convergence behaviour of other MM algorithms such as expectation maximization.

Remark 5. Corollary 1 implies that if a surrogate function is constructed to fulfill conditions [A1] and [A2], and if the next feasible point of the algorithm is always assigned as the minimizer of the surrogate function based on [A3], the resulting MM algorithm decreases $f(\mathbf{u})$ in each step. Although a weaker condition than [A3], i.e., $g(\mathbf{u}^{k+1}|\mathbf{u}^k) \leq g(\mathbf{u}^k|\mathbf{u}^k)$, is sufficient for a descending trend, we only consider MM algorithms in this paper; thus, condition [A3] is a crucial requirement. As we have shown in Theorem 1, [A3] is further required to prove guaranteed convergence of an MM algorithm.

We now show that three algorithms that have been proposed for solving the SBL cost function Eq. (12) can all be cast as instances of the MM framework invoking different majorization functions on $\mathcal{R}^{\text{II-x}}(\mathbf{X})$. For the convex-bounding based approach as well as the algorithm using MacKay updates, the full set of conditions [A1]–[A4] in Theorem 1 are proven. Due to the considerations made above, we, however, only prove Corollary 1 for the EM-based Champagne algorithm.

3.1.1. EM update as MM

It is known that the EM algorithm is a special case of MM framework using Jensen's inequality to construct the surrogate function (Sun et al., 2017; Wu et al., 2010). Here, we work out the specific surrogate function for the SBL cost function Eq. (12) (i.e., the negative log marginal likelihood).

As Wipf and Nagarajan have shown (Wipf and Nagarajan, 2009, Section III.A-1), the EM algorithm for Type-II problems consists of the following two parts: For the E-step, the posterior $p(\mathbf{X}|\mathbf{Y}, \gamma^k)$ is obtained given the value of γ at k -th iteration, γ^k . The M-step then solves:

$$\begin{aligned} \gamma^{k+1} &:= \arg \min_{\gamma} E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \left[-\log p(\mathbf{Y}, \mathbf{X}|\gamma) \right], \text{ where} \\ -\log p(\mathbf{Y}, \mathbf{X}|\gamma) &= \frac{T}{2} \log |\Gamma| + \frac{1}{2} \sum_{t=1}^T \bar{\mathbf{x}}(t)^\top \Gamma^{-1} \bar{\mathbf{x}}(t) \\ &+ \frac{T}{2} \log |\sigma^2 \mathbf{I}| + \sum_{t=1}^T \frac{1}{2\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}(t)\|_2^2, \end{aligned} \quad (24)$$

which leads to the update rule in Eq. (15).

Proposition 1. The EM based Champagne algorithm is an MM algorithm fulfilling Corollary 1, where the negative log-likelihood loss, $-\log p(\mathbf{Y}|\gamma)$, is majorized by the following surrogate function

$$\begin{aligned} \mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k) &= \frac{T}{2} \log |\Gamma| + E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \left[\frac{1}{2} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \Gamma^{-1} \bar{\mathbf{x}}^k(t) \right] \\ &+ \frac{T}{2} \log |\sigma^2 \mathbf{I}| + E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \left[\sum_{t=1}^T \frac{1}{2\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2 \right] \\ &+ E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\mathbf{Y}, \gamma^k). \end{aligned} \quad (25)$$

Proof. A detailed proof can be found in Appendix C. \square

Note that the EM algorithm is also equivalent to the restricted maximum likelihood (ReML) (Friston et al., 2002) and dynamic statistical parametric mapping (dSPM) approaches (Dale et al., 2000) for solving the sparse EEG/MEG inverse problem, which, thereby, can also be interpreted as instances of minimization-majorization.

3.1.2. Convex-bounding based approach as MM

We start by recalling the non-convex penalty $\mathcal{R}^{\text{II-x}}(\mathbf{X}, \gamma)$ as defined in Eq. (14):

$$\mathcal{R}^{\text{II-x}}(\mathbf{X}, \gamma) = \frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{x_n(t)^2}{\gamma_n} + \log |\Sigma_y|.$$

By setting $\mathbf{x} = \bar{\mathbf{x}}^k$ to the value obtained by the convex-bounding based method in the k -th iteration, the following holds:

Proposition 2. The convex-bounding based Champagne algorithm is an MM algorithm fulfilling Theorem 1, where $\mathcal{R}^{\text{II-x}}(\mathbf{X}, \gamma)$ is majorized by the following surrogate function:

$$\begin{aligned} \mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k) &= \left[\frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{\bar{x}_n^k(t)^2}{\gamma_n} \right] + \log |\Sigma_y^k| + \text{tr} \left[\left(\Sigma_y^k \right)^{-1} \Sigma_y \right] \\ &- \text{tr} \left[\left(\Sigma_y^k \right)^{-1} \Sigma_y^k \right]. \end{aligned} \quad (26)$$

Proof. A detailed proof is provided in Appendix D. \square

3.1.3. MacKay update as MM

Similar to convex-bounding, we can show that the MacKay updates for Champagne can be viewed as an MM algorithm.

Proposition 3. The Champagne variant employing MacKay updates is an MM algorithm fulfilling Theorem 1, where $\mathcal{R}^{\text{II-x}}(\mathbf{X}, \gamma)$ is majorized by $\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k)$.

Proof. The proof is similar to that of Proposition 2 and provided in Appendix E. \square

To summarize this section, we have shown that three popular strategies for solving the SBL problem in Eq. (12), namely the EM, the MacKay, and the convex bounding based approaches, can be characterized as MM algorithms. Importantly, this perspective provides a common framework for comparing different Champagne algorithms. For example, we can derive and compare certain characteristics of Champagne algorithms directly based on the properties of the majorization functions they employ. Conversely, it is also possible to design specific majorization functions that are optimal in a specific sense, leading to new source reconstruction algorithms.

4. LowSNR-brain source imaging (LowSNR-BSI)

Here, we assume a low-SNR regime, as it is common in BSI applications. SNR is defined in sensor space as signal power, $\mathbb{E}\{||\mathbf{L}\mathbf{x}(t)||^2\}$, divided by noise power, σ^2 : $\text{SNR} = \frac{\mathbb{E}\{||\mathbf{L}\mathbf{x}(t)||^2\}}{\sigma^2}$, and can be expressed in dB scale as $\text{SNR}_{\text{dB}} = 10\log_{10}(\text{SNR})$. In many practical applications, we are interested in solving the BSI problem for $\text{SNR}_{\text{dB}} \leq 0$; that is, when the noise power is comparable to the power of the signal or even larger. Although the algorithms presented in Sections 3.1.1–3.1.3 achieve satisfactory performance in terms of computational complexity, their reconstruction performance degrades significantly in low-SNR regimes. This behavior has been theoretically shown in (Khanna and Murthy, 2017a, Section VI-E) and has also been confirmed in several simulation studies (Cai et al., 2021; Owen et al., 2012).

In order to improve the performance of SBL in low-SNR settings, we propose a novel MM algorithm by constructing a surrogate function for Eq. (12) specifically for this setting. Based on (Haghighatshoar and Caire, 2017), we propose the following convex surrogate function:

$$\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\boldsymbol{\gamma}|\boldsymbol{\gamma}^k) = \text{tr}(\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top) + \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t). \quad (27)$$

The following proposition is based on results in (Haghighatshoar and Caire, 2017).

Proposition 4. *The surrogate function Eq. (27) majorizes the Type-II loss function Eq. (12) and results in an MM algorithm that fulfills Theorem 1. For $\text{SNR} \rightarrow 0$, Eq. (12) converges to Eq. (27):*

$$\mathcal{L}^{\text{II}}(\boldsymbol{\gamma}) = \mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\boldsymbol{\gamma}|\boldsymbol{\gamma}^k) + \mathcal{O}(\text{SNR}). \quad (28)$$

Proof. A detailed proof of this result is presented in Appendix F. The main idea is to first normalize the sensor and source covariance matrices by σ^2 and then consider the eigenvalue decomposition of $\text{tr}(\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top)$ as $\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top = \mathbf{U}\mathbf{P}\mathbf{U}^\top$ with $\mathbf{P} = \text{diag}(p_1, \dots, p_M)$. These two steps result in the following equality: $\log|\boldsymbol{\Sigma}_y| = \log|\mathbf{I} + \mathbf{U}\mathbf{P}\mathbf{U}^\top|$. Finally, the proof is completed by leveraging the concavity of the $\log(\cdot)$ function and using a Taylor expansion around the eigenvalues of $\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top$, i.e., p_i , for $i = 1, \dots, M$. \square

Note that, as a result of Proposition 4, the behaviour of the non-convex SBL cost function Eq. (12) is more and more well approximated in the vicinity of the current estimate by the proposed surrogate function Eq. (27) as the noise level increases, which sets it apart from existing surrogate functions. Therefore, the proposed bound is particularly suitable in low-SNR regimes.

In contrast to the original SBL cost function Eq. (12), the surrogate function Eq. (27) is convex and has unique minimum that can be found analytically in each iteration of the optimization. To find the optimal value of $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$, we first take the derivative of (27) with respect to each γ_n for $n = 1, \dots, N$, and then set it to zero, which yields the following closed-form solution for $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$:

$$\gamma_n^{k+1} := \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}_n^k(t)^2}{\mathbf{L}_n^\top \mathbf{L}_n}} \text{ for } n = 1, \dots, N. \quad (29)$$

A detailed derivation of Eq. (29) can be found in Appendix G. We call the algorithm obtained by iterating between (9)–(11) and (29) *LowSNR-Brain Source Imaging (LowSNR-BSI)*. In practice, values exactly equal to zero may not be obtained for the γ_n . Therefore, an *active-set* strategy is employed. Given a threshold γ_{thresh} , those variances γ_n for which $\gamma_n < \gamma_{\text{thresh}}$ holds are set to zero in each iteration of the algorithm. Algorithm 1 summarizes the steps of LowSNR-BSI. Table 1 allows for a direct comparison of the LowSNR-BSI update rule (last column) and the corresponding update rules of other Champagne variants derived within the MM framework.

Algorithm 1: LowSNR-BSI algorithm.

Input: The lead field matrix $\mathbf{L} \in \mathbb{R}^{M \times N}$, the measurement vectors $\mathbf{y}(t) \in \mathbb{R}^{M \times 1}$, $t = 1, \dots, T$, and the noise variance σ^2 .

Result: The estimated prior source variances $[\gamma_1, \dots, \gamma_N]^\top$, the posterior mean $\bar{\mathbf{x}}(t)$ and covariance $\boldsymbol{\Sigma}_x$ of the sources.

1 Set a random initial value for $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$ and construct $\boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma})$.

2 Calculate the statistical covariance $\boldsymbol{\Sigma}_y = \sigma^2 \mathbf{I} + \mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top$.

3 Initialize $k \leftarrow 1$

Repeat

4 Calculate the posterior mean as $\bar{\mathbf{x}}^k(t) = \boldsymbol{\Gamma}\mathbf{L}^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t)$.

5 Update γ_n for $n = 1, \dots, N$ based on Eq. (29).

6 Recalculate the active set of brain sources by selecting the values of γ_n that are greater than a pre-defined threshold:

$$\gamma_n > \gamma_{\text{thresh}}, \quad n = 1, \dots, N.$$

7 $k \leftarrow k + 1$

Until stopping condition is satisfied: $||\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k||_2^2 \leq \epsilon$ or $k = k_{\text{max}}$;

8 Calculate the posterior covariance as $\boldsymbol{\Sigma}_x = \boldsymbol{\Gamma} - \boldsymbol{\Gamma}\mathbf{L}^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{L}\boldsymbol{\Gamma}$.

5. Automatic estimation of the noise level

5.1. Adaptive noise learning

It is common practice to estimate the noise variance σ^2 from baseline data prior to solving the EEG/MEG inverse problem (Bijma et al., 2003; Cai et al., 2018; De Munck et al., 2002; Engemann and Gramfort, 2015; Huizenga et al., 2002; Jun et al., 2006; Plis et al., 2006). However, a baseline estimate may not always be available or may not be accurate enough, say, due to inherent non-stationarities in the data/experimental setup. Here, we argue that estimating the noise parameter from the to-be-reconstructed data can significantly improve the reconstruction performance even compared to a baseline estimate. To this end, we here derive data-driven update rules that allow us to tune estimate the noise variance, σ^2 within the source reconstruction procedure using the Champagne and LowSNR-BSI algorithms, where we build on prior work by Mika et al. (2001); Tipping (2001); Wipf and Rao (2007); Wu and Wipf (2012); Zhang, Rao, 2011. Practically we introduce the shortcut $\lambda = \sigma^2$ to underscore that λ is a tunable parameter whose estimate can substantially deviate from the baseline estimate in practice. We then treat λ as another model hyperparameter, similar to the source variances γ_n . Thus, in each step of learning cycles of the Champagne and LowSNR-BSI algorithms, we also minimize the loss function \mathcal{L}^{II} with respect to λ , where the remaining parameters $\boldsymbol{\Gamma}$ and $\boldsymbol{\Sigma}_x$ are fixed to the values obtained in the preceding iteration. This leads to the following theorem:

Theorem 2. *The minimization of $\mathcal{L}^{\text{II}}(\lambda)$ with respect to λ ,*

$$\lambda^* := \arg \min_{\lambda} \mathcal{L}^{\text{II}}(\lambda) = \arg \min_{\lambda} \left(\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) + \log|\boldsymbol{\Sigma}_y| \right),$$

yields the following update rule for λ at the $(k+1)$ -th iteration, assuming $\boldsymbol{\Gamma}^k$ and $\boldsymbol{\Sigma}_x^k$ be fixed values obtained in the (k) -th iteration:

$$\lambda^{k+1} := \frac{\frac{1}{T} \sum_{t=1}^T ||\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)||_2^2}{M - N^k + \text{tr}[(\boldsymbol{\Sigma}_x^k)(\boldsymbol{\Gamma}^k)^{-1}]}, \quad (30)$$

where N^k denotes the number of non-zero voxels identified at iteration k through an active-set strategy.

Proof. *A detailed proof can be found in Appendix H.* \square

As shown in Algorithm (1), our implementation uses an active-set strategy that only selects the non-zero voxels at each iteration based on a threshold. Therefore, at the initial steps of the algorithm, $N^k = N$ since all source variances are initialized randomly. But, when the algorithm

proceeds, the number of non-zero voxels decreases as a result of our active-set strategy, which results in smaller values for N^k .

5.2. Cross-validation strategies

In the previous section, we proposed to estimate the noise variance $\lambda = \sigma^2$ *in-sample* such that the SBL likelihood according to Eq. (12) was maximized, which led to an analytic update rule. As, under our assumption of homoscedastic sensor noise, λ is only a single scalar parameter, it moreover becomes feasible to make use of robust model selection techniques employing the concept of cross-validation (CV), whose aim it is to maximize the *out-of-sample* likelihood (Bishop, 2006; Hastie et al., 2009; Shalev-Shwartz and Ben-David, 2014). To this end, the data are split into two parts. On the so-called *training set*, the model parameters are fitted for a wide range of possible values of λ , which are fixed within each individual optimization. The likelihoods of the fitted models are then evaluated on the hold-out data parts, called the *test sets*. The choice of λ that maximizes the empirical likelihood on the test data is then used as an unbiased estimate of the noise variance. It is well-known from the field of machine learning that cross-validation effectively overcomes the problem of model overfitting in small samples. Here, we introduce two CV strategies employing different ways of splitting the data.

5.2.1. Temporal cross-validation

In temporal CV, the temporal sequence of the data samples is split into k different contiguous blocks (folds) (Blankertz et al., 2011; Lemm et al., 2011). Here, we use $k = 4$. Three folds form the training set, $\mathbf{Y}^{\text{train_temp}} \in \mathbb{R}^{M \times T^{\text{train_temp}}}$, on which we fit the Champagne and LowSNR-BSI models for a range of λ s. On the remaining fold, $\mathbf{Y}^{\text{test_temp}} \in \mathbb{R}^{M \times T^{\text{test_temp}}}$, the Type-II log-likelihood (c.f. Eqs. (12) and (13))

$$\begin{aligned} \mathcal{L}^{\text{II}}(\mathbf{Y}^{\text{train_temp}}, \mathbf{Y}^{\text{test_temp}}) &= \frac{1}{T} \sum_{t=1}^T \mathbf{y}^{\text{test_temp}}(t)^T \boldsymbol{\Sigma}_{\mathbf{y}^{\text{train_temp}}}^{-1} \mathbf{y}^{\text{test_temp}}(t) \\ &\quad + \log |\boldsymbol{\Sigma}_{\mathbf{y}^{\text{train_temp}}}| = \text{tr}(\mathbf{C}_{\mathbf{y}^{\text{test_temp}}} \boldsymbol{\Sigma}_{\mathbf{y}^{\text{train_temp}}}^{-1}) \\ &\quad + \log |\boldsymbol{\Sigma}_{\mathbf{y}^{\text{train_temp}}}| \end{aligned} \quad (31)$$

is then evaluated. Note that in Eq. (31) the model covariance $\boldsymbol{\Sigma}_{\mathbf{y}^{\text{train_temp}}}$ that has been determined on the training data $\mathbf{Y}^{\text{train_temp}}$ is combined with the empirical covariance of the hold-out data $\mathbf{Y}^{\text{test_temp}}$, which were not used during model fitting. Thus, Eq. (31) is the *out-of-sample* Type-II log-likelihood. It has been theoretically shown (Friedman et al., 2008; Khanna and Murthy, 2017a) that the Type-II log-likelihood function is a metric on the second-order information of the sensors closely related to the log-det Bregman divergence (discrepancy) between statistical (model) and empirical covariances (Bregman, 1967; James and Stein, 1992). The choice of λ that minimizes that discrepancy on hold-out data is, therefore, a sensible estimate for the true noise variance. We provide further details on the relation between the SBL likelihood and the log-det Bregman divergence in Appendix A.

5.2.2. Spatial cross-validation

In spatial CV, the data are not split into temporal segments but by dividing the available EEG/MEG sensors into the training and test sets. This variant has been proposed by Habermehl et al. (2014); Haufe et al. (2011). Here, we again use $k = 4$ folds, where we randomly assign 75% of the sensors to the training set, $\mathbf{Y}^{\text{train_spat}} \in \mathbb{R}^{M^{\text{train_spat}} \times T}$, and the remaining 25% to the test set, $\mathbf{Y}^{\text{test_spat}} \in \mathbb{R}^{M^{\text{test_spat}} \times T}$. On the training sensors, Champagne and LowSNR-BSI are fitted using the corresponding portion of the leadfield matrix, $\mathbf{L}^{\text{train_spat}}$, for the same range of λ s as used in temporal CV. The sources, $\mathbf{X}^{\text{train_spat}} \in \mathbb{R}^{N \times T}$, estimated from the fitted models are then mapped back to the sensor space, and the out-of-sample Type-I log-likelihood (c.f. Eq. (5)) is evaluated on the hold-out (test) sensors:

$$\begin{aligned} \mathcal{L}^{\text{I}}(\mathbf{Y}^{\text{train_spat}}, \mathbf{Y}^{\text{test_spat}}) &= \frac{1}{T} \sum_{t=1}^T \left\| \mathbf{y}^{\text{test_spat}}(t) - \mathbf{L}^{\text{test_spat}} \mathbf{x}^{\text{train_spat}}(t) \right\|_2^2 \\ &:= \left\| \mathbf{Y}^{\text{test_spat}} - \mathbf{L}^{\text{test_spat}} \mathbf{X}^{\text{train_spat}} \right\|_F^2. \end{aligned} \quad (32)$$

Note that, while the Type-II log-likelihood has an interpretation as a Bregman divergence between model and empirical covariance matrices, the Type-I log-likelihood is the Frobenius norm or mean-squared error (MSE) $\|\cdot\|_F^2$ of the model residuals, i.e., the average squared Euclidean distance between empirical and modeled observation vectors. Thus, while the Type-II likelihood compares model and observations in terms of their second-order statistics, the Type-I likelihood uses only first-order information. As in temporal CV, the value of λ that minimizes the MSE on the test sensors is selected as the final noise estimate.

6. Simulations

We conducted an extensive set of simulations, in which we compared the reconstruction performance of the proposed LowSNR-BSI algorithm to that of Champagne and two additional widely-used source reconstruction schemes for a range of different SNRs. We also tested impact of the proposed noise learning schemes (adaptive, temporal CV and spatial CV) on the source reconstruction performance compared to estimating the noise level from baseline data.

6.1. Pseudo-EEG signal generation

Forward modeling

Populations of pyramidal neurons in the cortical gray matter are known to be the main drivers of the EEG signal (Nunez et al., 2006). Here, we use a realistic volume conductor model of the human head to model the linear relationship between primary electrical source currents in these populations and the scalp surface potentials captured by EEG electrodes. The New York Head model (Huang et al., 2016) provides a segmentation of an average human head into six different tissue types. In this model, 2004 dipolar current sources were placed evenly on the cortical surface and 58 sensors were placed on the scalp according to the extended 10–20 system (Oostenveld and Praamstra, 2001). In accordance with the predominant orientation of pyramidal neuron assemblies, the orientation of all source currents was fixed to be perpendicular to the cortical surface, so that only scalar source amplitudes needed to be estimated. Finite-element modeling was used to compute the lead field matrix, $\mathbf{L} \in \mathbb{R}^{58 \times 2004}$, which serves as the forward model in our simulations.

Source generation

We simulated a sparse set of $N_0 = 3$ active sources, which were placed at random positions on the cortex. The temporal activity of each source was generated by a univariate linear autoregressive (AR) process, which models the activity at time t as a linear combination of the P past values:

$$x_i(t) = \sum_{p=1}^P a_i(p)x_i(t-p) + \xi_i(t), \text{ for } i = 1, 2, 3.$$

Here, $a_i(p)$ for $i = 1, 2, 3$ are linear AR coefficients, and P is the order of the AR model. The model residuals $\xi_i(\cdot)$ for $i = 1, 2, 3$ are also referred to as the innovation process; their variance determines the stability of the overall AR process. We here assume uncorrelated standard normal distributed innovations, which are independent for all sources. In the following, we use stable AR systems of order $P = 5$.

Noise model

To simulate the electrical neural activity of the underlying brain sources, $T = 20$ data points were sampled from the AR process described

above. Corresponding dipolar current sources were then placed at random locations, yielding sparse source activation vectors $\mathbf{x}(t)$. Source activations $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$ were mapped to the 58 EEG sensors through application of the lead field matrix \mathbf{L} :

$$\mathbf{Y}^{\text{signal}} = \mathbf{L}\mathbf{X} \quad (33)$$

Next, we added Gaussian white noise to the sensor-space signal. To this end, noise was randomly sampled from a standard normal distribution and normalized with respect to its Frobenius norm. A weighted sum of signal and noise contributions then yielded the pseudo-EEG signal

$$\mathbf{Y} = \mathbf{Y}^{\text{signal}} + \alpha \frac{\mathbf{Y}^{\text{noise}}}{\|\mathbf{Y}^{\text{noise}}\|_F}, \quad (34)$$

where α determines the signal-to-noise ratio in sensor space. For a given α , the noise variance is obtained as $\sigma^2 = 1/M \text{tr}[\Sigma_e]$, for $\Sigma_e = \text{Cov}[\alpha \frac{\mathbf{Y}^{\text{noise}}}{\|\mathbf{Y}^{\text{noise}}\|_F}]$, and the SNR (in dB) is calculated as $\text{SNR} = 20 \log_{10}(\|\mathbf{Y}^{\text{signal}}\|_F / \alpha)$. Since our goal is to investigate the effect of noise variance estimation on the performance of the proposed algorithms, we fixed the noise variance in each set of simulations so as to obtain distributions of performance metrics for a number of similar effective SNR values. We conducted four sets of simulations using $\alpha = \{2, 1.5, 1, 0.5\}$, corresponding to average noise variances of $\sigma^2 = \{37.4 \times 10^{-3}, 21.0 \times 10^{-3}, 9.4 \times 10^{-3}, 2.3 \times 10^{-3}\}$ and average SNRs of $\text{SNR} = \{0.33, 2.17, 4.87, 11.40\}$ (dB). Each set of simulations consists of 100 experiments, in which source locations and time series as well as noise realizations were randomly sampled.

In addition to the pseudo-EEG signal, a pseudo baseline measurement containing only noise but no signal was generated. The sole purpose of this measurement was to provide an empirical estimate of the noise variance as a baseline for our joint source reconstruction and noise estimation approaches, which estimate the same quantity from the summed pseudo-EEG signal. To ensure sufficiently precise baseline estimation, 300 noise samples were generated, normalized, and scaled by α as in Eq. (34) for each experiment.

6.2. Source reconstruction

We applied Champagne and LowSNR-BSI to the synthetic datasets described above. The variances of all voxels were initialized randomly by sampling from a standard normal distribution. The optimization programs were terminated either after reaching convergence (defined by a relative change of the Frobenius-norm of the reconstructed sources between subsequent iterations of less than 10^{-8}), or after reaching a maximum of $k_{\text{max}} = 3000$ iterations.

In each experiment, we evaluated the algorithms using 40 predefined choices of the noise variance ranging from $\lambda = 1/3\sigma^2$ to $\lambda = 30\sigma^2$. In addition, λ was estimated from data using the techniques introduced in Section 5. We observed that the variance estimated from baseline data, $\hat{\sigma}^2$ (averaged over all EEG channels) was typically almost identical to the ground-truth value $\lambda = \sigma^2$ used to simulate the data. The reconstruction performance obtained using this value was therefore included in the comparison as a baseline. Performance at baseline noise level was compared to the performance obtained using adaptive learning of the noise using Eq. (30) as well as using spatial or temporal cross-validation. Note that, for temporal CV, we generated $T = 80$ samples, so that we obtained 60 samples in each training set and 20 samples in each test fold. Due to the increased number of training samples, this method, therefore, has an advantage over the remaining ones. For spatial CV, due to the spatial blur introduced by volume conduction, there is a limit on how focal the measured sensor-space electrical potentials or magnetic fields can be, and the signal will usually be distributed over all sensors. Therefore, a setting in which all ‘signal-carrying’ electrodes will end up either in the training or test set is unlikely to occur in practice. Using, for example, $k = 4$ random splits, it is ensured that the training set will typically capture the signal pattern well. The test set in this approach is only used to

evaluate the out-of-sample likelihood on the remaining sensors, while no model fitting needs to take place. Therefore, missing certain aspects of the signal pattern in the test set does not pose a critical problem, especially if multiple splits are conducted.

Remark 6. The fact that real M/EEG data have time structure is acknowledged in our simulation setting by modeling source time courses as AR processes. The resulting samples of the training and test sets thereby become dependent. Technically, this violates the i.i.d. assumption underlying the theory of CV. However, one can argue that training and test sets are de-facto independent since the leakage from one set to another is small compared to the length of the data. In the spatial CV approach, in contrast, the sensors of the training and test sets are strongly dependent on another, because of the spatial blur introduced by volume conduction. Nevertheless, as we observe in Sections 6.4 and 7, spatial CV works very well both in simulations and real data analysis. This observation suggests that the cross-validation approach can work even if the i.i.d. assumption is violated, in line with previous literature (Habermehl et al., 2014; Hastie et al., 2009; Haufe et al., 2011; Kohavi et al., 1995).

In addition to Champagne and LowSNR-BSI, two non-SBL source reconstruction schemes were included for comparison. As an example of a sparse Type-I method based on ℓ_1 -norm minimization, S-FLEX (Haufe et al., 2011) was used. As spatial basis functions, unit impulses were used, so that the resulting estimate was identical to the so-called minimum-current estimate (Matsura and Okabe, 1995). In addition, the eLORETA estimate (Pascual-Marqui, 2007), a smooth inverse solution based on weighted ℓ_2^2 -norm minimization was used. eLORETA was used with 5% regularization, whereas S-FLEX was fitted so that the residual variance was consistent with the ground-truth noise level. Note that the 5% rule is chosen as it gives the best performance across a subset of regularization values ranging between 0.5% to 15%.

6.3. Evaluation metrics

Source reconstruction performance was evaluated according to the following metrics. First, the *earth mover’s distance* (EMD, Haufe et al. (2008); Rubner et al. (2000)) was used to quantify the spatial localization accuracy. The EMD metric measures the cost needed to transform two probability distributions, defined on the same metric domain, into each other. It was applied here to the $N \times 1$ amplitude distributions of the true and estimated sources, which were obtained by taking the voxel-wise ℓ_2 -norm along the time domain. EMD scores were normalized to be in $[0, 1]$. Second, the error in the reconstruction of the source time courses was measured. To this end, Pearson correlation between all pairs of simulated and reconstructed (i.e., those with non-zero activations) sources was measured. Each simulated source was matched to a reconstructed source based on maximum absolute correlation. Time course reconstruction error was then defined as one minus the average of these absolute correlations across sources. Finally, the runtime of the algorithms was measured in seconds (s).

6.4. Results

Fig. 1 shows the EMD (upper row), the time course reconstruction error (middle row) and the negative log-likelihood loss value (lower row) incurred by Champagne and LowSNR-BSI for two SNR settings (SNR = 0.33 dB and SNR = 11.40 dB). Four different schemes of estimating the noise level from data (estimation from baseline data, adaptive learning, spatial CV, and temporal CV) are compared. Note that we found previously that the ground-truth noise variance $\lambda = \sigma^2$ used in the simulation is generally accurately estimated from baseline data, which is referred to as ‘baseline’ in the figure, $\lambda = \hat{\sigma}^2$. Interestingly, however, this baseline is optimal only for LowSNR-BSI, and only with respect to temporal source reconstruction. For Champagne, and with respect to the spatial source reconstruction performance of LowSNR-BSI,

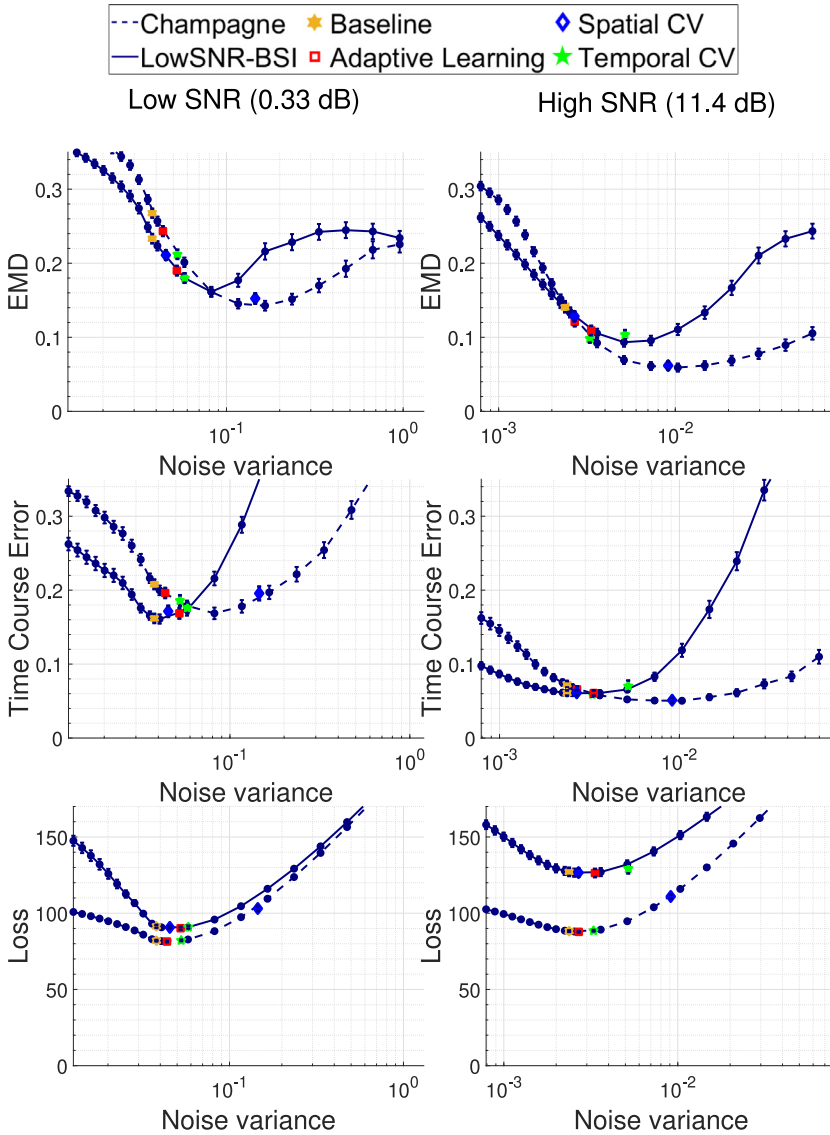


Fig. 1. Source reconstruction performance of Champagne and LowSNR-BSI in two different SNR regimes (low SNR: 0.33 dB, left column; high SNR: 11.4 dB, right column). Spatial reconstruction error is measured in terms of the earth-mover's distance, and is shown in the upper row, while time course reconstruction error is shown in the middle row. The lower row demonstrates the negative log-likelihood loss, SBL loss function Eq. (12), incurred by Champagne and LowSNR-BSI algorithms.

the choice of the baseline noise variance turns out to be suboptimal, as it is outperformed by all three proposed schemes that estimate the noise variance from the actual (task) data to be reconstructed ('Adaptive Learning', 'Spatial CV' and 'Temporal CV'). Interestingly, noise levels estimated using Spatial CV lead to near-optimal reconstruction performance in a broad variety of settings, in line with observations made in Habermehl et al. (2014); Haufe et al. (2011). All proposed noise learning schemes converge to points in the vicinity of the minimum of the SBL loss function Eq (12).

The EMD in our setting only depends on the spatial distribution of the sources. Therefore, the EMD is not able to fully capture potential advantages resulting from modeling temporal characteristics of the correlated EEG/MEG time courses. As a result, it is not highly aligned with the values of the loss. This explains the observed discrepancies between the loss function and the EMD values of Champagne and LowSNR-BSI in Fig. 1. To assess the reconstruction of the temporal characteristics of the brain sources, we also measure the time course error. All four variants of LowSNR-BSI algorithms not only outperform their Champagne counterparts but also approach the minimal achievable time course error. High EMD performance of Champagne with Spatial CV does not lead to high performance in terms of time course error as well as regarding the negative log-likelihood loss. For all algorithms, regularization values resulting in a smaller EMD metric can be found. However, this observation

does not imply a practical benefit of any algorithm as the ground-truth is unknown in real-world situations.

Fig. 2 further compares the source reconstruction performance of the four noise estimation variants separately for Champagne and LowSNR-BSI for a range of four SNR values. As already observed in Fig. 1, all three proposed approaches for noise variance estimation (adaptive learning, spatial CV, and temporal CV) lead to better source reconstruction performance than the estimation from baseline data. Overall, spatial CV for Champagne and temporal CV for LowSNR-BSI achieve the best combination of spatial and temporal reconstruction performance.

The superior performance of CV techniques, however, comes at the expense of higher computational complexity of the source reconstruction. As Fig. 2 demonstrates, using CV techniques with the specified numbers of folds increases the runtime of Champagne and LowSNR-BSI by approximately two orders of magnitude (10^3 s \sim 10^4 s) compared to the runtimes of eLORETA, S-FLEX, and the baseline and adaptive learning variants of Champagne and LowSNR-BSI (1 s \sim 10 s).

Fig. 3 provides an alternative depiction of the data presented in Fig. 2, which allows for a more direct comparison of Champagne and LowSNR-BSI. As benchmark algorithms, eLORETA (Pascual-Marqui, 2007) and S-FLEX (Haufe et al., 2011) are also included in the comparison. It can be seen that LowSNR-BSI in the baseline mode, using adaptive noise learning, and using temporal CV consistently outper-

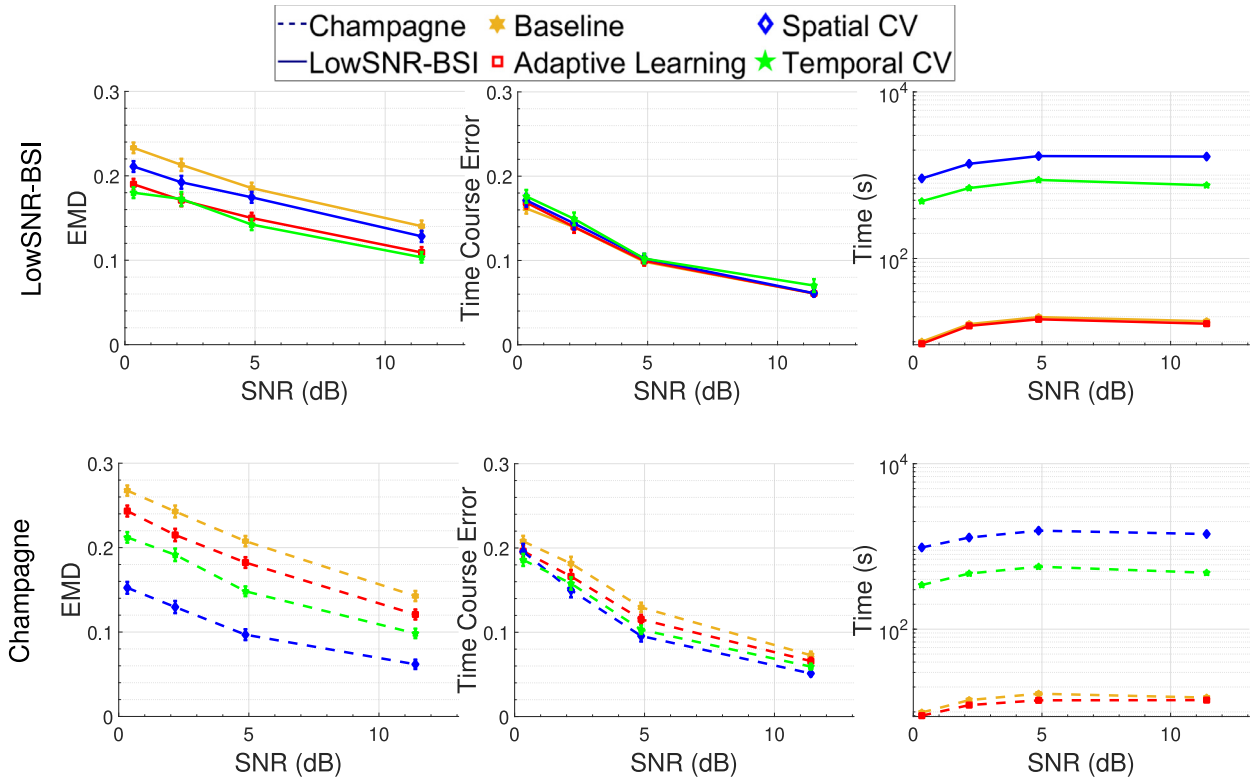


Fig. 2. Source reconstruction performance of four different variants of LowSNR-BSI (upper row) and Champagne (lower row). The noise variance was estimated from baseline data (ground truth), using adaptive learning, or using spatial or temporal cross-validation. Performance was evaluated for four SNRs (SNR = {0.33, 2.17, 4.87, 11.40} dB) and with respect to three different metrics (spatial reconstruction according to the earth-mover’s distance – left column, time course reconstruction error – middle column, and computational complexity according to the runtime (in seconds) – right column).

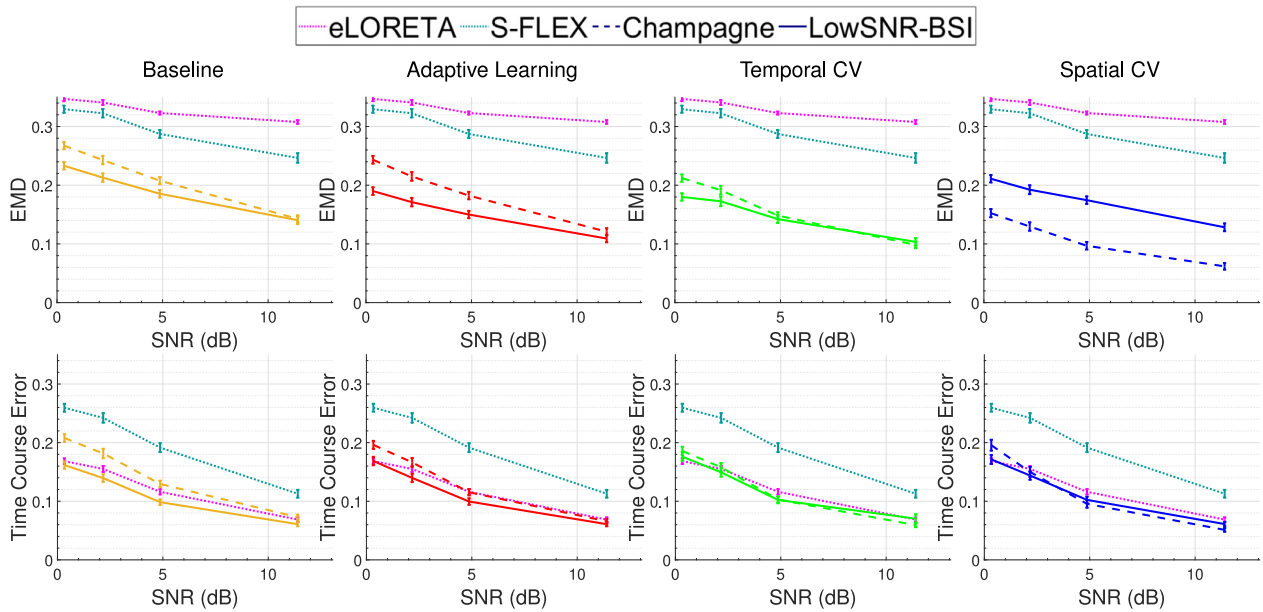


Fig. 3. Source reconstruction performance of Champagne (dashed line) and LowSNR-BSI (solid line) for four SNR values (SNR = {0.33, 2.17, 4.87, 11.40} dB). The noise variance was estimated from baseline data as well as using adaptive learning, spatial and temporal CV. Spatial reconstruction error was measured in terms of the earth-mover’s distance and is shown in the upper row, while time course reconstruction error is shown in the lower row.

forms Champagne in terms of spatial localization accuracy, in particular in low-SNR settings. This behavior indeed confirms the advantage of the surrogate function, $\mathcal{L}_{conv}^{Low-SNR}(\gamma|\gamma^k)$, which is designed to provide a better approximation of the non-convex SBL cost function in low-SNR regimes, as presented in Section 4. Consequently, as the SNR decreases, the gap between LowSNR-BSI and Champagne further increases. In terms of the

time course reconstruction error, LowSNR-BSI shows a similar improvement over Champagne when the SNR is low. However, the magnitude of this improvement is not as pronounced as observed for the EMD metric. The only setting in which Champagne consistently outperforms LowSNR-BSI is when spatial CV is used to estimate the noise variance, and spatial reconstruction performance is evaluated.

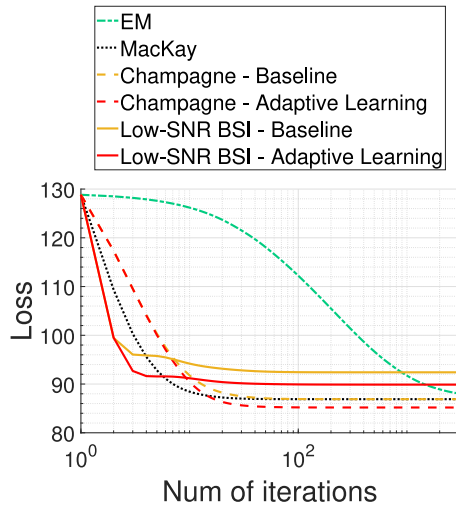


Fig. 4. Convergence behavior of LowSNR-BSI as well as Champagne using the standard (convex-bounding based) updates (Champagne) as well as EM and MacKay updates. For standard Champagne and LowSNR-BSI, the use of a fixed noise variance estimated from baseline data is compared with adaptive noise learning. LowSNR-BSI variants have faster convergence rate at early stages of the optimization procedure, but later converge to less optimal log-likelihood values. Adaptive learning variants of Champagne and LowSNR-BSI reach better log-likelihood values than their counterparts using a fixed noise variance estimated from baseline data.

Note that the LowSNR-BSI surrogate function in the baseline mode can provide a tight upper-bound for the original non-convex function only when SNR is equal to zero. For non-zero SNR, the current theory unfortunately does not apply; however, it is clear from our empirical results that the LowSNR-BSI surrogate function remains advantageous also in non-zero low-SNR regimes. As Fig. 3 demonstrates, we observe performance improvements of LowSNR-BSI over Champagne in the baseline mode for SNRs up to around 8 dB. After this point, the surrogate functions of LowSNR-BSI and Champagne both appear to be able to approximate the non-convex loss with a similar degree of precision; thus, their performance with respect to the evaluation metrics overlap as can be seen in the right column of Fig. 1.

It can further be observed that S-FLEX yields higher spatial localization accuracy (lower EMD) than eLORETA, while eLORETA yields higher temporal accuracy (lower time course error) than S-FLEX across all SNR values. With respect to spatial accuracy, both approaches, however, are consistently outperformed by Champagne and LowSNR-BSI. Note that the superior spatial reconstruction of sparsity-inducing algorithms (Champagne, LowSNR-BSI and S-FLEX) compared to eLORETA is expected here, because the simulated spatial distributions are indeed sparse. The superiority of SBL methods (Champagne, LowSNR-BSI) over S-FLEX that is observed here confirms observations and theoretical considerations made in Cai et al. (2021); Owen et al. (2012); Wipf et al. (2010). eLORETA shows comparable temporal reconstruction performance as LowSNR-BSI and Champagne, while S-FLEX is outperformed by all other methods.

The convergence behavior of the different SBL variants discussed and introduced in Sections 3–5 is illustrated in Fig. 4. LowSNR-BSI variants have faster convergence rates at the early stage of the optimization procedure compared to standard Champagne as well as Champagne with MacKay updates. They, however, reach lower negative log-likelihood values eventually, which indicates that they find better maxima of the model evidence. Furthermore, the adaptive-learning variants of Champagne and LowSNR-BSI reach lower negative log-likelihood values than their counterparts estimating the noise variance from baseline data, suggesting that learning the noise variance, or in other words overes-

timating the noise variance, improves the reconstruction performance through better model evidence maximization.

Note that the plots in Fig. 4 demonstrate the convergence behaviour of MM algorithms for only one single experiment. We conducted another experiment (see Appendix I), in which the simulation was carried out 100 times using different instances of source distributions and initializations. The final negative log-likelihood loss – attained after convergence – and runtimes of all methods were calculated. The median and interquartile ranges over 100 randomized experiments of these performance metrics are reported in Fig. 8, which confirms the observations made here.

7. Analysis of auditory evoked fields (AEF)

The MEG data used here were acquired in the Biomagnetic Imaging Laboratory at the University of California San Francisco (UCSF) with a CTF Omega 2000 whole-head MEG system from VSM MedTech (Coquitlam, BC, Canada) with 1200 Hz sampling rate. The neural responses of one subject to an Auditory Evoked Fields (AEF) stimulus were localized. The AEF response was elicited with single 600 ms duration tones (1 kHz) presented binaurally. The data were averaged across 120 trials (after the trials were time-aligned to the stimulus). The pre-stimulus window was selected to be -100 ms to 5 ms and the post-stimulus time window was selected to be 5 ms to 250 ms, where 0 ms is the onset of the tone. Further details on this dataset can be found in Cai et al. (2021); Dalal et al. (2011); Owen et al. (2012). The lead field for each subject was calculated with NUTMEG (<http://bil.ucsf.edu>) using a single-sphere head model (two spherical orientation lead fields) and an 8 mm voxel grid.

The results presented in Section 6 have been obtained for the scalar setting, where the orientation of the brain sources are assumed to be perpendicular to the surface of cortex and, hence, only the scalar deflection of each source along the fixed orientation needs to be estimated. In real data, surface normals are hard to estimate or even undefined in case of volumetric reconstructions. Consequently, we model each source here as a full 3-dimensional current vector. This is achieved by introducing three variance parameters for each source within the source covariance matrix, $\mathbf{\Gamma}^{3D} = \text{diag}(\gamma^{3D}) = [\gamma_1^x, \gamma_1^y, \gamma_1^z, \dots, \gamma_N^x, \gamma_N^y, \gamma_N^z]^T$. As all algorithms considered here model the source covariance matrix $\mathbf{\Gamma}$ to be diagonal, this extension can be readily implemented. Correspondingly, a full 3D lead-field matrix, $\mathbf{L}^{3D} \in \mathbb{R}^{M \times 3N}$, is used.

Fig. 5 shows the reconstructed sources of the AEF of one subject using conventional Champagne with pre-estimated $\lambda = \hat{\sigma}^2$, adaptive noise learning, and spatial CV. LowSNR-BSI with pre-estimated $\lambda = \hat{\sigma}^2$ was also included in the comparison. Shown in the top panel are the reconstructions at the time of the maximal deflection of the auditory N100 component (shown in bottom panel).

All reconstructions are able to correctly localize bilateral auditory activity to Heschel’s gyrus, which is the location of the primary auditory cortex. Note that an additional source in the midbrain, which is indicated by all three Champagne variants, is absent for LowSNR-BSI.

We tested the reconstruction performance of all methods for random subsets of 10, 20, 40, 60, and 100 trials. As Fig. 6 shows, the proposed noise learning variants of Champagne as well as LowSNR-BSI can correctly localize bilateral auditory activity to Heschel’s gyrus even when using as few as 10 trials. Focusing on the low-SNR regime, Fig. 7 shows seven reconstructions for random selections of 10 trials. LowSNR-BSI as well as all proposed noise learning variants of Champagne consistently show sources at the expected locations in the left and auditory cortices, where both cortices are jointly identified in the majority of experiments.

8. Discussion

We have provided a unifying theoretical platform for deriving different sparse Bayesian learning algorithms for electromagnetic brain imaging using the Majorization-Minimization (MM) framework. First, we

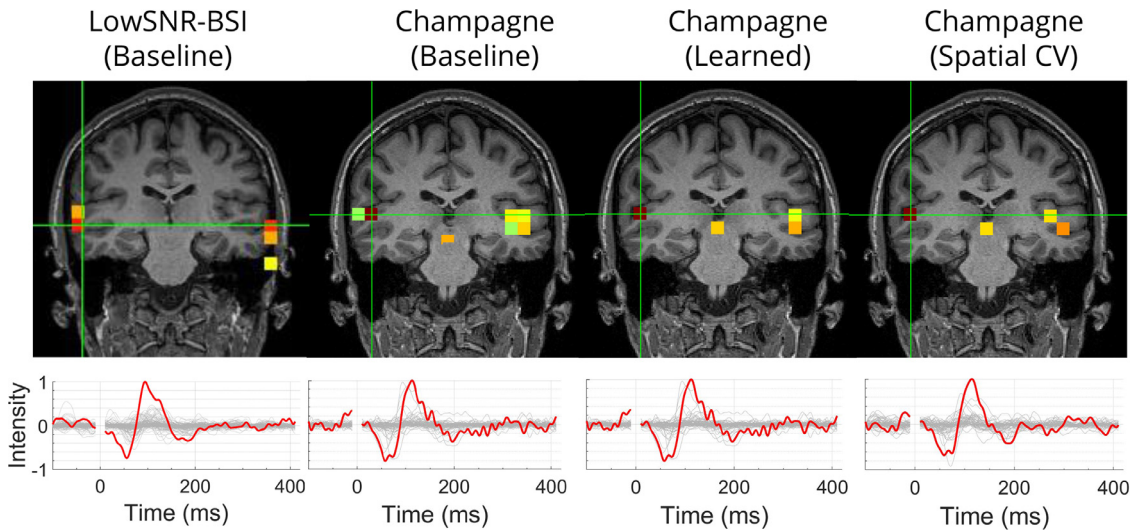


Fig. 5. Analysis of auditory evoked fields (AEF) of one subject using conventional Champagne with pre-estimated $\lambda = \hat{\sigma}^2$, adaptive noise learning, and spatial CV as well as LowSNR-BSI. Shown in the top panel are the reconstructions at the time of the maximal deflection of the auditory N100 component (shown in bottom panel). All reconstructions show sources at the expected locations in the left and right auditory cortex.

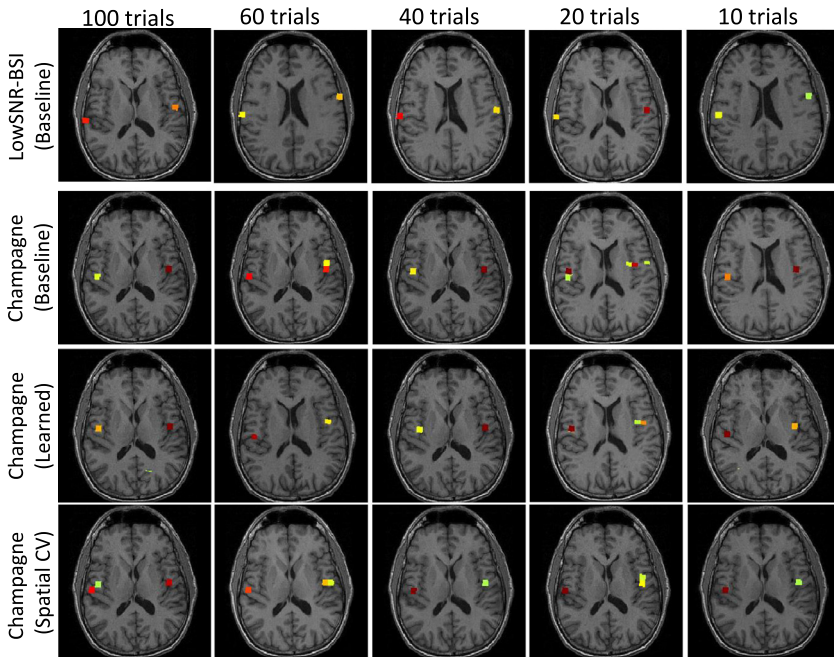


Fig. 6. Analysis of auditory evoked fields (AEF) of one subject using conventional Champagne with pre-estimated $\lambda = \hat{\sigma}^2$, adaptive noise learning, and spatial CV as well as LowSNR-BSI, tested with the number of trials limited to 10, 20, 40, 60, and 100. All proposed noise learning reconstructions show sources at the expected locations in the left and right auditory cortices.

demonstrated that the choice of upper bounds of the Type-II non-convex loss function within the MM framework influences the reconstruction performance and convergence rates of the resulting algorithms. Second, focusing on commonly occurring low-SNR settings, we derived a novel Type-II Bayesian algorithm, LowSNR-BSI, using a novel convex bounding MM function that converges to the original loss function as the SNR goes to zero. We demonstrated the advantage of LowSNR-BSI over existing benchmark algorithms including Champagne, eLORETA and S-FLEX. Consistent with the theoretical considerations, the advantage of LowSNR-BSI over Champagne decreases with increasing SNR. Third, we have derived an analytic solution that allows us to estimate the noise variance jointly within the source estimation procedure on the same (task-related) data that are used for the reconstruction. We have also adopted cross-validation schemes to empirically estimate the noise variance from hold-out data through a line search. We have proposed spatial and temporal CV schemes, where either subsets of EEG/MEG channels

or recorded samples are left out of the source reconstruction, and where the noise variance is selected as the minimizer of a divergence between model and hold-out data. We also demonstrate that precise knowledge of the noise variance is required in order to determine the optimal algorithm performance. Finally, according to our empirical results, all three proposed techniques for estimating the noise variance lead to superior source reconstruction performance compared to the setting in which the noise variance is estimated from baseline data.

8.1. Cross-validation vs. adaptive noise learning

Spatial CV for Champagne and Temporal CV for LowSNR-BSI achieved the best performances and are generally applicable to any distributed inverse solution. Their long computation time can, however, be challenging as their computational complexity is drastically higher (around two orders of magnitude) than using baseline data or adap-

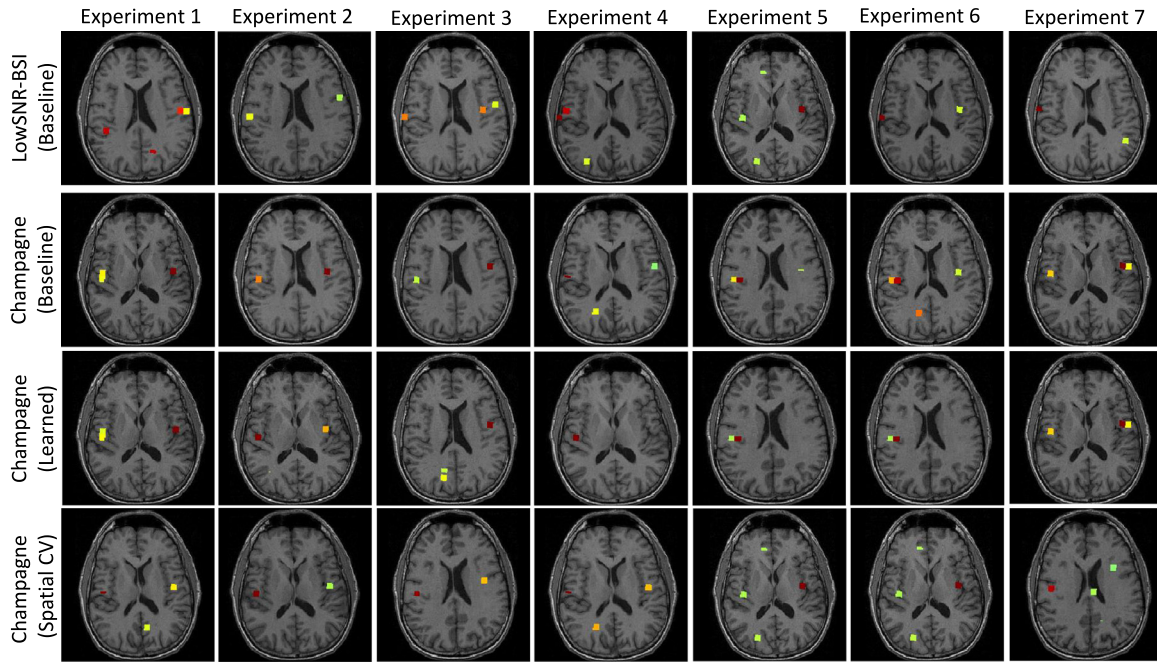


Fig. 7. Analysis of auditory evoked fields (AEF) of one subject using conventional Champagne with pre-estimated $\lambda = \hat{\sigma}^2$, adaptive noise learning, and spatial CV as well as LowSNR-BSI, tested with the number of trials limited to 10. Each column shows an experiment with a random selection of 10 trials. LowSNR-BSI as well as all proposed noise learning variants of Champagne always show sources at the expected locations in the left or right auditory cortex. In the majority of experiments, both cortices are jointly identified.

tive learning schemes. The high complexity of CV techniques is a potential limitation in settings where the efficiency of the algorithm or immediate access to the outcome is crucial. What is more, this approach quickly becomes infeasible if more than one parameter needs to be estimated through a grid search. In contrast, the computational complexity of the proposed noise level estimation scheme using adaptive learning is of the same order as the complexity of the baseline approach. Moreover, we have successfully extended this approach to the estimation of heteroscedastic noise, where a distinct variance is estimated for each M/EEG sensor (Cai et al., 2021). Hence, the adaptive-learning approach can be seen as an advancement of the baseline algorithm that combines performance improvement and computational efficiency. It is also worth noting that the computational complexity of CV techniques heavily relies on tunable parameters such as the number of folds/splits of the data and the total number of candidate points in the grid search.

8.2. Interpretation of Type-I and Type-II loss functions as divergences

We have pointed out (see Section 5.2 and Appendix A) that Type-I and Type-II Bayesian approaches implicitly use different metrics to compare the empirical sensor-space observations to the signal proportion explained by the reconstructed brain sources. Type-I approaches measure first-order differences between modeled and reconstructed time series using variants of the MSE, while Type-II approaches amount to using the log-det Bregman divergence to measure differences in the second-order statistics of the empirically observed and modeled data as summarized in the respective covariance matrices. While the connection between the Type-II loss function and the log-det Bregman divergence has been investigated and exploited in numerous forms such as *Stein's loss* (James and Stein, 1992) or the *graphical Lasso* (Friedman et al., 2008; Mazumder and Hastie, 2012; Ravikumar et al., 2011), and has found applications in disciplines such as information theory and metric learning (Davis et al., 2007; Zadeh et al., 2016), wireless communication (Khalilsarai et al., 2020), and signal processing (Khanna and Murthy, 2017a; 2017b; Wiesel et al., 2015), it has not received much attention in the BSI literature to the best of authors' knowledge. Here, we have used

this insight to devise a novel cross-validation scheme, temporal CV, in which model fit is measured in terms of the log-det Bregman divergence (or, Type-II likelihood) on held-out samples. In contrast, the previously introduced spatial CV uses the mean-squared error to measure out-of-sample model fit. Importantly, however, this difference does not imply that the application of spatial CV is restricted to Type-I approaches or that the use of temporal CV is restricted to Type-II approaches. Rather, both approaches are universally applicable. In fact, it is straightforward to evaluate the Type-I likelihood based on the source times series reconstructed with Type-II methods. Conversely, it is also possible to estimate the Type-II likelihood for Type-I approaches such as S-FLEX. Here, the model source and noise covariances are first estimated from the reconstructed sources as $\hat{\Gamma} = \text{Cov}[\mathbf{x}(t)]$ and $\hat{\sigma}^2 = 1/M \sum_m [\mathbf{C}_y - \mathbf{L}\hat{\Gamma}\mathbf{L}^T]_{[m,m]}$, after which Σ_y can be calculated. The optimal Type-I regularization parameter is then selected as the minimizer of $\mathcal{L}^{\text{II}}(\mathbf{Y}^{\text{train,temp}}, \mathbf{Y}^{\text{test,temp}})$ in Eq. (31).

8.3. Limitations and future work

One limiting assumption of the current work is that the activity of the sources is modeled to be independent across voxels, spatial orientations, and time samples. Analogously, the noise is assumed to be independent across times samples, and homoscedastic (independent with equal variance across sensors). These assumptions merely act as prior information whose purpose is to bias the inverse reconstruction towards solutions with lower complexity. Thus, they do not prevent the reconstruction of brain and noise sources with more complex structure if the observed data are inconsistent with these priors. On the other hand, modeling dependency structures that are in fact present in real data has the potential to substantially improve the source reconstruction. We have recently proposed adaptive noise learning algorithms that relax the rather unrealistic assumption of homoscedastic noise (Cai et al., 2021). Going further, it would be possible to also model spatial covariances of the sources between voxels and/or between source orientation within voxels, which would encode the realistic assumption that individual brain regions do not work in isolation. Similarly, the spatial covariance struc-

ture of the noise could be modeled in order to accommodate spatially distributed artifacts due to, for example, heart beat or line noise interference. Finally, electrophysiological data are known to possess a complex intrinsic autocorrelation structure, which is not modeled by the majority of existing BSI algorithms. We have recently proposed ways to also learn temporal correlations within the Type-II framework and have obtained promising results with respect to time course reconstruction (Hashemi and Haufe, 2018; Hashemi et al., 2021).

9. Conclusion

We have provided a unifying theoretical platform for deriving different sparse Bayesian learning algorithms for electromagnetic brain imaging using the Majorization-Minimization (MM) framework. This unification perspective not only provides a useful theoretical framework for comparing different algorithms in terms of their convergence behavior, but also provides a principled recipe for constructing novel algorithms with specific properties by designing appropriate bounds of the Bayesian marginal likelihood function. Building on MM principles, we then proposed a novel method called *LowSNR-BSI* that achieves favorable source reconstruction performance in low signal-to-noise-ratio settings. Recognizing the importance of noise estimation for algorithm performance, we present both analytical and cross-validation approaches for noise estimation. Empirically, we show that the monotonous convergence behavior predicted from MM theory is confirmed in numerical experiments. Using simulations, we further demonstrate the advantage of *LowSNR-BSI* over conventional Champagne in low-SNR regimes, and the advantage of learned noise levels over estimates derived from baseline data. To demonstrate the usefulness of our novel approach, we show neurophysiologically plausible source reconstructions on averaged auditory evoked potential data.

Our characterization of the Type-II likelihood as a divergence measure provides a novel perspective on the construction of BSI algorithms and might open new avenues of research in this field. It is conceivable that alternative divergence metrics can be used for solving the M/EEG source reconstruction problem in the future by modeling specific neurophysiologically valid aspects of similarity between data and model output. Promising metrics in that respect are information divergences such as Kullback-Leibler (KL) (Wei et al., 2020), Rényi (Khanna and Murthy, 2017b), Itakura-Saito (IS) (Févotte et al., 2009) and β divergences (Cichocki and Amari, 2010; Eguchi and Kato, 2010; Févotte and Idier, 2011; Samek et al., 2013) as well as transportation metrics such as the Wasserstein distance between empirical and statistical covariances (e.g., Gramfort et al., 2015; Janati et al., 2020; Peyré et al., 2019; Villani, 2008).

Although this paper focuses on electromagnetic brain source imaging, Type-II methods have also been successfully developed in other fields such as direction of arrival (DoA) and channel estimation in wireless communications (Gerstoft et al., 2016; Haghghatshoar and Caire, 2017; Khalilsarai et al., 2020; Prasad et al., 2015), Internet of Things (IoT) (Fengler et al., 2019a; 2019b), robust portfolio optimization in finance (Feng et al., 2016), covariance matching and estimation (Benfenati et al., 2020; Greenewald and Hero, 2015; Meriaux et al., 2020; Ollila et al., 2020; Ottersten et al., 1998; Tsiligkaridis et al., 2013; Werner et al., 2008; Zoubir et al., 2018), graph learning (Kumar et al., 2020), and brain functional imaging (Wei et al., 2020). The methods introduced in this work may also prove useful in these domains.

Data and code availability statement

The auditory evoked fields data used in this study will be made publicly available in deidentified and pre-processed form through a public data repository such as openneuro.org or www.nitrc.org. Moreover, MATLAB code used to analyze the AEF data as well code used in the simulation study will be uploaded to a publicly accessible GitHub repository. This includes the code to generate the synthetic data, implementations of the

brain source imaging (BSI) algorithms, and codes to evaluate source reconstruction performance.

Credit authorship contribution statement

Ali Hashemi: Conceptualization, Methodology, Software, Validation, Formal analysis, Investigation, Data curation, Writing - original draft, Writing - review & editing, Visualization. **Chang Cai:** Methodology, Data curation, Writing - review & editing. **Gitta Kutyniok:** Supervision, Writing - review & editing. **Klaus-Robert Müller:** Methodology, Writing - original draft, Writing - review & editing, Supervision, Funding acquisition. **Srikantan S. Nagarajan:** Conceptualization, Methodology, Formal analysis, Writing - original draft, Writing - review & editing, Visualization. **Stefan Haufe:** Conceptualization, Methodology, Validation, Formal analysis, Investigation, Data curation, Writing - original draft, Writing - review & editing, Visualization, Supervision, Project administration, Funding acquisition.

Acknowledgments

This result is part of a project that has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (Grant agreement No. 758985).

AH acknowledges scholarship support from the Machine Learning/Intelligent Data Analysis research group at Technische Universität Berlin. He further wishes to thank the Berlin International Graduate School in Model and Simulation based Research (BIMoS), the Berlin Mathematical School (BMS), and the Berlin Mathematics Research Center MATH+ for partial support. CC was supported by the National Natural Science Foundation of China under Grant 62007013. GK acknowledges partial support by the Bundesministerium für Bildung und Forschung (BMBF) through the Berliner Zentrum für Machine Learning (BZML), Project AP4, RTG DAEDALUS (RTG 2433), Projects P1 and P3, RTG BIOQIC (RTG 2260), Projects P4 and P9, and by the Berlin Mathematics Research Center MATH+, Projects EF1-1 and EF1-4. KRM was supported in part by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea Government (No. 2017-0-00451, Development of BCI based Brain and Cognitive Computing Technology for Recognizing User's Intentions using Deep Learning) and (No. 2019-0-00079, Artificial Intelligence Graduate School Program, Korea University), and by the German Ministry for Education and Research (BMBF) under Grants 01IS14013A-E, 01GQ1115, 01GQ0850, 01IS18025A, 031L0207D and 01IS18037A; the German Research Foundation (DFG) under Grant Math+, EXC 2046/1, Project ID 390685689. SSN was funded in part by National Institutes of Health grants (R01DC004855, R01EB022717, R01DC176960, R01DC010145, R01NS100440, R01AG062196, and R01DC013979), University of California MRPI MRP-17454755, the US Department of Defense grant (W81XWH-13-1-0494).

Appendix A. Bregman Divergence Formulation of the Type-II Loss Function

We start by recalling the definition of log-det Bregman matrix divergence - also known as Stein's loss (James and Stein, 1992) - between any two $M \times M$ positive semidefinite (PSD) matrices \mathbf{Q} and \mathbf{W} :

$$D_{\log\text{-det}}(\mathbf{Q}, \mathbf{W}) = \text{tr}(\mathbf{Q}\mathbf{W}^{-1}) - \log|\mathbf{Q}\mathbf{W}^{-1}| - M, \quad (35)$$

where the "log-det" Bregman matrix divergence in (35) is an special case of Bregman matrix divergence (Bregman, 1967), where $-\log|\cdot|$ is selected as a strictly convex function. By substituting \mathbf{C}_y and $\mathbf{\Sigma}_y$ in (35) instead of \mathbf{Q} and \mathbf{W} , the *log-det* Bregman matrix divergence can be written as follows (Davis et al., 2007; Friedman et al., 2008; Jalali et al., 2017; Khalilsarai et al., 2020; Khanna and Murthy, 2017a; Mazumder

and Hastie, 2012; Ravikumar et al., 2011; Tsiligkaridis and Hero, 2013; Zadeh et al., 2016):

$$\begin{aligned} D_{\log\text{-det}}(\mathbf{C}_y, \Sigma_y) &= \text{tr}(\mathbf{C}_y \Sigma_y^{-1}) - \log |\mathbf{C}_y \Sigma_y^{-1}| - M \\ &= \text{tr}(\mathbf{C}_y \Sigma_y^{-1}) + \log |\Sigma_y| - \log |\mathbf{C}_y| - M \\ &= \log |\Sigma_y| + \text{tr}(\mathbf{C}_y \Sigma_y^{-1}) + \underbrace{-\log |\mathbf{C}_y|}_{\text{const}} - M, \end{aligned} \quad (36)$$

where (36) is the same as (13) up to a constant. Note that $\log |\mathbf{C}_y|$ does not depend on γ and is, therefore, treated as a constant value here.

Appendix B. Proof of Corollary 1

Proof. To verify the descending trend in the MM framework, it is sufficient to show that $f(\mathbf{u}^{k+1}) \leq f(\mathbf{u}^k)$. To this end, we have $f(\mathbf{u}^{k+1}) \leq g(\mathbf{u}^{k+1}|\mathbf{u}^k)$ from condition [A2]. Condition [A3] further states that $g(\mathbf{u}^{k+1}|\mathbf{u}^k) \leq g(\mathbf{u}^k|\mathbf{u}^k)$, while $g(\mathbf{u}^k|\mathbf{u}^k) = f(\mathbf{u}^k)$ holds according to [A1]. Putting everything together, we have:

$$f(\mathbf{u}^{k+1}) \stackrel{[A2]}{\leq} g(\mathbf{u}^{k+1}|\mathbf{u}^k) \stackrel{[A3]}{\leq} g(\mathbf{u}^k|\mathbf{u}^k) \stackrel{[A1]}{=} f(\mathbf{u}^k),$$

which concludes the proof. \square

Appendix C. Proof of Proposition 1

Proof. We first show that the objective function of the M-step is derived by upper-bounding the negative log-likelihood, $-\log p(\mathbf{Y}|\gamma)$, using Jensen's inequality (J):

$$\begin{aligned} -\log p(\mathbf{Y}|\gamma) &= -\log E_{p(\mathbf{X}|\gamma)} p(\mathbf{Y}|\mathbf{X}, \gamma) = -\log E_{p(\mathbf{X}|\gamma)} \left(\frac{p(\mathbf{X}|\mathbf{Y}, \gamma^k) p(\mathbf{Y}|\mathbf{X}, \gamma)}{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \right) \\ &\stackrel{(I)}{=} -\log E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \left(\frac{p(\mathbf{Y}|\mathbf{X}, \gamma)}{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\gamma) \right) \\ &\stackrel{(J)}{\leq} -E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \log \left(\frac{p(\mathbf{Y}|\mathbf{X}, \gamma)}{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\gamma) \right) \\ &\stackrel{(II)}{=} -E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \log p(\mathbf{Y}, \mathbf{X}|\gamma) + \underbrace{E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \log p(\mathbf{X}|\mathbf{Y}, \gamma^k)}_{\text{const}} \\ &:= \mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k). \end{aligned} \quad (37)$$

The resulting bound is a majorizing function for $-\log p(\mathbf{Y}|\gamma)$, so that condition [A2] holds. Note that the term $E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\mathbf{Y}, \gamma^k)$ does not depend on γ and, therefore, does not influence the optimization. According to the definition of Jensen's inequality, the equality constraint – condition [A1] – holds if and only if the argument of the convex function is a constant. Therefore, to establish the equivalence of both sides of (J) when $\gamma = \gamma^k$, it is sufficient to show that the argument of the log function, $\frac{p(\mathbf{Y}|\mathbf{X}, \gamma)}{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\gamma)$, is constant when $\gamma = \gamma^k$. This can be verified by invoking Bayes rule:

$$\frac{p(\mathbf{Y}|\mathbf{X}, \gamma^k)}{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} p(\mathbf{X}|\gamma^k) = p(\mathbf{Y}|\gamma^k).$$

Since $p(\mathbf{Y}|\gamma^k)$ is a constant, equality condition [A1] holds.

After inserting the analytic form of $-\log p(\mathbf{Y}, \mathbf{X}|\gamma)$ in Eq. (24):

$$-\log p(\mathbf{Y}, \mathbf{X}|\gamma) = \frac{T}{2} \log |\Gamma| + \frac{1}{2} \sum_{t=1}^T \mathbf{x}(t)^\top \Gamma^{-1} \mathbf{x}(t) + \frac{T}{2} \log |2\sigma^2 \mathbf{I}| + \sum_{t=1}^T \frac{1}{\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\mathbf{x}(t)\|_2^2,$$

we are ready to prove that $\mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k)$ fulfills condition [A3]. We have:

$$\mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k) \propto \log |\Gamma| + E_{p(\mathbf{X}|\mathbf{Y}, \gamma^k)} \left[\frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \Gamma^{-1} \bar{\mathbf{x}}^k(t) \right] + \text{const}, \quad (38)$$

where const comprises all terms of Eq. (25) that are not a function of γ . To prove that $\mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k)$ satisfies condition [A3], we need to show

that $\mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k)$ reaches to its global minimum in each MM iteration. This can be easily guaranteed if Eq. (38) is convex. While the second term in (38) is convex, the first term, $\log |\Gamma|$, is in fact concave, which hampers conclusions concerning the convexity of their sum. However, we can use the concept of *geodesic convexity* or *g-convexity* from non-Euclidean and geometric optimization, which enables us to prove that any local minimum of Eq. (38) is actually a global minimum. For the sake of brevity, we will omit a detailed theoretical introduction of g-convexity, and only borrow the following required propositions, Propositions 5 and 6, from the literature (an interested reader can refer to (Wiesel et al., 2015, Chapter 1) for a gentle introduction to this topic, and to (Papadopoulos, 2005, Chapter 2) (Ben-Tal, 1977; Bonnabel and Sepulchre, 2009; Liberti, 2004; Moakher, 2005; Pallaschke and Rolewicz, 2013; Rapcsak, 1991; Vishnoi, 2018) for more in-depth technical details). Now, we state the following preliminary results: \square

Proposition 5. The function $\log |\Gamma|$ is g-convex in Γ , where Γ belongs to the manifold of positive definite (PD) matrices.

Proof. A detailed proof can be found in (Wiesel et al., 2015, Lemma. 1.13). The main idea is to leverage the geodesic $\mathbf{Q}_q = \mathbf{V}\mathbf{D}^q\mathbf{V}^\top$, $q \in [0, 1]$ between two matrices, $\mathbf{Q}_0 = \mathbf{V}\mathbf{V}^\top$ and $\mathbf{Q}_1 = \mathbf{V}\mathbf{D}\mathbf{V}^\top$, in order to transfer the problem into the following form:

$$f(\mathbf{Q}_q) = \log |\mathbf{V}\mathbf{D}^q\mathbf{V}^\top| = 2 \log |\mathbf{V}| + q \log |\mathbf{D}|,$$

where $f(\mathbf{Q}_q)$ is a linear function and, therefore, convex in q . \square

Remark 7. The log-determinant function is concave in classical Euclidean analysis. However, Proposition 5 demonstrates that it is g-convex with respect to the PD manifold.

Proposition 6. Any local minimum of a g-convex function over a g-convex set is a global minimum.

Proof. A detailed proof is presented in (Rapcsak, 1991, Theorem 2.1).

Given that g-convexity is an extension of classical convexity to non-Euclidean geometry, it is straightforward to show that all convex functions are also g-convex, where the geodesics between pairs of matrices are simply line segments. Therefore, given Proposition 5, we can conclude that Eq. (38) is g-convex; hence, any local minimum of $\mathcal{L}_{\text{EM}}^k(\gamma|\gamma^k)$ is a global minimum according to Proposition 6. This proves that condition [A3] is fulfilled and completes the proof of Proposition 1. \square

Appendix D. Proof of Proposition 2

Proof. We start by recalling $\mathcal{R}^{\Pi-x}(\mathbf{X}, \gamma)$ in Eq. (14):

$$\mathcal{R}^{\Pi-x}(\mathbf{X}, \gamma) = \frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{x_n(t)^2}{\gamma_n} + \log |\Sigma_y|.$$

Based on (Sun et al., 2017, Example 2), [A2] can be directly inferred from the concavity of the log-determinant function and its first-order Taylor expansion around the value from the previous iteration, Σ_y^k , which leads to the following inequality:

$$\begin{aligned} \log |\Sigma_y| &\leq \log |\Sigma_y^k| + \text{tr} \left[\left(\Sigma_y^k \right)^{-1} (\Sigma_y - \Sigma_y^k) \right] \\ &= \log |\Sigma_y^k| + \text{tr} \left[\left(\Sigma_y^k \right)^{-1} \Sigma_y \right] - \text{tr} \left[\left(\Sigma_y^k \right)^{-1} \Sigma_y^k \right]. \end{aligned} \quad (39)$$

Note that the first and last term in (39) do not depend on γ ; hence, they can be ignored in the optimization procedure. Conditions [A1] and [A4] are automatically satisfied by construction because the majorizing function is obtained through a Taylor expansion around Σ_y^k . Concretely, [A1] is satisfied because the equality in Eq. (39) holds for $\Sigma_y = \Sigma_y^k$. Similarly, [A4] is satisfied because the gradient of $\log |\Sigma_y|$ at point Σ_y^k , $\left(\Sigma_y^k \right)^{-1}$, defines the linear Taylor approximation $\log |\Sigma_y^k| +$

$\text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\left(\boldsymbol{\Sigma}_y - \boldsymbol{\Sigma}_y^k\right)\right]$. Thus, both gradients coincide in $\boldsymbol{\Sigma}_y^k$ by construction. Now, we show that [A3] can be satisfied easily using standard optimization algorithms by proving that $\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k)$ is a convex function with respect to γ . To this end, we rewrite Eq. (26):

$$\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k) = \left[\frac{1}{T} \sum_{t=1}^T \sum_{n=1}^N \frac{\bar{x}_n^k(t)^2}{\gamma_n}\right] + \log|\boldsymbol{\Sigma}_y^k| + \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\boldsymbol{\Sigma}_y\right] - \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\boldsymbol{\Sigma}_y^k\right],$$

as follows:

$$\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k) = \text{diag}[\mathbf{U}]\gamma^{-1} + \text{diag}[\mathbf{V}]\gamma + \text{const}, \quad (40)$$

where $\mathbf{U} := \frac{1}{T} \sum_{t=1}^T [\bar{\mathbf{x}}^k(t) \bar{\mathbf{x}}^k(t)^\top]$ and $\mathbf{V} := \mathbf{L}^\top \left(\boldsymbol{\Sigma}_y^k\right)^{-1} \mathbf{L}$ are defined as parameters that do not depend on γ . The term const also collects constant terms in (39), i.e. $\text{const} := \log|\boldsymbol{\Sigma}_y^k| + \sigma^2 \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\right] - M$. Besides, $\gamma^{-1} = [\gamma_1^{-1}, \dots, \gamma_N^{-1}]^\top$ is defined as the element-wise inversion of γ . The convexity of $\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k)$ can be directly inferred from the convexity of $\text{diag}[\mathbf{U}]\gamma^{-1}$ and $\text{diag}[\mathbf{V}]\gamma$ with respect to γ (Boyd and Vandenberghe, 2004, Chapter. 3). The convexity of $\mathcal{R}_{\text{conv}}^k(\gamma|\gamma^k)$, which ensures that condition [A3] can be satisfied using standard optimization, along with fulfillment of conditions [A1], [A2] and [A4], ensure that Theorem 1 holds.

In order to establish the equivalence of the MM algorithm using the majorization function Eq. (26) and the convex-bounding based Champagne variant presented in Section 2.3.2, we here decompose $\boldsymbol{\Sigma}_y$ into rank-one matrices as introduced in (Sun et al., 2016). The first term of Eq. (26) can be reformulated as follows:

$$\begin{aligned} \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\boldsymbol{\Sigma}_y\right] &= \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\left(\sigma^2\mathbf{I} + \mathbf{L}\mathbf{L}^\top\right)\right] \\ &= \text{tr}\left[\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\tilde{\mathbf{L}}\tilde{\mathbf{L}}^\top\right] = \text{diag}\left[\tilde{\mathbf{L}}^\top\left(\boldsymbol{\Sigma}_y^k\right)^{-1}\tilde{\mathbf{L}}\right] \tilde{\boldsymbol{\gamma}}, \end{aligned} \quad (41)$$

where $\tilde{\mathbf{L}} = \text{diag}(\gamma_1, \dots, \gamma_N, \sigma^2, \dots, \sigma^2)$, and $\tilde{\mathbf{L}} = [\mathbf{L}, \mathbf{I}]$. Since we are optimizing Eq. (26) with respect to γ_n , for $n = 1, \dots, N$, the elements of $\tilde{\mathbf{L}}$ and $\tilde{\boldsymbol{\gamma}}$ related to the sensor noise σ^2 vanish. Thus, by inserting Eq. (41) into Eq. (26), taking the derivative with respect to γ_n , for $n = 1, \dots, N$, and setting it to zero,

$$\begin{aligned} \frac{\partial}{\partial \gamma_n} \left(\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2 \gamma_n^{-1} + \left[\mathbf{L}_n^\top \left(\boldsymbol{\Sigma}_y^k\right)^{-1} \mathbf{L}_n \right] \gamma_n \right) \\ = -\frac{1}{(\gamma_n)^2} \left(\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2 \right) + \left[\mathbf{L}_n^\top \left(\boldsymbol{\Sigma}_y^k\right)^{-1} \mathbf{L}_n \right] \\ = 0 \quad \text{for } n = 1, \dots, N, \end{aligned}$$

where \mathbf{L}_n denotes the n -th column of the lead field matrix, we obtain an update rule in terms of the original variables $\boldsymbol{\Gamma}$ and \mathbf{L} :

$$\gamma_n^{k+1} := \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2}{\mathbf{L}_n^\top \left(\boldsymbol{\Sigma}_y^k\right)^{-1} \mathbf{L}_n}}, \quad (42)$$

which is identical to the update rule of the convex-bounding based approach discussed in Section 2.3.2, Eqs. (17)–(19). \square

Appendix E. Proof of Proposition 3

Proof. The proof that conditions [A1]–[A4] are satisfied is directly analogous to that of Proposition 2; therefore, it is omitted here. The equivalence of the Champagne variant based on MacKay updates (Wipf and Nagarajan, 2009, Section III.A-2) presented in Section 2.3.3 and the solution derived within the MM framework can be derived by transforming the update rule Eq. (42) into a fixed-point iteration of the form $\gamma^{k+1} = f(\gamma^k)$, which is an alternative way of minimizing the same surrogate function (Eq. (26)). By squaring the left and right hand sides of

Eq. (42), one can divide both sides by γ_n^{k+1} and re-interpret the term on the right hand side as the estimate from the previous (k -th) iteration:

$$\gamma_n^{k+1} := \left[\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2 \right] \left(\gamma_n^k \mathbf{L}_n^\top \left(\boldsymbol{\Sigma}_y^k\right)^{-1} \mathbf{L}_n \right)^{-1} \quad (43)$$

for $n = 1, \dots, N$. This is indeed identical to the MacKay update in Eq. (22), which concludes the proof. \square

Appendix F. Proof of Proposition 4

Proof. (following (Haghighatshoar and Caire, 2017, Appendix C-A)) Without loss of generality, we here consider the case $\sigma^2 = 1$, which can be obtained by normalizing the sensor and source covariance matrices by σ^2 : $\boldsymbol{\Gamma} \leftarrow \boldsymbol{\Gamma}/\sigma^2$, $\boldsymbol{\Sigma}_y \leftarrow \boldsymbol{\Sigma}_y/\sigma^2 = \mathbf{I} + \mathbf{L}\mathbf{L}^\top$. Also, due to the concavity of the $\log(\cdot)$ function and by using a Taylor expansion around point a , we have:

$$\log(x) = \log a + \frac{x}{a} - 1 + \mathcal{O}(x), \quad \forall a > 0. \quad (44)$$

Assuming that $\mathbf{L}\mathbf{L}^\top$ has an eigenvalue decomposition $\mathbf{L}\mathbf{L}^\top = \mathbf{U}\mathbf{P}\mathbf{U}^\top$ with $\mathbf{P} = \text{diag}(p_1, \dots, p_M)$, the majorizing function $\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k)$ as well as Eq. (28) are derived as follows:

$$\begin{aligned} \log|\boldsymbol{\Sigma}_y| &= \log|\mathbf{I} + \mathbf{U}\mathbf{P}\mathbf{U}^\top| \stackrel{\text{(I)}}{=} \sum_{i=1}^M \log(1 + p_i) \stackrel{\text{(II)}}{=} \sum_{i=1}^M p_i + \mathcal{O}(p_i) \\ &= \text{tr}(\mathbf{L}\mathbf{L}^\top) + \mathcal{O}(\text{SNR}), \end{aligned} \quad (45)$$

where the p_i , for $i = 1, \dots, M$ denote the diagonal elements of \mathbf{P} , which are equivalent to the eigenvalues of $\mathbf{L}\mathbf{L}^\top$. The term $\mathcal{O}(p_i)$ represents the second and higher-order residuals of the Taylor expansion. Note that (45)-(I) is obtained by expanding \mathbf{P} over its diagonal elements, while (45)-(II) is derived by exploiting the concavity of the $\log(\cdot)$ function and its first-order Taylor expansion around $a = 1$ based on Eq. (44). Given the eigenvalue decomposition of $\mathbf{L}\mathbf{L}^\top = \mathbf{U}\mathbf{P}\mathbf{U}^\top$ and the normalization with respect to the noise variance, the sum over all eigenvalues of $\mathbf{L}\mathbf{L}^\top$, i.e., $\sum_{i=1}^M p_i$, represents the ratio between the power of the signal and the power of the noise; hence, one can replace $\sum_{i=1}^M \mathcal{O}(p_i)$ in Eq. (45) with $\mathcal{O}(\text{SNR})$. To elaborate this more, note that given $\sigma^2 = 1$, we have $\text{SNR} \propto \mathbb{E}\{\|\mathbf{L}\mathbf{x}(t)\|^2\} = \text{tr}(\mathbf{L}\mathbf{L}^\top)$, where $\mathbb{E}\{\mathbf{x}(t)\mathbf{x}(t)^\top\} = \boldsymbol{\Gamma} = \text{diag}([\gamma_1, \dots, \gamma_N]^\top)$ due to the independence between voxels. Therefore, $\text{SNR} \propto \text{tr}(\mathbf{L}\mathbf{L}^\top) = \sum_{i=1}^M p_i$, as the sum of the eigenvalues of a matrix is equal to its trace.

As we have shown that $\log|\boldsymbol{\Sigma}_y| = \text{tr}(\mathbf{L}\mathbf{L}^\top) + \mathcal{O}(\text{SNR})$, condition [A2] holds and $\mathcal{L}^{\text{H}}(\gamma)$ converges to $\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k)$ when $\text{SNR} \rightarrow 0$.

Moreover, as Eq. (28) is constructed using a linear Taylor approximation, [A1] and [A4] hold due to the same arguments made in the proof of Proposition 2. It remains to be shown that condition [A3] can be easily fulfilled due to the convexity of $\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k)$. To this end, we exploit the following key relationship between the sensor and source space covariances:

$$\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) = \frac{1}{T} \sum_{t=1}^T \left[\frac{1}{\lambda} \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2 + \bar{\mathbf{x}}^k(t)^\top \boldsymbol{\Gamma}^{-1} \bar{\mathbf{x}}^k(t) \right]. \quad (46)$$

By replacing $\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t)$ in Eq. (27) with its source space equivalence in (46), we have:

$$\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k) = \text{tr}(\mathbf{L}\mathbf{L}^\top) + \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \boldsymbol{\Gamma}^{-1} \bar{\mathbf{x}}^k(t) + \text{const}, \quad (47)$$

where const denotes the terms that do not depend on γ . Reformulating (47) as

$$\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k) = \text{diag}[\mathbf{W}]\gamma + \text{diag}[\mathbf{Q}]\gamma^{-1} + \text{const},$$

with $\mathbf{W} := \mathbf{L}^\top \mathbf{L}$, $\mathbf{Q} := \frac{1}{T} \sum_{t=1}^T [\bar{\mathbf{x}}^k(t) \bar{\mathbf{x}}^k(t)^\top]$ and $\gamma^{-1} = [\gamma_1^{-1}, \dots, \gamma_N^{-1}]^\top$ proves the convexity of $\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\gamma|\gamma^k)$ using the same arguments made for proving convexity in Proposition 2. Thus, we have shown that conditions [A1]–[A4] hold, which concludes the proof. \square

Appendix G. Detailed Derivation of the LowSNR-BSI Algorithm

To find the optimal value of $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$, we take the derivative of $\mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\boldsymbol{\gamma}|\boldsymbol{\gamma}^k)$ in (27) with respect to each γ_n for $n = 1, \dots, N$:

$$\begin{aligned} \frac{\partial}{\partial \gamma_n} \mathcal{L}_{\text{conv}}^{\text{Low-SNR}}(\boldsymbol{\gamma}|\boldsymbol{\gamma}^k) &= \frac{\partial}{\partial \gamma_n} \left[\text{tr}(\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top) + \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) \right] \\ &\stackrel{\text{(I)}}{=} \frac{\partial}{\partial \gamma_n} \left[\sum_{n=1}^N \gamma_n \mathbf{L}_n^\top \mathbf{L}_n + \frac{\partial}{\partial \gamma_n} \left[\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) \right] \right] \\ &\stackrel{\text{(II)}}{=} \mathbf{L}_n^\top \mathbf{L}_n + \frac{\partial}{\partial \gamma_n} \sum_{t=1}^T \frac{1}{T} \left[\frac{1}{\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2 + \bar{\mathbf{x}}^k(t)^\top \boldsymbol{\Gamma}^{-1} \bar{\mathbf{x}}^k(t) \right] \\ &\stackrel{\text{(III)}}{=} \mathbf{L}_n^\top \mathbf{L}_n + \left[\frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \left(\frac{\partial}{\partial \gamma_n} \boldsymbol{\Gamma}^{-1} \right) \bar{\mathbf{x}}^k(t) \right] \\ &= \mathbf{L}_n^\top \mathbf{L}_n + \left(-\frac{1}{\gamma_n^2} \right) \left[\frac{1}{T} \sum_{t=1}^T \left(\bar{\mathbf{x}}_n^k(t) \right)^2 \right], \end{aligned} \quad (48)$$

where Eq. (48)-I is derived based on a *sum-of-rank-one matrices* reformulation of the term $\text{tr}(\mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top)$ by exploiting the diagonal structure of $\boldsymbol{\Gamma}$. Equality (48)-II is the direct implication of the duality between $\boldsymbol{\gamma}$ -space and \mathbf{X} -space that has been pointed out in (14). Finally, $\frac{1}{\sigma^2} \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2$ does not appear in (48)-III and is ignored since it does not depend on $\boldsymbol{\gamma}$. Setting the derivative in Eq. (48) to zero yields the following closed-form update for $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$:

$$\gamma_n^{k+1} := \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T \left(\bar{\mathbf{x}}_n^k(t) \right)^2}{\mathbf{L}_n^\top \mathbf{L}_n}} \text{ for } n = 1, \dots, N,$$

which is identical to the update rule in Eq. (29). This completes the derivation of the LowSNR-BSI algorithm.

Appendix H. Proof of Theorem 2

Proof. We start by taking the derivative of $\mathcal{L}^\text{II}(\lambda)$ with respect to λ :

$$\frac{\partial}{\partial \lambda} \mathcal{L}^\text{II}(\lambda) = \frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_y|) + \frac{\partial}{\partial \lambda} \left[\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) \right]. \quad (49)$$

We first calculate the first term, $\frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_y|)$. Using the matrix inversion equality

$$\log|\boldsymbol{\Sigma}_y| = \log|\lambda \mathbf{I} + \mathbf{L}\boldsymbol{\Gamma}\mathbf{L}^\top| = \log\left| \frac{1}{\lambda} \mathbf{L}^\top \mathbf{L} + \boldsymbol{\Gamma}^{-1} \right| + \log|\boldsymbol{\Gamma}| + \log|\lambda \mathbf{I}|,$$

we have

$$\frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_y|) = \frac{\partial}{\partial \lambda} \left(M \log \lambda + \log\left| \frac{1}{\lambda} \mathbf{L}^\top \mathbf{L} + \boldsymbol{\Gamma}^{-1} \right| \right) = \frac{\partial}{\partial \lambda} \left(M \log \lambda + \log|\boldsymbol{\Sigma}_x^{-1}| \right),$$

where the term $\log|\boldsymbol{\Gamma}|$ is omitted since it does not depend on λ . Then, the derivative of $\log|\boldsymbol{\Sigma}_y|$ with respect to λ can be obtained as follows:

$$\frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_y|) = \frac{M}{\lambda} - \left(\frac{1}{\lambda^2} \right) \text{tr}[\boldsymbol{\Sigma}_x \mathbf{L}^\top \mathbf{L}], \quad (50)$$

where the second term in (50) is derived according to the equality $\boldsymbol{\Sigma}_x^{-1} = (\boldsymbol{\Gamma}^{-1} + \frac{1}{\lambda} \mathbf{L}^\top \mathbf{L})$, which holds for the inverse of the posterior covariance in Eq. (10) (Sekihara and Nagarajan, 2015, Chapter 4):

$$\frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_x^{-1}|) = \text{tr} \left[\boldsymbol{\Sigma}_x \frac{\partial}{\partial \lambda} \boldsymbol{\Sigma}_x^{-1} \right] = \text{tr} \left[\boldsymbol{\Sigma}_x \frac{\partial}{\partial \lambda} \left(\boldsymbol{\Gamma}^{-1} + \frac{1}{\lambda} \mathbf{L}^\top \mathbf{L} \right) \right]$$

$$= \text{tr} \left[\boldsymbol{\Sigma}_x \frac{\partial}{\partial \lambda} \left(\frac{1}{\lambda} \mathbf{L}^\top \mathbf{L} \right) \right] = - \left(\frac{1}{\lambda^2} \right) \text{tr}[\boldsymbol{\Sigma}_x \mathbf{L}^\top \mathbf{L}].$$

In the next step, we calculate the derivative of the second term in Eq. (49) using the following key relation between the sensor and source space covariances presented in Appendix F. Given (46), we have

$$\frac{\partial}{\partial \lambda} \left[\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) \right] = \left(-\frac{1}{\lambda^2} \right) \frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2, \quad (51)$$

where the term $\bar{\mathbf{x}}^k(t)^\top \boldsymbol{\Gamma}^{-1} \bar{\mathbf{x}}^k(t)$ is neglected since it does not depend on λ . Let $\boldsymbol{\Gamma}^k$ and $\boldsymbol{\Sigma}_x^k$ be fixed values obtained in the (k) -th iteration. Then, by substituting Eqs. (50) and (51) into Eq. (49), we have:

$$\begin{aligned} \frac{\partial}{\partial \lambda} \mathcal{L}^\text{II}(\lambda) &= \frac{\partial}{\partial \lambda} (\log|\boldsymbol{\Sigma}_y|) + \frac{\partial}{\partial \lambda} \left[\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \boldsymbol{\Sigma}_y^{-1} \mathbf{y}(t) \right] \\ &= \frac{M}{\lambda} - \left(\frac{1}{\lambda^2} \right) \text{tr}[\boldsymbol{\Sigma}_x^k \mathbf{L}^\top \mathbf{L}] + \left(-\frac{1}{\lambda^2} \right) \frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2. \end{aligned} \quad (52)$$

By expressing $\text{tr}[\boldsymbol{\Sigma}_x^k \mathbf{L}^\top \mathbf{L}]$ in terms of the values at the (k) -th iteration according to the following matrix equality (Zhang, Rao, 2011):

$$\text{tr}[\boldsymbol{\Sigma}_x^k \mathbf{L}^\top \mathbf{L}] = \text{tr}[\boldsymbol{\Sigma}_x^k \lambda^k (\boldsymbol{\Sigma}_x^k)^{-1} - (\boldsymbol{\Gamma}^k)^{-1}] = \text{tr}[\lambda^k \mathbf{I}_{N^k}] - \text{tr}[\lambda^k (\boldsymbol{\Sigma}_x^k) (\boldsymbol{\Gamma}^k)^{-1}],$$

Eq. (52) can be reformulated as follows:

$$\frac{\partial}{\partial \lambda^k} \mathcal{L}^\text{II}(\lambda^k) = \frac{M}{\lambda^k} - \frac{1}{(\lambda^k)^2} \text{tr}[\lambda^k \mathbf{I}_{N^k}] + \frac{1}{(\lambda^k)^2} \text{tr}[\lambda^k (\boldsymbol{\Sigma}_x^k) (\boldsymbol{\Gamma}^k)^{-1}] - \frac{1}{(\lambda^k)^2} \frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2.$$

Note that N^k denotes the number of non-zero voxels at the (k) -th iteration. Now by setting the derivative to zero, the update rule for λ at the $(k+1)$ -th iteration is obtained as

$$\lambda^{k+1} := \frac{\frac{1}{T} \sum_{t=1}^T \|\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t)\|_2^2}{M - N^k + \text{tr}[(\boldsymbol{\Sigma}_x^k) (\boldsymbol{\Gamma}^k)^{-1}]}.$$

This completes the proof. \square

Appendix I. A Statistical Analysis of Computational Complexity and Convergence Behaviour of MM Methods

Here, we conducted an experiment, in which the simulation presented in Fig. 4 was carried out 100 times using different instances of source distributions and initializations. The final negative log-likelihood loss – attained after convergence – and runtimes of all methods were calculated. The median and inter-quartile ranges over 100 randomized experiments of these performance metrics are reported in Fig. 8.

As demonstrated in Fig. 4, the EM algorithm indeed needs a larger number of iterations for convergence than its peer MM variants, which eventually results in longer runtimes and higher computational complexity if we measure runtime in units of seconds, demonstrated in Fig. 8-(A). The overall computation complexity in each iteration of the EM, however, is comparable to the other MM variants. Even though an

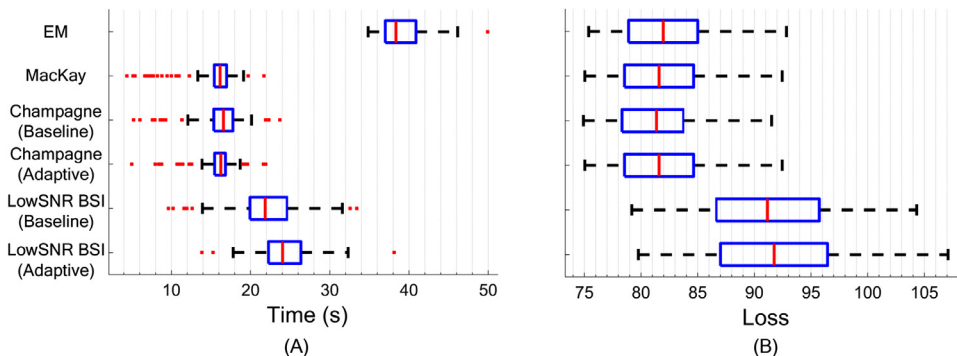


Fig. 8. (A) Runtime and (B) Type-II negative log marginal likelihood loss attained after convergence of different variants of LowSNR-BSI and Champagne, as well as using EM and MacKay updates. Shown are the median and inter-quartile ranges over 100 randomized experiments.

additional operation for calculating the posterior matrix of the sources, Σ_x , is involved in each iteration of the EM algorithm – which operates in the high-dimensional source space, efficient implementation techniques can drastically reduce the computational complexity of this operation, e.g., from $\mathcal{O}(N^2)$ to $\mathcal{O}(N)$, since only the main diagonal elements of Σ_x are required in the update rule, i.e., $[\Sigma_x^k]_{n,n}$. Therefore, the overall computational complexity of the EM algorithm at each iteration is dominated by $\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2$, which is a common term in all other MM-based approaches, e.g., convex bounding, MacKay, and LowSNR-BSI. Interested readers can refer to (Zumer, Attias, Sekihara, Nagarajan, 2007) for a computational analysis of EM and other Type-II methods. Fig. 8-(B) also depicts the median and inter-quartile ranges of the final negative log-likelihood loss – attained after convergence – of different variants of LowSNR-BSI and Champagne as well as Champagne using EM and MacKay updates.

References

- Baillet, S., Mosher, J.C., Leahy, R.M., 2001. Electromagnetic brain mapping. *IEEE Signal Process Mag* 18 (6), 14–30.
- Bauschke, H.H., Combettes, P.L., 2017. Fenchel–Rockafellar Duality. In: *Convex Analysis and Monotone Operator Theory in Hilbert Spaces*. Springer, pp. 247–262.
- Bekhti, Y., Lucka, F., Salmon, J., Gramfort, A., 2018. A hierarchical bayesian perspective on majorization-minimization for non-convex sparse regression: application to m/EEG source imaging. *Inverse Probl* 34 (8), 085010.
- Ben-Tal, A., 1977. On generalized means and generalized convex functions. *J Optim Theory Appl* 21 (1), 1–13.
- Benfenati, A., Chouzenoux, E., Pesquet, J.-C., 2020. Proximal approaches for matrix optimization problems: application to robust precision matrix estimation. *Signal Processing* 169, 107417.
- Benidis, K., Feng, Y., Palomar, D.P., et al., 2018. Optimization methods for financial index tracking: from theory to practice. *Foundations and Trends® in Optimization* 3 (3), 171–279.
- Bijma, F., De Munck, J.C., Huizenga, H.M., Heethaar, R.M., 2003. A mathematical approach to the temporal stationarity of background noise in MEG/EEG measurements. *Neuroimage* 20 (1), 233–243.
- Bishop, C.M., 2006. *Pattern recognition and machine learning*. Springer.
- Blankertz, B., Lemm, S., Treder, M., Haufe, S., Müller, K.-R., 2011. Single-trial analysis and classification of ERP components – a tutorial. *Neuroimage* 56 (2), 814–825.
- Bonnabel, S., Sepulchre, R., 2009. Riemannian metric and geometric mean for positive semidefinite matrices of fixed rank. *SIAM J. Matrix Anal. Appl.* 31 (3), 1055–1070.
- Boyd, S.P., Vandenberghe, L., 2004. *Convex optimization*. Cambridge university press.
- Bregman, L.M., 1967. The relaxation method of finding the common point of convex sets and its application to the solution of problems in convex programming. *USSR Computational Mathematics and Mathematical Physics* 7 (3), 200–217.
- Cai, C., Hashemi, A., Diwakar, M., Haufe, S., Sekihara, K., Nagarajan, S.S., 2021. Robust estimation of noise for electromagnetic brain imaging with the champagne algorithm. *Neuroimage* 225, 117411.
- Cai, C., Sekihara, K., Nagarajan, S.S., 2018. Hierarchical multiscale bayesian algorithm for robust MEG/EEG source reconstruction. *Neuroimage* 183, 698–715.
- Calvetti, D., Pascarella, A., Pitolli, F., Somersalo, E., Vantaggi, B., 2019. Brain activity mapping from MEG data via a hierarchical bayesian algorithm with automatic depth weighting. *Brain Topogr* 32 (3), 363–393.
- Calvetti, D., Somersalo, E., 2018. Inverse problems: from regularization to bayesian inference. *Wiley Interdiscip. Rev. Comput. Stat.* 10 (3), e1247.
- Castañó-Candamil, S., Höhne, J., Martínez-Vargas, J.-D., An, X.-W., Castellanos-Domínguez, G., Haufe, S., 2015. Solving the EEG inverse problem based on space–time–frequency structured sparsity constraints. *Neuroimage* 118, 598–612.
- Cichocki, A., Amari, S.-i., 2010. Families of alpha-beta-and gamma-divergences: flexible and robust measures of similarities. *Entropy* 12 (6), 1532–1568.
- Dalal, S.S., Zumer, J.M., Guggisberg, A.G., Trumppis, M., Wong, D.D., Sekihara, K., Nagarajan, S.S., 2011. MEG/EEG Source reconstruction, statistical evaluation, and visualization with NUTMEG. *Comput Intell Neurosci* 2011.
- Dale, A.M., Liu, A.K., Fischl, B.R., Buckner, R.L., Belliveau, J.W., Lewine, J.D., Halgren, E., 2000. Dynamic statistical parametric mapping: combining fMRI and MEG for high-resolution imaging of cortical activity. *Neuron* 26 (1), 55–67.
- Dale, A.M., Sereno, M.I., 1993. Improved localization of cortical activity by combining EEG and MEG with MRI cortical surface reconstruction: a linear approach. *J Cogn Neurosci* 5 (2), 162–176.
- Davis, J.V., Kulis, B., Jain, P., Sra, S., Dhillon, I.S., 2007. Information-theoretic metric learning. In: *Proceedings of the 24th International Conference on Machine Learning*, pp. 209–216.
- De Munck, J.C., Huizenga, H.M., Waldorp, L.J., Heethaar, R., 2002. Estimating stationary dipoles from MEG/EEG data contaminated with spatially and temporally correlated background noise. *IEEE Trans. Signal Process.* 50 (7), 1565–1572.
- Dempster, A.P., Laird, N.M., Rubin, D.B., 1977. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society: Series B (Methodological)* 39 (1), 1–22.
- Eguchi, S., Kato, S., 2010. Entropy and divergence associated with power function and the statistical application. *Entropy* 12 (2), 262–274.
- Engemann, D.A., Gramfort, A., 2015. Automated model selection in covariance estimation and spatial whitening of MEG and EEG signals. *Neuroimage* 108, 328–342.
- Feng, Y., Palomar, D.P., et al., 2016. A signal processing perspective on financial engineering. *Foundations and Trends® in Signal Processing* 9 (1–2), 1–231.
- Fengler, A., Caire, G., Jung, P., Haghigatshoar, S., 2019. Massive MIMO unsourced random access. *arXiv preprint arXiv:1901.00828*.
- Fengler, A., Haghigatshoar, S., Jung, P., Caire, G., 2019. Non-Bayesian activity detection, large-scale fading coefficient estimation, and unsourced random access with a massive MIMO receiver. *arXiv preprint arXiv:1910.11266*.
- Févotte, C., 2011. Majorization-minimization algorithm for smooth Itakura-Saito nonnegative matrix factorization. In: *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, pp. 1980–1983.
- Févotte, C., Bertin, N., Durrieu, J.-L., 2009. Nonnegative matrix factorization with the itakura-saito divergence: with application to music analysis. *Neural Comput* 21 (3), 793–830.
- Févotte, C., Idier, J., 2011. Algorithms for nonnegative matrix factorization with the β -divergence. *Neural Comput* 23 (9), 2421–2456.
- Friedman, J., Hastie, T., Tibshirani, R., 2008. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics* 9 (3), 432–441.
- Friston, K.J., Penny, W., Phillips, C., Kiebel, S., Hinton, G., Ashburner, J., 2002. Classical and bayesian inference in neuroimaging: theory. *Neuroimage* 16 (2), 465–483.
- Gerstoft, P., Mecklenbräuker, C.F., Xenaki, A., Nannuru, S., 2016. Multisnapshot sparse bayesian learning for DOA. *IEEE Signal Process Lett* 23 (10), 1469–1473.
- Gorodnitsky, L.F., George, J.S., Rao, B.D., 1995. Neuromagnetic source imaging with FOCUSS: a recursive weighted minimum norm algorithm. *Electroencephalogr Clin Neurophysiol* 95 (4), 231–251.
- Gramfort, A., 2009. Mapping, timing and tracking cortical activations with MEG and EEG: Methods and application to human vision. *Ecole nationale supérieure des telecommunications-ENST*.
- Gramfort, A., Kowalski, M., Hämäläinen, M., 2012. Mixed-norm estimates for the m/EEG inverse problem using accelerated gradient methods. *Phys Med Biol* 57 (7), 1937.
- Gramfort, A., Peyré, G., Cuturi, M., 2015. Fast optimal transport averaging of neuroimaging data. In: *2016 International Conference on Information Processing in Medical Imaging*. Springer, pp. 261–272.
- Gramfort, A., Strohmeier, D., Haueisen, J., Hämäläinen, M.S., Kowalski, M., 2013. Time-frequency mixed-norm estimates: sparse m/EEG imaging with non-stationary source activations. *Neuroimage* 70, 410–422.
- Greenewald, K., Hero, A.O., 2015. Robust kronecker product PCA for spatio-temporal covariance estimation. *IEEE Trans. Signal Process.* 63 (23), 6368–6378.
- Habermehl, C., Steinbrink, J.M., Müller, K.-R., Haufe, S., 2014. Optimizing the regularization for image reconstruction of cerebral diffuse optical tomography. *J Biomed Opt* 19 (9), 096006.
- Haghigatshoar, S., Caire, G., 2017. Massive MIMO channel subspace estimation from low-dimensional projections. *IEEE Trans. Signal Process.* 65 (2), 303–318.
- Hämäläinen, M., Hari, R., Ilmoniemi, R.J., Knuutila, J., Lounasmaa, O.V., 1993. Magnetoencephalography theory, instrumentation, and applications to noninvasive studies of the working human brain. *Rev Mod Phys* 65 (2), 413.
- Hämäläinen, M.S., Ilmoniemi, R.J., 1994. Interpreting magnetic fields of the brain: minimum norm estimates. *Medical & Biological Engineering & Computing* 32 (1), 35–42.
- Hashemi, A., Haufe, S., 2018. Improving EEG source localization through spatio-temporal sparse bayesian learning. In: *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE, pp. 1935–1939.
- Hashemi, A., Nagarajan, S.S., Müller, K.-R., Haufe, S., 2021. Spatio-temporal brain source imaging using sparse bayesian learning: mathematical guarantees and trade-off. Preprint.
- Hastie, T., Tibshirani, R., Friedman, J., 2009. *The elements of statistical learning: Data mining, inference, and prediction*. Springer Science & Business Media.
- Haufe, S., Nikulin, V.V., Ziehe, A., Müller, K.-R., Nolte, G., 2008. Combining sparsity and rotational invariance in EEG/MEG source reconstruction. *Neuroimage* 42 (2), 726–738.
- Haufe, S., Tomioka, R., Dickhaus, T., Sannelli, C., Blankertz, B., Nolte, G., Müller, K.-R., 2011. Large-scale EEG/MEG source localization with spatial flexibility. *Neuroimage* 54 (2), 851–859.
- Huang, Y., Parra, L.C., Haufe, S., 2016. The new york head – a precise standardized volume conductor model for EEG source localization and tES targeting. *Neuroimage* 140, 150–162.
- Huizenga, H.M., De Munck, J.C., Waldorp, L.J., Grasman, R.P., 2002. Spatiotemporal EEG/MEG source analysis based on a parametric noise covariance model. *IEEE Trans. Biomed. Eng.* 49 (6), 533–539.
- Hunter, D.R., Lange, K., 2004. A tutorial on MM algorithms. *Am Stat* 58 (1), 30–37.
- Jacobson, M.W., Fessler, J.A., 2007. An expanded theoretical treatment of iteration-dependent majorize-minimize algorithms. *IEEE Trans. Image Process.* 16 (10), 2411–2422.
- Jalali, A., Saunderson, J., Fazel, M., Hassibi, B., 2017. Error bounds for Bregman denoising and structured natural parameter estimation. In: *2017 IEEE International Symposium on Information Theory (ISIT)*. IEEE, pp. 2273–2277.
- James, W., Stein, C., 1992. *Estimation with Quadratic Loss*. In: *Breakthroughs in Statistics*. Springer, pp. 443–460.
- Janati, H., Bazeille, T., Thirion, B., Cuturi, M., Gramfort, A., 2020. Multi-subject MEG/EEG source imaging with sparse multi-task regression. *Neuroimage* 220, 116847.
- Jun, S.C., Plis, S.M., Ranken, D.M., Schmidt, D.M., 2006. Spatiotemporal noise covariance estimation from limited empirical magnetoencephalographic data. *Physics in Medicine & Biology* 51 (21), 5549.
- Khalilsarai, M.B., Yang, T., Haghigatshoar, S., Caire, G., 2020. Structured channel covariance estimation from limited samples in Massive MIMO. In: *ICC 2020-2020 IEEE International Conference on Communications (ICC)*. IEEE, pp. 1–7.

- Khanna, S., Murthy, C.R., 2017. On the support recovery of jointly sparse gaussian sources using sparse bayesian learning. arXiv preprint arXiv:1703.04930.
- Khanna, S., Murthy, C.R., 2017. Rényi divergence based covariance matching pursuit of joint sparse support. In: 18th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC). IEEE, pp. 1–5.
- Kohavi, R., et al., 1995. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: IJcai, 14. Montreal, Canada, pp. 1137–1145.
- Kumar, S., Ying, J., de Miranda Cardoso, J.V., Palomar, D.P., 2020. A unified framework for structured graph learning via spectral constraints. *Journal of Machine Learning Research* 21 (22), 1–60.
- Lemm, S., Blankertz, B., Dickhaus, T., Müller, K.-R., 2011. Introduction to machine learning for brain imaging. *Neuroimage* 56 (2), 387–399.
- Liberti, L., 2004. On a class of nonconvex problems where all local minima are global. *Publications de l'Institut Mathématique* 76 (90), 101–109.
- Lipp, T., Boyd, S., 2016. Variations and extension of the convex–concave procedure. *Optimization and Engineering* 17 (2), 263–287.
- Luessi, M., Hämläinen, M.S., Solo, V., 2013. Sparse component selection with application to MEG source localization. In: 2013 IEEE 10th International Symposium on Biomedical Imaging. IEEE, pp. 556–559.
- Matsuura, K., Okabe, Y., 1995. Selective minimum-norm solution of the biomagnetic inverse problem. *IEEE Trans. Biomed. Eng.* 42 (6), 608–615.
- Mazumder, R., Hastie, T., 2012. The graphical lasso: new insights and alternatives. *Electron J Stat* 6, 2125.
- Meriaux, B., Ren, C., Breloy, A., El Korso, M.N., Forster, P., 2020. Matched and mismatched estimation of kronecker product of linearly structured scatter matrices under elliptical distributions. *IEEE Trans. Signal Process.*
- Mika, S., Rätsch, G., Müller, K.-R., 2001. A mathematical programming approach to the kernel fisher algorithm. *Adv Neural Inf Process Syst* 591–597.
- Moakher, M., 2005. A differential geometric approach to the geometric mean of symmetric positive-definite matrices. *SIAM J. Matrix Anal. Appl.* 26 (3), 735–747.
- Nunez, P.L., Srinivasan, R., et al., 2006. Electric fields of the brain: The neurophysics of EEG. Oxford University Press, USA.
- Oguz-Ekim, P., Gomes, J.P., Xavier, J., Oliveira, P., 2011. Robust localization of nodes and time-recursive tracking in sensor networks using noisy range measurements. *IEEE Trans. Signal Process.* 59 (8), 3930–3942.
- Ollila, E., Palomar, D.P., Pascal, F., 2020. Shrinking the eigenvalues of m-estimators of covariance matrix. *IEEE Trans. Signal Process.*
- Oostenveld, R., Praamstra, P., 2001. The five percent electrode system for high-resolution EEG and ERP measurements. *Clinical Neurophysiology: Official Journal of the International Federation of Clinical Neurophysiology* 112 (4), 713–719. doi:10.1016/s1388-2457(00)00527-7.
- Ottersten, B., Stoica, P., Roy, R., 1998. Covariance matching estimation techniques for array signal processing applications. *Digit Signal Process* 8 (3), 185–210.
- Owen, J.P., Wipf, D.P., Attias, H.T., Sekihara, K., Nagarajan, S.S., 2012. Performance evaluation of the champagne source reconstruction algorithm on simulated and real m/EEG data. *Neuroimage* 60 (1), 305–323.
- Pallaschke, D.E., Rolewicz, S., 2013. Foundations of mathematical optimization: Convex analysis without linearity, 388. Springer Science & Business Media.
- Papadopoulos, A., 2005. Metric spaces, convexity and nonpositive curvature, 6. European Mathematical Society.
- Pascual-Marqui, R.D., 2007. Discrete, 3D distributed, linear imaging methods of electric neuronal activity. part 1: exact, zero error localization. arXiv preprint arXiv:0710.3341.
- Pascual-Marqui, R.D., Michel, C.M., Lehmann, D., 1994. Low resolution electromagnetic tomography: a new method for localizing electrical activity in the brain. *International Journal of Psychophysiology* 18 (1), 49–65.
- Pascual-Marqui, R.D., et al., 2002. Standardized low-resolution brain electromagnetic tomography (sloreta): technical details. *Methods Find Exp Clin Pharmacol* 24 (Suppl D), 5–12.
- Peyré, G., Cuturi, M., et al., 2019. Computational optimal transport: with applications to data science. *Foundations and Trends® in Machine Learning* 11 (5–6), 355–607.
- Plis, S.M., Schmidt, D.M., Jun, S.C., Ranken, D.M., 2006. A generalized spatiotemporal covariance model for stationary background in analysis of MEG data. In: 2006 International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE, pp. 3680–3683.
- Prasad, R., Murthy, C.R., Rao, B.D., 2015. Joint channel estimation and data detection in MIMO-OFDM systems: a sparse bayesian learning approach. *IEEE Trans. Signal Process.* 63 (20), 5369–5382.
- Rapcsak, T., 1991. Geodesic convexity in nonlinear optimization. *J Optim Theory Appl* 69 (1), 169–183.
- Ravikumar, P., Wainwright, M.J., Raskutti, G., Yu, B., et al., 2011. High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electron J Stat* 5, 935–980.
- Razaviyayn, M., Hong, M., Luo, Z.-Q., 2013. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM J. Optim.* 23 (2), 1126–1153.
- Rockafellar, R.T., 1970. Convex analysis. Princeton University Press.
- Rubner, Y., Tomasi, C., Guibas, L.J., 2000. The earth mover's distance as a metric for image retrieval. *Int J Comput Vis* 40 (2), 99–121.
- Samek, W., Kawanabe, M., Müller, K.-R., 2013. Divergence-based framework for common spatial patterns algorithms. *IEEE Rev Biomed Eng* 7, 50–72.
- Seeger, M.W., Wipf, D.P., 2010. Variational bayesian inference techniques. *IEEE Signal Process Mag* 27 (6), 81–91.
- Sekihara, K., Nagarajan, S.S., 2015. Electromagnetic brain imaging: A bayesian perspective. Springer.
- Shalev-Shwartz, S., Ben-David, S., 2014. Understanding machine learning: From theory to algorithms. Cambridge University Press.
- Shen, K., Yu, W., Zhao, L., Palomar, D.P., 2019. Optimization of MIMO device-to-Device networks via matrix fractional programming: A Minorization–Maximization approach. *IEEE/ACM Trans. Networking* 27 (5), 2164–2177.
- Strohmeier, D., Bekhti, Y., Hauelsen, J., Gramfort, A., 2016. The iterative reweighted mixed-norm estimate for spatio-temporal MEG/EEG source reconstruction. *IEEE Trans Med Imaging* 35 (10), 2218–2228.
- Strohmeier, D., Gramfort, A., Hauelsen, J., 2015. MEG/EEG source imaging with a non-convex penalty in the time-frequency domain. In: Pattern Recognition in NeuroImaging (PRNI), 2015 International Workshop on. IEEE, pp. 21–24.
- Stuart, A.M., 2010. Inverse problems: a bayesian perspective. *Acta Numerica* 19, 451–559.
- Sun, Y., Babu, P., Palomar, D.P., 2016. Robust estimation of structured covariance matrix for heavy-tailed elliptical distributions. *IEEE Trans. Signal Process.* 64 (14), 3576–3590.
- Sun, Y., Babu, P., Palomar, D.P., 2017. Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Trans. Signal Process.* 65 (3), 794–816.
- Tipping, M.E., 2000. The relevance vector machine. In: *Advances in Neural Information Processing Systems*, pp. 652–658.
- Tipping, M.E., 2001. Sparse bayesian learning and the relevance vector machine. *Journal of Machine Learning Research* 1 (Jun), 211–244.
- Trujillo-Barreto, N.J., Aubert-Vázquez, E., Valdés-Sosa, P.A., 2004. Bayesian model averaging in EEG/MEG imaging. *Neuroimage* 21 (4), 1300–1319.
- Tsiligkaridis, T., Hero, A.O., 2013. Covariance estimation in high dimensions via kronecker product expansions. *IEEE Trans. Signal Process.* 61 (21), 5347–5360.
- Tsiligkaridis, T., Hero III, A.O., Zhou, S., 2013. On convergence of kronecker graphical lasso algorithms. *IEEE Trans. Signal Process.* 61 (7), 1743–1755.
- Villani, C., 2008. Optimal transport: Old and new, 338. Springer Science & Business Media.
- Vishnoi, N.K., 2018. Geodesic convex optimization: differentiation on manifolds, geodesics, and convexity. arXiv preprint arXiv:1806.06373.
- Wei, H., Jafarian, A., Zeidman, P., Litvak, V., Razi, A., Hu, D., Friston, K.J., 2020. Bayesian fusion and multimodal DCM for EEG and fMRI. *Neuroimage* 211, 116595.
- Werner, K., Jansson, M., Stoica, P., 2008. On estimation of covariance matrices with kronecker product structure. *IEEE Trans. Signal Process.* 56 (2), 478–491.
- Wiesel, A., Zhang, T., et al., 2015. Structured robust covariance estimation. *Foundations and Trends® in Signal Processing* 8 (3), 127–216.
- Wipf, D., Nagarajan, S., 2009. A unified bayesian framework for MEG/EEG source imaging. *Neuroimage* 44 (3), 947–966.
- Wipf, D., Nagarajan, S., 2010. Iterative reweighted ℓ_1 and ℓ_2 methods for finding sparse solutions. *IEEE J Sel Top Signal Process* 4 (2), 317–329.
- Wipf, D.P., Owen, J.P., Attias, H.T., Sekihara, K., Nagarajan, S.S., 2010. Robust bayesian estimation of the location, orientation, and time course of multiple correlated neural sources using MEG. *Neuroimage* 49 (1), 641–655.
- Wipf, D.P., Rao, B.D., 2004. Sparse bayesian learning for basis selection. *IEEE Trans. Signal Process.* 52 (8), 2153–2164.
- Wipf, D.P., Rao, B.D., 2007. An empirical bayesian strategy for solving the simultaneous sparse approximation problem. *IEEE Trans. Signal Process.* 55 (7), 3704–3716.
- Wipf, D.P., Rao, B.D., Nagarajan, S., 2011. Latent variable bayesian models for promoting sparsity. *IEEE Trans. Inf. Theory* 57 (9), 6236–6255.
- Wu, C.J., 1983. On the convergence properties of the EM algorithm. *The Annals of Statistics* 95–103.
- Wu, T.T., Lange, K., et al., 2010. The MM alternative to EM. *Statistical Science* 25 (4), 492–505.
- Wu, W., Nagarajan, S., Chen, Z., 2016. Bayesian machine learning: EEG\MEG signal processing measurements. *IEEE Signal Process Mag* 33 (1), 14–36.
- Wu, Y., Wipf, D.P., 2012. Dual-space analysis of the sparse linear model. In: *Advances in Neural Information Processing Systems*, pp. 1745–1753.
- Yuille, A.L., Rangarajan, A., 2003. The concave-convex procedure. *Neural Comput* 15 (4), 915–936.
- Zadeh, P., Hosseini, R., Sra, S., 2016. Geometric mean metric learning. In: *International Conference on Machine Learning*, pp. 2464–2471.
- Zhang, Z., Rao, B.D., 2011. Sparse signal recovery with temporally correlated source vectors using sparse bayesian learning. *IEEE J Sel Top Signal Process* 5 (5), 912–926.
- Zoubir, A.M., Koivunen, V., Ollila, E., Muma, M., 2018. Robust statistics for signal processing. Cambridge University Press.
- Zumer, J.M., Attias, H.T., Sekihara, K., Nagarajan, S.S., 2007. A probabilistic algorithm integrating source localization and noise suppression for MEG and EEG data. *Neuroimage* 37 (1), 102–115.