



**HAL**  
open science

# LEMUR: Latent EM Unsupervised Regression for Sparse Inverse Problems

Pierre Barbault, Matthieu Kowalski, Charles Soussen

► **To cite this version:**

Pierre Barbault, Matthieu Kowalski, Charles Soussen. LEMUR: Latent EM Unsupervised Regression for Sparse Inverse Problems. 2024. hal-04542061

**HAL Id: hal-04542061**

**<https://hal.science/hal-04542061>**

Preprint submitted on 11 Apr 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# LEMUR: Latent EM Unsupervised Regression for Sparse Inverse Problems

Pierre Barbault, Matthieu Kowalski and Charles Soussen

**Abstract**—Most methods for sparse signal recovery require to set one or several hyperparameters. We propose an unsupervised method to estimate the parameters of a Bernoulli-Gaussian (BG) model describing sparse signals. The proposed method is first derived for denoising problems, based on a maximum likelihood (ML) approach. Then, an extension to general inverse problems is achieved through a latent variable formulation. Two expectation-maximization (EM) algorithms are then proposed to estimate the signal together with the BG model parameters. Combining these two approaches leads to the proposed LEMUR algorithm. All proposed algorithms are then evaluated on extensive simulations in terms of ability to recover the parameters and provide accurate sparse signal estimates.

**Index Terms**—Bernoulli-Gaussian, Expectation-Maximization, Sparsity, Hyperparameter estimation.

## I. INTRODUCTION

Linear inverse problems are prevalent in various scientific and engineering fields, encompassing scenarios where the objective is to reconstruct an unknown signal or image from noisy measurements. Many signals of interest, such as natural images and audio signals, exhibit sparse representation, which can be exploited for reconstruction from incomplete and noisy measurements. Sparsity is typically achieved using a linear transform or dictionary that allows the signal to be represented with sparse synthesis coefficients. In the following, we assume, without loss of generality, that the signal of interest is sparse with respect to the canonical basis.

The relation between the noisy measurement vector  $\mathbf{y} \in \mathbb{R}^M$  and the unknown sparse signal  $\mathbf{x} \in \mathbb{R}^N$  is given by the linear operator  $\mathbf{H} \in \mathbb{R}^{M \times N}$ , with additive white Gaussian noise  $\mathbf{e}$  having variance  $\sigma_e^2$  and zero mean:

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{e}. \quad (1)$$

Regularized least-squares is a popular approach for solving sparse inverse problems. The Lasso [1] or Basis Pursuit Denoising [2] is a common approach which makes use of an  $\ell_1$  constraint or penalty. While  $\ell_1$  regularization is prevalent, an ideal measure of sparsity is the so-called  $\ell_0$ -“norm”, which counts the number of non-zero coefficients in  $\mathbf{x}$ . The  $\ell_0$ -regularized least-squares problems can then be solved using various methods including proximal descent algorithms [3] such as the Fast Iterative Shrinkage/Thresholding algorithm (FISTA) [4]. These methods were initially developed for convex regularizers. However, they can be generalized to

the non-convex case as long as the proximal operator can be efficiently computed. For the  $\ell_0$  regularized problem, the Iterative Hard Thresholding algorithm (IHT) [5] can be used. Alternatives to IHT are greedy algorithms such as Orthogonal Matching Pursuit [6] and Single Best Replacement [7]. While these methods are effective in finding a sub-optimal solution, exact solvers based mixed integer programming have recently been proposed [8], [9] for moderate size problems, together with efficient branch and bound implementations.

The main shortcoming of these methods is the need to tune at least one hyperparameter, balancing the data-fidelity and sparsity-promoting terms. For the Lasso, the hyperparameter can be chosen using the Stein Unbiased Risk Estimator (SURE) in the denoising case [10], [11] by providing an unbiased estimate of the mean squared error between the estimated signal and the actual signal. This approach has been extended to inverse problems using the SUGAR [12] or Generalized-SURE (GSURE) [13]. However, SURE cannot be applied to  $\ell_0$  regularization as it is not weakly differentiable. In the denoising case, the SURE-based Coordinate-wise RELaxed (SCORE) [14] has been proposed, but its extension to inverse problems is not straightforward.

In the Bayesian framework, regularized least-squares can be interpreted as a Maximum *a Posteriori* (MAP) approach, where the regularization term is the neg-likelihood of a prior distribution over the signal coefficients. Pereyra *et al.* [15] proposed to use a Gamma prior on the hyperparameter when the negative log-likelihood of the prior corresponds to a 1-homogeneous regularizer. This method can be efficiently optimized. Among Bayesian methods for sparsity, Sparse Bayesian Learning (SBL) is a popular approach that uses a Gaussian prior on the signal coefficients with a free variance for each coefficient. Variance estimation is achieved through marginalization with respect to the signal  $\mathbf{x}$ , leading to sparse estimation. A related model proposed by Calvetti and Somersalo [16] uses a Gamma prior on the variances.

One advantage of the Bayesian approach is the ability to use hierarchical Bayes models, where a hyperprior can be chosen for the hyperparameters. Markov Chain Monte Carlo (MCMC) methods with hierarchical Bayes models provide a fully unsupervised approach for hyperparameter selection. While these methods can be computationally expensive, they can accurately estimate the hyperparameters.

The Bernoulli-Gaussian (BG) model is a popular choice for modeling sparsity in signal processing and is strongly related to the  $\ell_0$  regularized problem through the MAP estimator [7]. In [17], the authors proposed a stochastic-EM procedure to estimate the sparse solution of a deconvolution problem. However, a classical EM procedure becomes intractable because of

P. Barbault and C. Soussen are with L2S, Université Paris-Saclay, CNRS, CentraleSupélec, Gif-sur-Yvette, France.

M. Kowalski is with Inria, Université Paris-Saclay, CNRS, Laboratoire Interdisciplinaire des Sciences du Numériques, Gif-sur-Yvette, France

This work was supported by the French National Agency for Research through the BMWs project (ANR-20-CE45-0018).

the combinatorial nature of the related optimization problem. Following the EM approach developed in [18] that relies on an additional latent variable, [19] proposed an iterative procedure to estimate the parameters of a BG model when the matrix  $\mathbf{H}$  is a union of unitary dictionaries. However, this approach leads to a non-convergent estimate even in the denoising case, as shown in [20], when using the chosen prior on the hyperparameters. Moreover, preliminary experiments have shown that this method is not robust for the general inverse problem. MCMC approaches for the BG sparse model do not scale well in practice [21], although recent efforts have been made to improve their efficiency [22], [23].

**Contributions and outline of the paper.** This article is an extension of our conference paper [24]. The primary purpose is to provide an estimation of the BG parameters in the inverse problem framework together with the sparse signal reconstruction. On the contrary, only the parameters were estimated in [24]. Here, we aim to provide a fully unsupervised iterative/shrinkage thresholding algorithm. In Section II, we introduce the statistical models and derive the expression of likelihood functionals. Section III is dedicated to the denoising case. We first propose a joint MAP/Maximum Likelihood method to estimate both sparse signal and hyperparameters. Additionally, we derive an EM algorithm, which is initialized using the method of moments. In Section IV, we elaborate the extension to general inverse problems using the latent variable formulation proposed in [18]. We explore the choice of hidden variables and propose two distinct EM approaches. Finally, Section V includes a comprehensive experimental evaluation of the proposed approaches.

## II. BG MODEL

The sparse source signal  $\mathbf{x}$  is modeled using the Bernoulli-Gaussian (BG) process. Assuming that the signal coefficients  $x_n$  are independent and identically distributed,  $x_n$  is modeled as a mixture of a Gaussian distribution and the Dirac  $\delta$  distribution:

$$x_n \sim \mathcal{BG}(p, \sigma_x^2) = p\mathcal{N}(0, \sigma_x^2) + (1-p)\delta(x_n) \quad (2)$$

where  $p$  and  $\sigma_x^2$  refer to the sparsity level and the variance of the nonzero entries, respectively. Alternatively, one may introduce the sequence of binary variables  $s_n$  such that

$$s_n = \begin{cases} 1 & \text{if } x_n \neq 0, \\ 0 & \text{if } x_n = 0. \end{cases} \quad (3)$$

$s_n$  follows the Bernoulli distribution of parameter  $p$ , and the conditional distribution of  $x_n$  given  $s_n$  reads  $(x_n|s_n) \sim \mathcal{N}(0, \sigma_x^2)$  if  $s_n = 1$  and  $(x_n|s_n) = 0$ .

In Eq. (1), the independent noise vector  $\mathbf{e}$  is assumed to be white and Gaussian with variance  $\sigma_e^2$ , thus we have:

$$\mathbf{y}|\mathbf{x} \sim \mathcal{N}(\mathbf{H}\mathbf{x}, \sigma_e^2 \mathbf{I}_M) \quad (4)$$

where  $\mathbf{I}_M$  is the  $M \times M$  identity matrix.

Due to the presence of the Dirac distribution in Eq. (2), one cannot directly define the maximum *a posteriori* (MAP) estimator for  $\mathbf{x}$  using this equation. However, by introducing the indicator variable  $s_n$  in Eq. (3) and observing that

$p(x_n \neq 0|s_n = 1) = 1$  and  $p(x_n = 0|s_n = 0) = 1$ , we can establish a joint-MAP formulation in terms of  $(\mathbf{s}, \mathbf{x})$ . Ultimately, we demonstrate that this joint estimator relies solely on  $\mathbf{x}$ , and we denote the latter estimator as the MAP estimator for  $\mathbf{x}$ . Denoting the set of model hyperparameters as

$$\theta = (p, \sigma_x^2, \sigma_e^2), \quad (5)$$

this estimator is expressed hereafter.

**Proposition 1.** *Let  $\mathbf{x}$  follow the BG model Eq. (2) and  $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{e}$ . The MAP estimator of  $\mathbf{x}$  is the solution of the following optimization problem:*

$$\underset{\mathbf{x}}{\operatorname{argmin}} \frac{1}{2\sigma_e^2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \frac{\|\mathbf{x}\|_2^2}{2\sigma_x^2} + \lambda \|\mathbf{x}\|_0 + C(\sigma_e^2, p)$$

with

$$\lambda(\sigma_x^2, p) = \log \left( \sqrt{2\pi\sigma_x^2} \frac{1-p}{p} \right)$$

and

$$C(\sigma_e^2, p) = N \log \left( \frac{\sqrt{2\pi\sigma_e^2}}{1-p} \right)$$

*Proof.* The proof is given in Appendix A.  $\square$

The expression of  $C(\sigma_e^2, p)$  is given for reasons that will become clear in the next sections devoted to estimating the parameters. As far as we know, it is the first time that this MAP estimator is properly derived. In [20], this MAP was given in the denoising case with no proof. One can notice that this posterior slightly differs from the one given in [7]. The latter cost function is similar with  $\lambda = \log \left( \frac{1-p}{p} \right)$ . This difference comes from the chosen model to induce the BG random variable  $\mathbf{x}$ . In [7],  $\mathbf{x}$  is written as the product of two *independent* random variables  $\mathbf{s}$  and  $\mathbf{r}$ , which model the support and the value of the coefficients, respectively. Hence, the derivation is based on the factorization  $p(\mathbf{s}, \mathbf{r}) = p(\mathbf{s})p(\mathbf{r})$ . Hence, the related MAP criterion differs from the one stated in Proposition 1.

## III. PARAMETER ESTIMATION IN THE DENOISING CASE

This section deals with parameter estimation in the *denoising* case. This case is modeled by  $\mathbf{H} = \mathbf{I}_M$ , hence

$$\mathbf{y} = \mathbf{x} + \mathbf{e}. \quad (6)$$

Here, the distribution of  $\mathbf{y}$  is a simple mixture of two centered Gaussians with variances  $\sigma_e^2$  and  $\sigma_e^2 + \sigma_x^2$ , and respective weights  $(1-p)$  and  $p$ . In other words,  $y_n$  are i.i.d. random variables such that

$$y_n \sim p\mathcal{N}(0, \sigma_x^2 + \sigma_e^2) + (1-p)\mathcal{N}(0, \sigma_e^2). \quad (7)$$

In Section III-A, we derive an EM procedure for ML estimation of parameters  $\theta$  and then propose two methods of reconstruction of the sparse signal  $\mathbf{x}$  knowing  $\theta$ . In Section III-B, we propose a joint estimator of  $\mathbf{x}$  and  $\theta$ . We further derive the method of moments, which provides an initial estimate of  $\theta$  (Section III-C).

---

**Algorithm 1:** Denoising problem: EM algorithm for parameter estimation [24].

---

**Result:**  $p, \sigma_x^2, \sigma_e^2$

**Input:**  $t = 0, p^{(0)}, (\sigma_x^2)^{(0)}, (\sigma_e^2)^{(0)}$

**while not converged do**

update  $\phi_n^{(t)}$  with Eq. (8);  
 update  $\mu^{(t)}$  and  $\nu^{(t)}$  with Eq. (9);  
 update  $p^{(t+1)}, (\sigma_x^2)^{(t+1)}, (\sigma_e^2)^{(t+1)}$  with Eqs. (10)  
 to (12);  
 $t = t + 1$ ;

**end**

---

### A. Maximum Likelihood estimation

1) *Estimation of hyperparameters:* ML estimation amounts to maximizing  $p(\mathbf{y}|\theta)$  with respect to  $\theta$ .  $\mathbf{y}$  being a mixture of two Gaussians, it is a typical case where the EM procedure can be derived [25]. In [24], we derived an EM algorithm using  $\mathbf{x}$  as hidden variable. We briefly summarize this procedure for the sake of completeness.

Denoting by  $\theta^{(t)} = \{p^{(t)}, (\sigma_x^2)^{(t)}, (\sigma_e^2)^{(t)}\}$  the estimated parameters at iteration  $t$ , let

$$\begin{aligned} \phi_n^{(t)} &= p(s_n = 1|y_n, \theta^{(t)}) \\ &= \frac{\frac{p}{\sqrt{\sigma_x^2 + \sigma_e^2}} e^{-\frac{y_n^2}{2(\sigma_x^2 + \sigma_e^2)}}}{\frac{p}{\sqrt{\sigma_x^2 + \sigma_e^2}} e^{-\frac{y_n^2}{2(\sigma_x^2 + \sigma_e^2)}} + \frac{1-p}{\sqrt{\sigma_e^2}} e^{-\frac{y_n^2}{2\sigma_e^2}}} \end{aligned} \quad (8)$$

(where the latter equation results from Bayes' rule) and

$$\mu^{(t)} = \frac{(\sigma_x^2)^{(t)}}{(\sigma_x^2)^{(t)} + (\sigma_e^2)^{(t)}}, \quad \nu^{(t)} = \frac{(\sigma_x^2)^{(t)}(\sigma_e^2)^{(t)}}{(\sigma_x^2)^{(t)} + (\sigma_e^2)^{(t)}}. \quad (9)$$

According to [24], the expectation  $\mathbb{E}_{\mathbf{x}|\mathbf{y}, \theta^{(t)}}[\log p(\mathbf{y}, \mathbf{x}|\theta)]$  reaches a maximum value when  $\theta = \theta^{(t+1)}$ , with

$$p^{(t+1)} = \frac{1}{N} \sum_{n=1}^N \phi_n^{(t)} \quad (10)$$

$$(\sigma_x^2)^{(t+1)} = \nu^{(t)} + \frac{(\mu^{(t)})^2}{N p^{(t+1)}} \sum_{n=1}^N y_n^2 \phi_n^{(t)} \quad (11)$$

$$(\sigma_e^2)^{(t+1)} = \frac{1}{N} \sum_{n=1}^N y_n^2 - \frac{2\mu^{(t)}}{N} \sum_{n=1}^N y_n^2 \phi_n^{(t)} + p^{(t+1)}(\sigma_x^2)^{(t+1)}. \quad (12)$$

The resulting EM algorithm is summarized in Alg. 1. Since EM algorithms are local maximization schemes, the choice of initial solution is an important issue. This point will be further discussed in Section III-C.

2) *Estimation of the sparse signal:* Once the hyperparameters  $\theta$  have been retrieved using Algorithm 1, one can reconstruct the signal  $\mathbf{x}$  using two alternative strategies.

The first strategy amounts to computing the posterior mean estimate  $\hat{\mathbf{x}} = \mathbb{E}_{\mathbf{x}|\mathbf{y}, \theta}[\mathbf{x}]$ , which is also the Bayesian estimator minimizing the Mean Squared Error  $\mathbb{E}_{\mathbf{x}|\mathbf{y}, \theta}[\|\hat{\mathbf{x}} - \mathbf{x}\|^2]$  (with respect to  $\mathbf{x}$ ). For denoising problems, the posterior mean is

defined componentwise. Using the total probability rule, we get

$$\begin{aligned} \mathbb{E}_{x_n|y_n, \theta}[x_n] &= p(s_n = 1|y_n, \theta) \times \mathbb{E}_{x_n|y_n, s_n=1, \theta}[x_n] \\ &\quad + p(s_n = 0|y_n, \theta) \times 0 \\ &= \phi_n \mu y_n \end{aligned} \quad (13)$$

where  $\phi_n$  and  $\mu$  refer to the estimates in Eq. (8)-Eq. (9) computed with parameter  $\theta$  obtained at the final iteration of Alg. 1 (for more details, see [24]).

The second strategy is based on the MAP reconstruction of the signal support, that is

$$s_n = 1 \Leftrightarrow p(s_n = 1|y_n, \theta) > p(s_n = 0|y_n, \theta). \quad (14)$$

where  $p(s_n = 1|y_n, \theta)$  is given in Eq. (8) and  $p(s_n = 0|y_n, \theta) = 1 - p(s_n = 1|y_n, \theta)$ . Up to a few rearrangements, we find that  $s_n = 1$  if and only if

$$y_n^2 > 2\sigma_e^2 \frac{\sigma_e^2 + \sigma_x^2}{\sigma_x^2} \log \left( \frac{1-p}{p} \sqrt{1 + \frac{\sigma_x^2}{\sigma_e^2}} \right). \quad (15)$$

Once the support is reconstructed using Eq. (15), we propose to compute the mean of the posterior distribution of  $\mathbf{x}$  conditionally to  $\mathbf{s}$  and  $\theta$ . This estimator can be easily computed given that the latter distribution is Gaussian:

$$\mathbb{E}_{x_n|y_n, s_n, \theta}[x_n] = \begin{cases} \frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} y_n & \text{if } s_n = 1, \\ 0 & \text{if } s_n = 0. \end{cases} \quad (16)$$

The reader can refer to [26] for a study and comparison of both estimates Eq. (13) and Eq. (16) in denoising problems.

### B. Joint estimation of $\mathbf{x}$ and $\theta$

As one is usually primarily interested in the estimation of signal  $\mathbf{x}$ , it seems natural to investigate the following joint estimation problem

$$(\hat{\mathbf{x}}, \hat{\theta}) = \underset{\mathbf{x}, \theta}{\operatorname{argmin}} J(\mathbf{x}, \theta) \quad (17)$$

where

$$J(\mathbf{x}, \theta) = \frac{1}{2\sigma_e^2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \frac{\|\mathbf{x}\|_2^2}{2\sigma_x^2} + \lambda(\sigma_e^2, p) \|\mathbf{x}\|_0 + C(\sigma_e^2, p) \quad (18)$$

with  $\lambda(\sigma_e^2, p)$  and  $C(\sigma_e^2, p)$  defined as in Proposition 1. This joint optimization problem can be interpreted as a joint MAP estimator of  $(\mathbf{x}, \theta)$  using a non-informative (uniform) prior on  $\theta$  following the approach of Proposition 1 where the support  $\mathbf{s}$  is inferred together with  $\mathbf{x}$ .

We first notice that for a given set of parameters  $\theta$ , the minimization of  $J$  with respect to  $\mathbf{x}$  leads to

$$\forall n, \hat{x}_n = \begin{cases} \frac{\sigma_x^2}{\sigma_x^2 + \sigma_e^2} y_n & \text{if } y_n^2 > T(p, \sigma_x^2, \sigma_e^2), \\ 0 & \text{otherwise} \end{cases} \quad (19)$$

where the threshold  $T(p, \sigma_x^2, \sigma_e^2)$  is given by

$$T(p, \sigma_x^2, \sigma_e^2) = 2\sigma_e^2 \frac{\sigma_x^2 + \sigma_e^2}{\sigma_x^2} \log \left( \sqrt{2\pi\sigma_x^2} \frac{1-p}{p} \right). \quad (20)$$

By substituting  $\hat{\mathbf{x}}$  into Eq. (17), it has been shown [20] that the functional to minimize can be rewritten as a function

depending on the threshold  $T$  and the parameters  $\theta$ . The problem reduces to minimizing with respect to  $\theta$  and  $T$  the following functional

$$\begin{aligned} \tilde{J}(\theta, T) = & \frac{1}{2\sigma_e^2} \sum_{\{n: y_n^2 > T\}} (y_n - x_n)^2 + \frac{1}{2\sigma_x^2} \sum_{\{n: y_n^2 > T\}} x_n^2 \\ & + \lambda(\sigma_e^2, p) \sum_{\{n: y_n^2 > T\}} 1 + C(\sigma_e^2, p) \end{aligned} \quad (21)$$

subject to

$$T = 2\sigma_e^2 \frac{\sigma_x^2 + \sigma_e^2}{\sigma_x^2} \log \left( \sqrt{2\pi\sigma_x^2} \frac{1-p}{p} \right). \quad (22)$$

For a given threshold value  $T$ , the minimization with respect to the parameters  $\theta$  leads to:

$$\hat{p} = \frac{1}{N} \sum_{\{n: y_n^2 > T\}} 1 \quad (23)$$

$$\hat{\sigma}_x^2 = \frac{\sum_{\{n: y_n^2 > T\}} y_n^2}{(1+\gamma)^2 \sum_{\{n: y_n^2 > T\}} 1} \quad (24)$$

$$\hat{\sigma}_e^2 = \gamma \sigma_x^2 \quad (25)$$

where  $\gamma$  is the smallest root of

$$\gamma^2(1+\rho) + \gamma \left( 2 - \rho \frac{N}{\sum_{\{n: y_n^2 > T\}} 1} \right) + 1 = 0 \quad (26)$$

with

$$\rho = \frac{\sum_{\{n: y_n^2 > T\}} y_n^2}{\sum_{\{n: y_n^2 \leq T\}} y_n^2}. \quad (27)$$

To minimize Eq. (21), one can remark that, because we are in a finite-dimensional discrete setting, the quantities  $\sum_{n: y_n^2 > T} 1$ ,  $\sum_{n: y_n^2 \leq T} y_n^2$  and  $\sum_{n: y_n^2 > T} y_n^2$  have a piecewise constant dependency with respect to  $T$ . Indeed, denoting by  $\tilde{\mathbf{y}}$  the vector  $\mathbf{y}$  sorted in decreasing order, then the previous quantities remain constant on  $[\tilde{y}_n^2, \tilde{y}_{n+1}^2[$ . Overall, there are only  $N$  possible values for  $\hat{p}$ ,  $\hat{\sigma}_x^2$  and  $\hat{\sigma}_e^2$ . Hence, one just needs to compute these  $N$  possible parameter values, estimate the corresponding  $\hat{\mathbf{x}}$  and compute the related value of  $J$  in Eq. (17). Then, the parameters yielding the minimum value of  $J$  are selected. This procedure is summarized in Alg. 2.

It is worth noticing that when  $N$  tends to infinity, one recovers the expression of the parameters given in [20]. One of the main results in [20] is that the estimation Eq. (23) is necessarily biased and non consistent. Moreover, it turns out that this method reduces to a Classification-EM (CEM) procedure [27] applied to a mixture of two centered Gaussians, also known to produce a biased estimator.

### C. Estimation by the method of moments

Algs. 1 and 2 presented above are local optimization procedures. Therefore, they are sensitive to the initial  $\theta$  value.

---

### Algorithm 2: Joint Estimation of $\mathbf{x}$ and $\theta$

---

**Result:**  $\hat{\mathbf{x}}, \hat{\theta} = (\hat{p}, \hat{\sigma}_x^2, \hat{\sigma}_e^2)$

**Input:**  $\mathbf{y}$

**for** all  $k \in \{1, \dots, N\}$  **do**

    Compute  $\theta^k = (\hat{p}, \hat{\sigma}_x^2, \hat{\sigma}_e^2)$  according to Eqs. (23) to (25) with  $T = y_k^2$ ;

    Compute  $\mathbf{x}^k$  according to Eq. (19) with  $T(\hat{p}, \hat{\sigma}_x^2, \hat{\sigma}_e^2)$ ;

    Compute  $J^k = J(\mathbf{x}^k, \theta^k)$ ;

**end**

Set  $\kappa = \operatorname{argmin}_k J^k$ ;

Set  $\hat{\theta} = \theta^\kappa$ ;

Set  $\hat{\mathbf{x}} = \mathbf{x}^\kappa$ ;

---

Hereafter, we propose an estimator of  $\theta$  based the method of moments that will be further used to initialize Algs. 1 and 2.

**Theorem 1.** Let  $\mathbf{y} = \{y_1, \dots, y_N\}$  such that, for all  $n$ ,  $y_n$  are i.i.d. distributed according to Eq. (7). Define the following quantities:

$$m_2 = \frac{1}{N} \sum_{n=1}^N y_n^2 \quad m_4 = \frac{1}{3N} \sum_{n=1}^N y_n^4 \quad m_6 = \frac{1}{15N} \sum_{n=1}^N y_n^6$$

$$A = \frac{m_6 - m_2 m_4}{2(m_4 - m_2^2)} \quad B = \frac{m_2 m_6 - m_4^2}{m_4 - m_2^2}.$$

Then, the method of moments applied to  $\theta = (p, \sigma_x^2, \sigma_e^2)$  yields the following estimator:

$$\begin{aligned} \hat{\sigma}_e^2 &= A - \sqrt{A^2 - B} \\ \hat{\sigma}_x^2 &= \frac{m_4 - (\hat{\sigma}_e^2)^2}{m_2 - \hat{\sigma}_e^2} - 2\hat{\sigma}_e^2 \\ \hat{p} &= \frac{m_2 - \hat{\sigma}_e^2}{\hat{\sigma}_x^2} \end{aligned}$$

Moreover, this estimator is consistent.

*Proof.* The derivation of the method of moments is given in Appendix B. Consistency comes from the continuity of the functions used to build the estimator [28, Chp. 9].  $\square$

Although the estimator is consistent, the method of moments is known to yield biased estimators. In practice, it could happen that the conditions  $m_2^2 < m_4$  and  $m_4^2 < m_2 m_6$  (which imply that  $B > 0$  and then  $\hat{\sigma}_e^2$  given above is well-defined, see Appendix B) are not met for specific observations  $\mathbf{y}$ . Therefore, we apply an empirical procedure which iteratively removes one entry at a time in  $\mathbf{y}$  so as to increase as much as possible the value of

$$\min \{0; m_4 - m_2^2\} + \min \{0; m_2 m_6 - m_4^2\}. \quad (28)$$

This process is stopped when  $m_2^2 < m_4$  and  $m_4^2 < m_2 m_6$ .

#### IV. EXTENSION TO GENERAL INVERSE PROBLEMS

In this section, we extend the hyperparameter estimation procedure of Section III to general inverse problems (1). The case where  $\mathbf{H}$  is orthonormal boils down to a denoising problem with vector  $\mathbf{H}^T \mathbf{y}$  as input, since  $\|\mathbf{y} - \mathbf{H}\mathbf{x}\|^2 = \|\mathbf{H}^T \mathbf{y} - \mathbf{x}\|^2$ . When  $\mathbf{H}$  is not orthonormal, the distribution of  $\mathbf{y}$  is no longer a mixture of two Gaussians. Then, the standard EM procedure cannot be applied anymore. Indeed, the computation and optimization of the likelihood function is known to be untractable, see [17].

We resort to the re-parametrization proposed in [18], involving a so-called latent variable  $\mathbf{z}$ . Eq. (1) rewrites:

$$\mathbf{y} = \mathbf{H}\mathbf{z} + \mathbf{b} \quad (29)$$

$$\mathbf{z} = \mathbf{x} + \mathbf{n} \quad (30)$$

where  $\mathbf{n} \sim \mathcal{N}(0, \sigma_n^2 \mathbf{I}_N)$  and  $\mathbf{b} \sim \mathcal{N}(0, \Gamma_b)$  are independent random vectors. Since  $\mathbf{e} = \mathbf{b} + \mathbf{H}\mathbf{n}$ , the covariance of the ‘‘outer’’ noise  $\mathbf{b}$  must satisfy

$$\Gamma_b = \sigma_e^2 \mathbf{I}_M - \sigma_n^2 \mathbf{H}\mathbf{H}^T \quad (31)$$

with  $\sigma_n^2 < \frac{\sigma_e^2}{\|\mathbf{H}\mathbf{H}^T\|}$  to yield a non-degenerate distribution, where  $\|\cdot\|$  refers to the spectral norm of a matrix.

One can notice from (29) that the distribution of  $\mathbf{y}|\mathbf{z}, \mathbf{x}$  is independent of  $\mathbf{x}$ . Moreover, from Bayes’ rule, we have  $p(\mathbf{x}|\mathbf{y}, \mathbf{z}) \propto p(\mathbf{y}|\mathbf{x}, \mathbf{z})p(\mathbf{x}|\mathbf{z})$  (where  $\propto$  denotes proportional-ity). It follows that

$$p(\mathbf{y}|\mathbf{z}, \mathbf{x}) = p(\mathbf{y}|\mathbf{z}) \quad \text{and} \quad p(\mathbf{x}|\mathbf{y}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z}). \quad (32)$$

With the reparametrization (30), the parameters to be estimated are now  $(p, \sigma_x^2, \sigma_e^2, \sigma_n^2)$ . Hereafter, we will set

$$\sigma_n^2 = \alpha \sigma_e^2 \quad (33)$$

where  $\alpha$  is a predefined parameter satisfying  $\alpha < \frac{1}{\|\mathbf{H}\mathbf{H}^T\|}$ . Therefore, the remaining free parameters to be estimated are  $\theta = (p, \sigma_x^2, \sigma_n^2)$ .

Because

$$p(\mathbf{y}|\mathbf{z}, \mathbf{x}, \theta) \propto p(\mathbf{y}|\mathbf{z}, \theta) p(\mathbf{z}|\mathbf{x}, \theta) \quad (34)$$

and both probability density functions appearing in the latter equation are Gaussian, the posterior distribution of  $\mathbf{z}$  given  $(\mathbf{y}, \mathbf{x})$  is a multivariate Gaussian distribution written (see e.g., [28, Chap. 10])

$$\mathbf{z}|\mathbf{y}, \mathbf{x}, \theta \sim \mathcal{N}(\mu_z, \Gamma_z) \quad (35)$$

with

$$\mu_z = \mathbf{x} + \alpha \mathbf{H}^T (\mathbf{y} - \mathbf{H}\mathbf{x}) \quad (36)$$

$$\Gamma_z = \sigma_n^2 (\mathbf{I} - \alpha \mathbf{H}^T \mathbf{H}) \quad (37)$$

Notice that the computation of  $\mu_z$  reads as a gradient descent iteration on the  $\ell_2$ -loss  $\frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2$  with step  $\alpha$ .

We provide two Bayesian estimators. In Section IV-A, using a noninformative uniform prior on the parameters, we derive the joint MAP estimator of  $(\mathbf{x}, \theta)$ :

$$(\hat{\mathbf{x}}, \hat{\theta}) = \underset{\mathbf{x}, \theta}{\operatorname{argmax}} p(\mathbf{x}, \theta | \mathbf{y}) \quad (38)$$

In Section IV-B, we exhibit the joint MAP estimator of  $(\mathbf{z}, \theta)$ :

$$(\hat{\mathbf{z}}, \hat{\theta}) = \underset{\mathbf{z}, \theta}{\operatorname{argmax}} p(\mathbf{z}, \theta | \mathbf{y}) \quad (39)$$

already proposed in our preliminary conference publication [24]. In the denoising case (and by extension, when  $\mathbf{H}$  is orthogonal), the former yields a biased estimator of  $\theta$  (see Section III-B). The later reduces to the ML estimator described in Section III-A. It yields a consistent estimator of  $\theta$  [24].

##### A. Joint estimation of $(\mathbf{x}, \theta)$ using $\mathbf{z}$ as hidden variable

Adapting [18] to the BG model, let us derive the EM algorithm using  $\mathbf{z}$  as hidden variable.

a) *E-Step*: Exploiting the decoupling property in (32), the E-step of (38) reads:

$$\begin{aligned} Q(\mathbf{x}, \theta | \mathbf{x}^{(t)}, \theta^{(t)}) &= \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{x}, \theta, \mathbf{z} | \mathbf{y})] \\ &= \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{y} | \mathbf{z}, \theta) - \log p(\mathbf{z}, \mathbf{x}, \theta)] + C_1 \\ &= \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{y} | \mathbf{z}, \theta) - \log p(\mathbf{x}, \theta | \mathbf{z})] + C_2 \end{aligned} \quad (40)$$

where  $C_1$  and  $C_2$  do not depend on  $\mathbf{x}$  nor  $\theta$ .

Let us now express both terms within (40). According to Eq. (29), we have  $(\mathbf{y} | \mathbf{z}, \theta) \sim \mathcal{N}(\mathbf{H}\mathbf{z}, \Gamma_b)$ . Using elementary identities on the expectation of quadratic forms, we get

$$\begin{aligned} \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{y} | \mathbf{z}, \theta)] &= \frac{1}{2} \operatorname{Trace}[\Gamma_z^{(t)} \mathbf{H}^T \Gamma_b^{-1} \mathbf{H}] \\ &\quad + \frac{1}{2} (\mathbf{y} - \mathbf{H}\mathbf{z}^{(t)})^T \Gamma_b^{-1} (\mathbf{y} - \mathbf{H}\mathbf{z}^{(t)}) + \frac{1}{2} \log |2\pi \Gamma_b| \end{aligned} \quad (41)$$

where  $\mathbf{z}^{(t)} = \mathbf{x}^{(t)} + \alpha \mathbf{H}^T (\mathbf{y} - \mathbf{H}\mathbf{x}^{(t)})$  and  $\Gamma_z^{(t)} = (\sigma_n^2)^{(t)} (\mathbf{I} - \alpha \mathbf{H}^T \mathbf{H})$  are defined from Eqs. (36) and (37).

In order to derive the second term of (40), we need to express  $-\log p(\mathbf{x}, \theta | \mathbf{z})$ . Since  $\mathbf{z}$  reads as a noisy version of  $\mathbf{x}$  (see (30)), one can use Proposition 1 in the denoising case.  $-\log p(\mathbf{x}, \theta | \mathbf{z})$  thus reads

$$\frac{1}{2\sigma_n^2} \|\mathbf{z} - \mathbf{x}\|_2^2 + \frac{1}{2\sigma_x^2} \|\mathbf{x}\|_2^2 + \lambda(\sigma_x^2, p) \|\mathbf{x}\|_0 + C(\sigma_n^2, p). \quad (42)$$

It follows from Eq. (35) that

$$\begin{aligned} \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{x}, \theta | \mathbf{z})] &= \frac{\|\mathbf{z}^{(t)} - \mathbf{x}\|_2^2 + \operatorname{Trace}[\Gamma_z^{(t)}]}{2\sigma_n^2} \\ &\quad + \frac{1}{2\sigma_x^2} \|\mathbf{x}\|_2^2 + \lambda(\sigma_x^2, p) \|\mathbf{x}\|_0 + C(\sigma_n^2, p). \end{aligned} \quad (43)$$

b) *M-Step*: Summarizing, the criterion  $Q(\mathbf{x}, \theta | \mathbf{x}^{(t)}, \theta^{(t)})$  reads as the sum of (41) and (43), plus a constant independent of  $\mathbf{x}$  and  $\theta$ . One can notice that (43) depends on both  $\mathbf{x}$  and  $\theta$ , whereas (41) solely depends on  $\sigma_e^2$ . Indeed, all vectors and matrices within (41) do not depend on  $\theta$ , except for  $\Gamma_b = \sigma_e^2 (\mathbf{I} - \alpha \mathbf{H}\mathbf{H}^T)$ , see Eqs. (31) and (33).

Following these dependencies and taking into account (33), the minimization scheme is as follows:

- 1)  $(\mathbf{x}^{(t+1)}, \theta^{(t+1)}) = \underset{\mathbf{x}, \theta}{\operatorname{argmin}} \mathbb{E}_{\mathbf{z}|\mathbf{y}, \mathbf{x}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{z}, \mathbf{x}, \theta)]$
- 2)  $(\sigma_e^2)^{(t+1)} = \frac{(\sigma_n^2)^{(t+1)}}{\alpha}$

The first step is very similar to the joint denoising problem of [Section III-B](#) applied to  $\mathbf{z}^{(t)}$ , with the extra term  $\frac{\text{Trace}[\Gamma_z^{(t)}]}{2\sigma_n^2}$ . Hence, one can derive an algorithm very similar to [Alg. 2](#) to estimate jointly  $\mathbf{x}^{(t+1)}$  and  $\theta^{(t+1)}$ . The only difference is related to the constant  $\rho$  in [Eq. \(27\)](#), which must be set as follows:

$$\rho = \frac{\sum_{\{n:(z_n^{(t)})^2 > T\}} (z_n^{(t)})^2}{\sum_{\{n:(z_n^{(t)})^2 \leq T\}} (z_n^{(t)})^2 + \text{Trace}[\Gamma_z^{(t)}]}. \quad (44)$$

The whole EM procedure is summarized in [Alg. 3](#), where parameter  $\alpha$  is set to the limit value  $\frac{1}{\|\mathbf{H}\mathbf{H}^T\|}$ .

---

**Algorithm 3:** Joint estimation of  $(\mathbf{x}, \theta)$

---

**Result:**  $\hat{\mathbf{x}}, \hat{\theta}$

**Input:**  $t = 0, \mathbf{z}^{(t)} = \mathbf{0}, \mathbf{x}^{(t)} = \mathbf{0}, \alpha = \frac{1}{\|\mathbf{H}\mathbf{H}^T\|}$

**while** Not converged **do**

$\mathbf{z}^{(t)} = \mathbf{x}^{(t)} + \alpha\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x}^{(t)});$   
 $\Gamma_z^{(t)} = \sigma_n^2(\mathbf{I} - \alpha\mathbf{H}^T\mathbf{H});$   
 Estimate  $\mathbf{x}^{(t)}$  and  $\theta^{(t)}$  by applying [Alg. 2](#) with  $\mathbf{y} = \mathbf{z}^{(t)}$  and the modified expression (44) for  $\rho$ ;  
 $(\sigma_e^2)^{(t+1)} = \frac{(\sigma_n^2)^{(t+1)}}{\alpha};$   
 $t = t + 1;$

**end**

---

*B. Joint estimation of  $\mathbf{z}, \theta$  using  $\mathbf{x}$  as hidden variable*

We first briefly summarize the procedure described in our conference paper [\[24\]](#). We use  $\mathbf{z}$  as the variable to be estimated and consider  $\mathbf{x}$  as the hidden variable that is marginalized out. Using a uniform prior on  $\theta$  and exploiting the property [\(32\)](#), the E-step reads

$$\begin{aligned} \mathbb{E}_{\mathbf{x}|\mathbf{y}, \mathbf{z}^{(t)}, \theta^{(t)}} [-\log p(\mathbf{x}, \mathbf{z}, \theta|\mathbf{y})] &= -\log p(\mathbf{y}|\mathbf{z}) \\ &\quad - \mathbb{E}_{\mathbf{x}|\mathbf{z}^{(t)}, \theta^{(t)}} [\log p(\mathbf{z}, \mathbf{x}|\theta)] + C \end{aligned} \quad (45)$$

where  $C$  does not depend on  $\mathbf{z}$  nor  $\theta$ . Similar to [Section IV-A](#), we set  $\sigma_n^2 = \alpha\sigma_e^2$ . The resulting EM algorithm alternates between updates of  $\theta$  and  $\mathbf{z}$ , as detailed in [\[24\]](#). At iteration  $t$ ,  $\mathbf{z}$  is estimated according to

$$\mathbf{z}^{(t+1)} = \mathbf{x}^{(t)} + \alpha\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x}^{(t)}). \quad (46)$$

Then, the update of  $\theta$  reduces to the ML estimate of the denoising problem in  $\mathbf{z}^{(t)}$  discussed in [Section III-A](#).

The EM procedure is given in [Alg. 4](#), which provides estimates of the model parameters  $p$  and  $\sigma_x^2$  as well as the noise variance  $\sigma_n^2 = \alpha\sigma_e^2$ . By construction, the algorithm also provides an estimate for  $\mathbf{x}$ , which is the posterior mean of the denoising problem for  $\mathbf{z}$ .

The EM that uses  $\mathbf{z}$  as hidden variable is proved to be asymptotically biased [\[29\]](#) when it comes to parameter estimation but appears to be stable in practice, that is, not very sensitive to initial settings. On the other hand, the EM algorithm that uses  $\mathbf{x}$  as hidden variable is unbiased at the

---

**Algorithm 4:** Joint estimation of  $(\mathbf{z}, \theta)$

---

**Result:**  $\hat{\mathbf{z}}, \hat{\theta}$

**Input:**  $t = 0, \mathbf{z}^{(t)} = \mathbf{0}, \mathbf{x}^{(t)} = \mathbf{0}, \alpha = \frac{1}{\|\mathbf{H}\mathbf{H}^T\|}$

**while** not converged **do**

$\mathbf{z}^{(t+1)} = \mathbf{x}^{(t)} + \alpha\mathbf{H}^T(\mathbf{y} - \mathbf{H}\mathbf{x}^{(t)});$   
 Estimate  $\theta^{(t+1)}$  using [Alg. 1](#) with  $\mathbf{y} = \mathbf{z}^{(t+1)};$   
 Estimate  $\mathbf{x}^{(t+1)}$  using [Eq. \(13\)](#) with  $\mathbf{y} = \mathbf{z}^{(t+1)};$   
 $t = t + 1;$

**end**

---

optimum but very sensitive to initialization. Therefore, we propose the complete **Latent EM Unsupervised Regression (LEMUR)** procedure, where [Alg. 4](#) is initialized by [Alg. 3](#).

## V. EXPERIMENTAL STUDY

In this section, we present an assessment of the proposed methods for parameter estimation and sparse signal recovery. To do so, we generate synthetic BG signals  $\mathbf{x}$  following model [Eq. \(2\)](#) for varying degrees of sparsity controlled by the Bernoulli parameter  $p \in \{0.01, 0.05, 0.1\}$ . The variance of the non-zero coefficients is set to  $\sigma_x^2 = 1$ . We then generate several observations  $\mathbf{y}$  following the direct model [\(1\)](#) using various matrix operators  $\mathbf{H}$ . In addition to the simple denoising model, where  $\mathbf{H} = \mathbf{I}_M$  with  $M = 900$ , we generate correlated Gaussian random matrices of size  $300 \times 900$ . The covariance matrix  $\mathbf{C}$  is such that

$$C_{ij} = w^{|i-j|} \quad (47)$$

where  $w \in (0, 1)$  is the correlation degree. The matrices  $\mathbf{H}$  are then normalized so as to have unit norm columns. When  $w = 0$ , the covariance matrix is simply the identity matrix, which makes of  $\mathbf{H}$  a pure Gaussian random matrix as found in compressive sensing. For each  $\mathbf{y}$ , we generate a white Gaussian noise vector with variance

$$\sigma_e^2 = \|\mathbf{H}\mathbf{x}\| \times 10^{-\text{iSNR}/10} \quad (48)$$

where iSNR stands for the input Signal-to-Noise Ratio.

We monitor the estimation of both parameters and sources. Denoting by  $(\hat{p}, \hat{\sigma}_e^2, \hat{\sigma}_x^2)$  the estimated hyperparameters, the metric used to quantify the quality of parameter estimation is the absolute error for  $p \in \{0.01, 0.05, 0.1\}$  and  $\sigma_x^2 = 1$ :

$$\text{error}_p = |p - \hat{p}| \quad \text{error}_{\sigma_x^2} = |\sigma_x^2 - \hat{\sigma}_x^2|. \quad (49)$$

The value of  $\sigma_e^2$  being fixed to reach the target iSNR, as it depends directly on all the parameters as well as the matrix  $\mathbf{H}$ , we define the relative error:

$$\text{error}_{\sigma_e^2} = \frac{|\sigma_e^2 - \hat{\sigma}_e^2|}{\sigma_e^2}. \quad (50)$$

The results are averaged on 100 realizations of  $(\mathbf{x}, \mathbf{e})$  for each setup  $(\mathbf{H}, p, \text{iSNR})$ .

Let us denote by  $\hat{\mathbf{x}}$  an estimate of the original signal  $\mathbf{x}$ . To quantify the quality of source recovery, we define the output Signal-to-Noise Ratio (oSNR):

$$\text{oSNR} = 10 \log_{10} \left( \frac{\|\mathbf{x}\|_2^2}{\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2} \right) \quad (51)$$

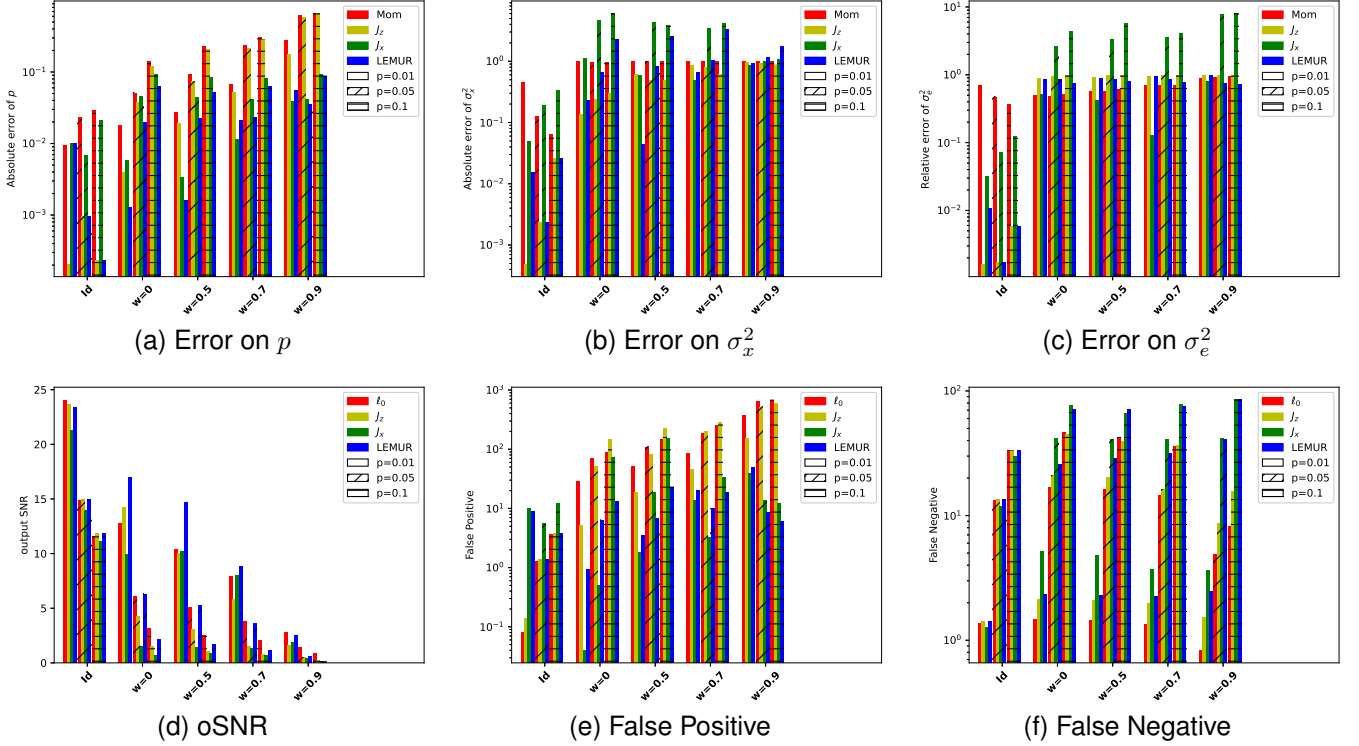


Fig. 1. iSNR = 5dB. Top: Parameter estimation (error<sub>p</sub>, error<sub>σ<sub>x</sub><sup>2</sup></sub>, error<sub>σ<sub>e</sub><sup>2</sup></sub> as defined in Eq. (49) and Eq. (50)), Bottom: output SNR as defined in Eq. (51), False Negative and False positive, Eq. (52)

We also count the number of False positives (Fp) and False negatives (Fn) in the support of  $\hat{x}$ , that is

$$\text{Fp} = \#\{\hat{x}_n \neq 0 | x_n = 0\}, \quad \text{Fn} = \#\{\hat{x}_n = 0 | x_n \neq 0\}. \quad (52)$$

The algorithms are all initialized using the method of moments described in Section III-C, with  $\mathbf{H}^T \mathbf{y}$  as input. The initial choice of the estimated sources is the zero vector. The output of Alg. 4 being a posterior mean estimate, the recovered signal  $\mathbf{z}$  is not sparse. Hence, we choose to estimate the support of the signal by considering the coordinates such that  $p(x_n \neq 0 | z_n)$  (given by Eq. (8) with  $y_n \leftarrow z_n$ ) is greater than  $\frac{1}{2}$ .

The algorithms are designated with shorthand notations. Alg. 3 will be referred to as  $J_x$  for joint estimation of  $\mathbf{x}$  and  $\theta$ . Similarly, Alg. 4 is referred to as  $J_z$ . We consider the IHT algorithm [5] dedicated to the minimization of the cost function  $\frac{1}{2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_0$  which can be thought of as a limit form of the cost function appearing in Proposition 1 when  $\sigma_x^2$  tends to infinity. The latter is a supervised algorithm, where  $\lambda$  is set empirically. IHT is run repeatedly for 100 decreasing  $\lambda$ -values defined on a logarithmic grid with warm-restart. We keep the IHT output yielding the better SNR score. The last considered algorithm is LEMUR, in which Alg. 4 is initialized using the output of Alg. 3, as described in Section IV-B.

The experiments are run with Python on a laptop with an Intel Core i7 CPU at 1.8 GHz with 16 GB of RAM<sup>1</sup>.

Figs. 1 to 3 display the performance of algorithms in terms of parameter and signal recovery for specific input signal-to-

noise ratios (iSNR) of 5, 10, and 20 dB, respectively. The top part of each figure illustrates the absolute error on parameters  $p$  and  $\sigma_x^2$ , and the relative error on  $\sigma_e^2$ . The bottom part of the figures showcases metrics related to signal recovery, that is, the output SNR, the number of false negatives and of false positives. This evaluation is repeated for different values of  $p$  and matrix operators  $\mathbf{H}$ , including the identity matrix and random matrices with varying coefficients  $w$  (see Eq. (47)). For parameter estimation, the proposed EM algorithms are compared against the methods of moments (in red). In terms of signal recovery, the three algorithms  $J_x$ ,  $J_z$  and LEMUR are compared against the outcomes of the supervised IHT algorithm (in red). The behavior remains consistent across varying iSNRs. For clarity, we will focus our analysis primarily on Fig. 2 corresponding to an iSNR of 10dB.

#### A. Parameter estimation

We first evaluate the quality of BG parameter estimation, as displayed in Figs. 2a to 2c.

Regarding the recovery of parameter  $p$ , LEMUR consistently outperforms other methods except in the denoising scenario with  $p = 0.1$ , where  $J_z$  is exceptionally good, especially for high iSNRs. It has been observed that LEMUR performs better in estimating  $\sigma_x^2$  for low values of  $p$  and large  $w$ . On the other hand, in the compressed sensing scenario with  $w = 0$ ,  $J_z$  is found to be competitive. It is worth noting that  $J_z$  outperforms other methods when  $p = 0.1$ . When focusing on  $\sigma_e^2$  estimation, the method of moments exhibits strong performance. LEMUR consistently outperforms  $J_z$  and  $J_x$  across all cases.

<sup>1</sup>Code is available at <https://github.com/PBrlt/BGBox>



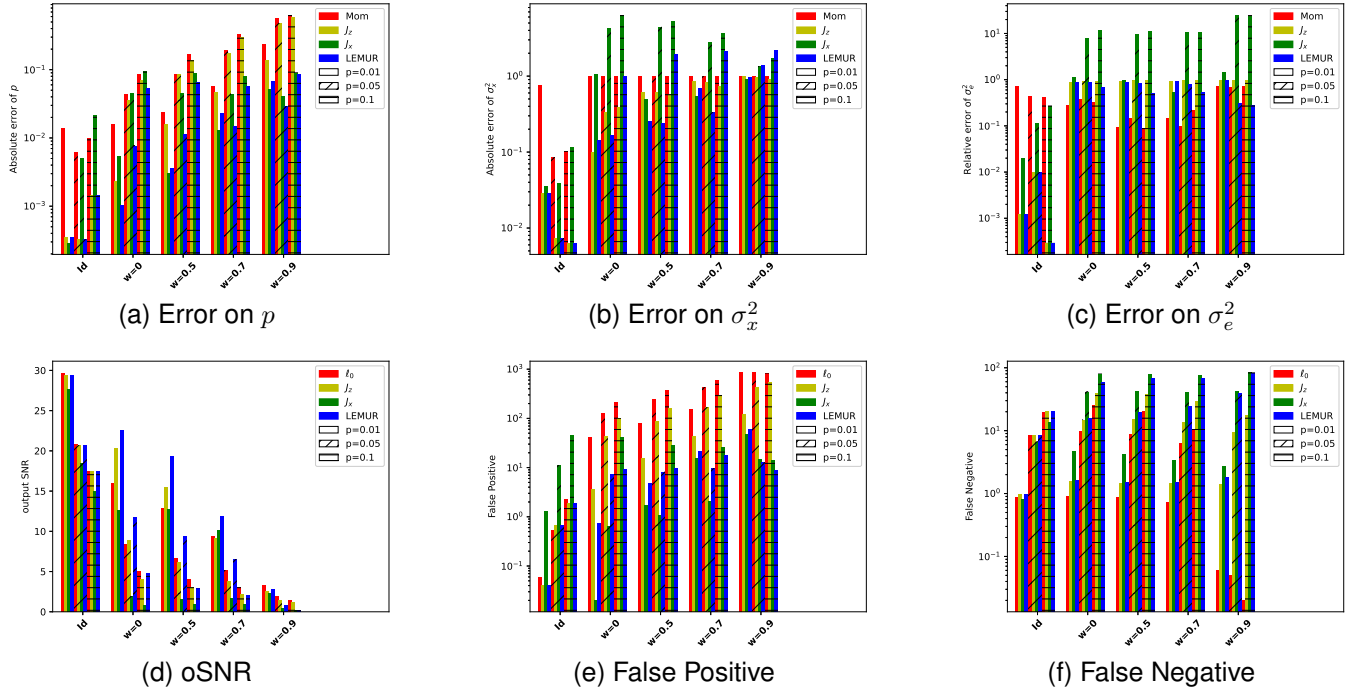


Fig. 2. iSNR = 10dB. Top: Parameter estimation ( $\text{error}_p$ ,  $\text{error}_{\sigma_x^2}$  and  $\text{error}_{\sigma_e^2}$ ), Bottom: output SNR, False Negatives and False positives.

It is essential to highlight that the methods of moments and the  $J_x$  method are prone to generating biased estimators, as discussed in Section III, which leads to lower accuracy and performance in estimating parameters, especially in denoising scenarios. Across different scenarios and iSNR levels, LEMUR is the most robust and competitive approach for parameter estimation, even with varying sparsity levels of the actual solution.

### B. Signal estimation

We now focus on evaluating the quality of signal and support recovery, see Figs. 2d to 2f. Identifying the support of a sparse signal indeed appears to be an important aspect that complements the signal estimation process.

The oSNR metric is used to evaluate the recovery of a signal. Among all approaches,  $J_x$  consistently gives the weakest performance when it comes to denoising situations. On the other hand, the LEMUR method outperforms all other methods. However, when the dictionary matrix is highly correlated ( $w = 0.9$ ), there is an exception where supervised IHT exhibits superior performance followed closely by  $J_z$ .

The estimation quality of the support is measured by the False Positive (where a coefficient is mistakenly considered present in the signal) and False Negative (where a true coefficient is missed) scores. Among various scenarios, LEMUR has been found to be the most effective method regarding False Positive rates. However, Lemur yields an increases rate False Negative score. This behavior is actually desirable, since some dictionary atoms with small weights  $x_i$  are likely to be lost in the observation noise. One thus cannot expect to detect such signal entries with high confidence.

## VI. CONCLUSION

In addressing the challenges inherent to sparse signal recovery, particularly the delicate balance required between data-fidelity and sparsity promotion through hyperparameter tuning, the LEMUR algorithm emerges as a key solution. Building upon the established Iterative Shrinkage/Thresholding algorithm used in solving Lasso/BPDN and related variants, this work introduces a novel extension to an unsupervised setting. Focusing on parameter estimation within the Bernoulli-Gaussian model, the methodology leverages maximum likelihood estimation for denoising tasks, further expanding these techniques to general inverse problems through latent variable formulations.

We proposed two distinct expectation-maximization approaches to address signal and BG model parameter estimation. While one approach displays robustness to initialization but shows bias even in simple denoising scenarios, the second approach suffers from sensitivity to initialization. LEMUR leverages on both EM strategies to yield stable estimates of both hyperparameters and signals, showcasing superior performance in output SNR and False Positive Rate recovery across various scenarios.

Noteworthy is the fact that LEMUR is operates in a fully unsupervised manner, positioning it favorably compared to the optimally (empirically) tuned IHT algorithm. Moreover, by capitalizing on the iterative shrinkage/thresholding algorithm, LEMUR ensures computational efficiency without imposing the prohibitive demands associated with MCMC methods.

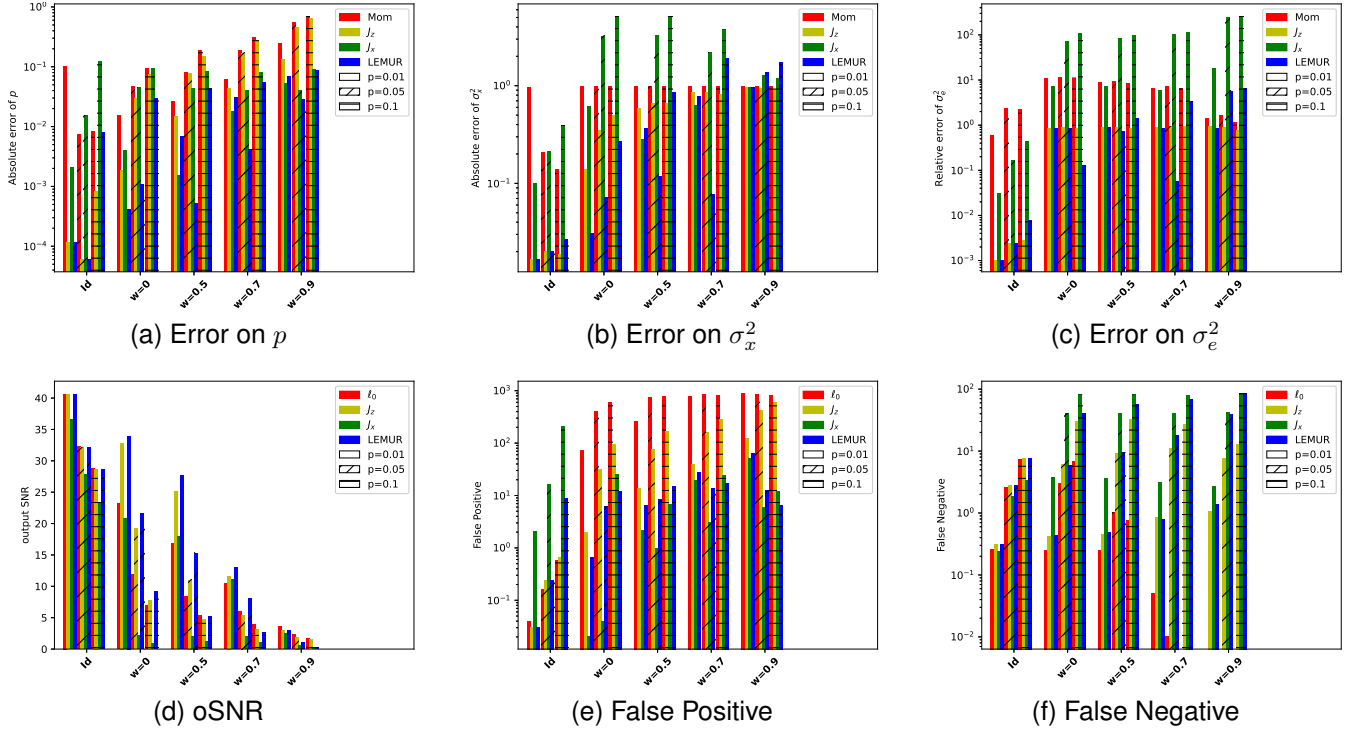


Fig. 3. iSNR = 20dB. Top: Parameter estimation ( $\text{error}_p$ ,  $\text{error}_{\sigma_x^2}$  and  $\text{error}_{\sigma_e^2}$ ), Bottom: output SNR, False Negatives and False Positives.

### APPENDIX A PROOF OF PROPOSITION 1

Using the indicator variable  $\mathbf{s}$ , and after observing that  $p(\mathbf{y}|\mathbf{x}, \mathbf{s}, \theta) = p(\mathbf{y}|\mathbf{x}, \theta)$ , we can define the joint-MAP estimator of  $(\mathbf{x}, \mathbf{s})$  as

$$\begin{aligned} (\hat{\mathbf{x}}, \hat{\mathbf{s}}) &= \underset{\mathbf{x}, \mathbf{s}}{\text{argmax}} p(\mathbf{x}, \mathbf{s}|\mathbf{y}, \theta) \\ &= \underset{\mathbf{x}, \mathbf{s}}{\text{argmax}} p(\mathbf{y}|\mathbf{x}, \theta) p(\mathbf{x}, \mathbf{s}|\theta). \end{aligned} \quad (\text{A.1})$$

With the fact that  $p(x_n \neq 0 | s_n = 1) = 1$  and  $p(x_n = 0 | s_n = 0) = 1$ , we can write

$$p(\mathbf{x}, \mathbf{s}|\theta) = \prod_{n=1}^N p(x_n, s_n|\theta) \quad (\text{A.2})$$

$$\begin{aligned} &= \prod_{\{n:s_n=0\}} p(x_n = 0 | s_n = 0, \theta) p(s_n = 0|\theta) \\ &\times \prod_{\{n:s_n=1\}} p(x_n | s_n = 1, \theta) p(s_n = 1|\theta) \end{aligned} \quad (\text{A.3})$$

$$= \prod_{\{n:s_n=0\}} (1-p) \prod_{\{n:s_n=1\}} \left( \frac{p}{\sqrt{2\pi\sigma_x^2}} e^{-\frac{x_n^2}{2\sigma_x^2}} \right) \quad (\text{A.4})$$

$$= \prod_{n=1}^N \left( \frac{1}{\sqrt{2\pi\sigma_x^2}} e^{-\frac{x_n^2}{2\sigma_x^2}} \right)^{s_n} p^{s_n} (1-p)^{1-s_n} \quad (\text{A.5})$$

Then, using the fact that for all  $n$ ,  $s_n x_n = x_n$  and  $\sum_n s_n =$

$\|\mathbf{x}\|_0$ , we can write

$$\begin{aligned} -\log p(\mathbf{x}, \mathbf{s}, \theta) &= \frac{1}{2\sigma_x^2} \|\mathbf{x}\|^2 + \log \left( \sqrt{2\pi\sigma_x^2} \frac{1-p}{p} \right)^{\|\mathbf{x}\|_0} \\ &\quad - N \log(1-p) \end{aligned} \quad (\text{A.6})$$

which does not depend on  $\mathbf{s}$  anymore. The conclusion follows from Eq. (A.1) using

$$-\log p(\mathbf{y}|\mathbf{x}, \theta) = \frac{1}{2\sigma_e^2} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2 + N \log(\sqrt{2\pi\sigma_e^2}).$$

### APPENDIX B PROOF OF THEOREM 1

The first three non-zero moments of the mixture of two centered Gaussians, see (7), are given by

$$\mathbb{E}\{y_n^2\} = p\sigma_x^2 + \sigma_e^2 \quad (\text{B.7})$$

$$\mathbb{E}\{y_n^4\} = 3p(\sigma_x^2 + \sigma_e^2)^2 + 3(1-p)(\sigma_e^2)^2 \quad (\text{B.8})$$

$$\mathbb{E}\{y_n^6\} = 15p(\sigma_x^2 + \sigma_e^2)^3 + 15(1-p)(\sigma_e^2)^3. \quad (\text{B.9})$$

Let us define the empirical moment estimators as follows:

$$m_2 = \frac{1}{N} \sum_{n=1}^N y_n^2, \quad m_4 = \frac{1}{3N} \sum_{n=1}^N y_n^4, \quad m_6 = \frac{1}{15N} \sum_{n=1}^N y_n^6$$

Denoting by  $\hat{p}, \hat{\sigma}_e^2, \hat{\sigma}_x^2$  the estimates of  $p, \sigma_e^2, \sigma_x^2$  given by the method of moments, we have:

$$m_2 = \hat{p}(\hat{\sigma}_x^2 + \hat{\sigma}_e^2) + (1-p)\sigma_e^2 \quad (\text{B.10a})$$

$$= \hat{p}\hat{\sigma}_x^2 + \hat{\sigma}_e^2 \quad (\text{B.10b})$$

$$m_4 = \hat{p}(\hat{\sigma}_x^2 + \hat{\sigma}_e^2)^2 + (1 - \hat{p})(\hat{\sigma}_e^2)^2 \quad (\text{B.11a})$$

$$= \hat{p}\hat{\sigma}_x^2(\hat{\sigma}_x^2 + 2\hat{\sigma}_e^2) + (\hat{\sigma}_e^2)^2 \quad (\text{B.11b})$$

$$m_6 = \hat{p}(\hat{\sigma}_x^2 + \hat{\sigma}_e^2)^3 + (1 - \hat{p})(\hat{\sigma}_e^2)^3 \quad (\text{B.12})$$

Hereafter, we will use the simplified notations  $p, \sigma_e^2, \sigma_x^2$  instead of  $\hat{p}, \hat{\sigma}_e^2, \hat{\sigma}_x^2$  for improved readability. Eq. (B.10b) yields:

$$p = \frac{m_2 - \sigma_e^2}{\sigma_x^2}. \quad (\text{B.13})$$

Since  $0 < p < 1$ , one can check that the following inequalities are satisfied:

$$\sigma_e^2 < m_2 \quad (\text{B.14})$$

$$m_2 - \sigma_e^2 < \sigma_x^2 \quad (\text{B.15})$$

$$m_4^2 < m_2 m_6. \quad (\text{B.16})$$

The latter is obtained from Eq. (B.10a), Eq. (B.11a) and Eq. (B.12):

$$m_4^2 - m_2 m_6 = p(p - 1)\sigma_e^2(\sigma_x^2 + \sigma_e^2)(\sigma_x^2)^2 < 0. \quad (\text{B.17})$$

Plugging Eq. (B.13) in Eq. (B.11b), one has

$$m_4 = (m_2 - \sigma_e^2)(\sigma_x^2 + 2\sigma_e^2) + (\sigma_e^2)^2 \quad (\text{B.18})$$

hence

$$\sigma_x^2 = \frac{m_4 - (\sigma_e^2)^2}{m_2 - \sigma_e^2} - 2\sigma_e^2. \quad (\text{B.19})$$

Also, combining Eq. (B.10b) with Eq. (B.11b), we get  $m_2^2 - m_4 = \sigma_x^2(p^2 - p) < 0$ , so

$$m_2^2 < m_4. \quad (\text{B.20})$$

Using Eqs. (B.13) and (B.19) within Eq. (B.12), and after some tedious algebra, we get

$$(m_4 - m_2^2)(\sigma_e^2)^2 + (m_2 m_4 - m_6)\sigma_e^2 + (m_2 m_6 - m_4^2) = 0 \quad (\text{B.21})$$

Hence, we obtain the following quadratic equation in  $\sigma_e^2$ :

$$\varphi(\sigma_e^2) := (\sigma_e^2)^2 + \frac{m_2 m_4 - m_6}{m_4 - m_2^2} \sigma_e^2 + \frac{m_2 m_6 - m_4^2}{m_4 - m_2^2} = 0. \quad (\text{B.22})$$

According to (B.14), we need to have one root of  $\varphi$  in between 0 and  $m_2$ . On the one hand,

$$\begin{aligned} \varphi(\sigma_e^2 = m_2) &= m_2^2 + \frac{m_2^2 m_4 - m_2 m_6}{m_4 - m_2^2} + \frac{m_2 m_6 - m_4^2}{m_4 - m_2^2} \\ &= m_2^2 - m_4 < 0 \end{aligned} \quad (\text{B.23})$$

using Eq. (B.20). On the other hand,

$$\varphi(0) = \frac{m_2 m_6 - m_4^2}{m_4 - m_2^2} > 0 \quad (\text{B.24})$$

using Eq. (B.16). It follows that  $\varphi(\sigma_e^2) = 0$  admits one root in  $(0, m_2)$ .

Since  $m_4 - m_2^2 > 0$ , Descartes's rule of signs ensures that the weight of  $\sigma_e^2$  within Eq. (B.22) is negative. Using

Eq. (B.16), we get  $m_2 m_4 - m_6 < 0$ . Then, the root in  $(0, m_2)$  is given by

$$\sigma_e^2 = A - \sqrt{A^2 - B} \quad (\text{B.25})$$

with

$$A = \frac{m_6 - m_2 m_4}{2(m_4 - m_2^2)} \text{ and } B = \frac{m_2 m_6 - m_4^2}{m_4 - m_2^2}. \quad (\text{B.26})$$

## REFERENCES

- [1] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1, pp. 267–288, 1996.
- [2] S. S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM review*, vol. 43, no. 1, pp. 129–159, 2001.
- [3] P. L. Combettes and V. R. Wajs, "Signal recovery by proximal forward-backward splitting," *SIAM Journal on Multiscale Modeling and Simulation*, vol. 4, no. 4, pp. 1168–1200, 2005.
- [4] A. Beck and M. Teboulle, "A fast iterative shrinkage-thresholding algorithm for linear inverse problems," *SIAM Journal on Imaging Science*, vol. 2, no. 1, pp. 183–202, 2009.
- [5] T. Blumensath and M. E. Davies, "Iterative hard thresholding for compressed sensing," *Appl. Comp. Harmonic Anal.*, vol. 27, no. 3, pp. 265–274, 2009.
- [6] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, "Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition," in *Proc. 27th Asilomar Conf. on Signals, Systems and Computers*, vol. 1, Nov. 1993, pp. 40–44.
- [7] C. Soussen, J. Idier, D. Brie, and J. Duan, "From Bernoulli-Gaussian deconvolution to sparse signal restoration," *IEEE Transactions on Signal Processing*, vol. 59, no. 10, pp. 4572–4584, 2011.
- [8] S. Bourguignon, J. Ninin, H. Carfantan, and M. Mongeau, "Exact sparse approximation problems via mixed-integer programming: Formulations and computational performance," *IEEE Transactions on Signal Processing*, vol. 64, no. 6, pp. 1405–1419, 2015.
- [9] H. Hazimeh, R. Mazumder, and A. Saab, "Sparse regression at scale: Branch-and-bound rooted in first-order optimization," *Mathematical Programming*, vol. 196, no. 1-2, pp. 347–388, 2022.
- [10] D. L. Donoho and I. M. Johnstone, "Adapting to unknown smoothness via wavelet shrinkage," *Journal of the American Statistical Association*, vol. 90, no. 432, pp. 1200–1224, 1995.
- [11] T. Blu and F. Luisier, "The SURE-LET approach to image denoising," *IEEE Trans. Image Process.*, vol. 16, no. 11, pp. 2778–2786, 2007.
- [12] C.-A. Deledalle, S. Vaiteer, J. Fadili, and G. Peyré, "Stein unbiased gradient estimator of the risk (sugar) for multiple parameter selection," *SIAM Journal on Imaging Sciences*, vol. 7, no. 4, pp. 2448–2487, 2014.
- [13] R. Giryes, M. Elad, and Y. C. Eldar, "The projected GSURE for automatic parameter tuning in iterative shrinkage methods," *Applied and Computational Harmonic Analysis*, vol. 30, no. 3, pp. 407–422, 2011.
- [14] C.-A. Deledalle, G. Peyré, and J. Fadili, "Stein consistent risk estimator (score) for hard thresholding," *arXiv preprint arXiv:1301.5874*, 2013.
- [15] M. Pereyra, J. M. Bioucas-Dias, and M. A. Figueiredo, "Maximum-a-posteriori estimation with unknown regularisation parameters," in *23rd European Signal Processing Conference (EUSIPCO)*, 2015, pp. 230–234.
- [16] D. Calvetti, E. Somersalo, and A. Strang, "Hierarchical Bayesian models and sparsity:  $\ell_2$ -magic," *Inverse Problems*, vol. 35, no. 3, p. 035003, 2019.
- [17] F. Champagnat, Y. Goussard, and J. Idier, "Unsupervised deconvolution of sparse spike trains using stochastic approximation," *IEEE Transactions on Signal Processing*, vol. 44, no. 12, pp. 2988–2998, Dec. 1996.
- [18] M. A. Figueiredo and R. D. Nowak, "An EM algorithm for wavelet-based image restoration," *IEEE Transactions on Image Processing*, vol. 12, no. 8, pp. 906–916, 2003.
- [19] M. Kowalski and T. Rodet, "An unsupervised algorithm for hybrid/morphological signal decomposition," in *Proc. IEEE ICASSP*, 2011, pp. 4112–4115.
- [20] E. Gassiat, F. Monfront, and Y. Goussard, "On simultaneous signal estimation and parameter identification using a generalized likelihood approach," *IEEE transactions on information theory*, vol. 38, no. 1, pp. 157–162, 1992.
- [21] D. Ge, J. Idier, and E. Le Carpentier, "A new MCMC algorithm for blind Bernoulli-Gaussian deconvolution," in *Proc. 16th Eur. Sig. Proc. Conf.*, 2008, pp. 1–5.

- [22] M. Amrouche, H. Carfantan, and J. Idier, "Efficient sampling of bernoulli-gaussian-mixtures for sparse signal restoration." *IEEE Transactions on Signal Processing*, vol. 70, pp. 5578–5591, 2022.
- [23] L. Chaari, J.-Y. Tournier, and C. Chaux, "Sparse signal recovery using a Bernoulli generalized Gaussian prior," in *23rd European Signal Processing Conference (EUSIPCO)*. IEEE, 2015, pp. 1711–1715.
- [24] P. Barbault, M. Kowalski, and C. Soussen, "Parameter estimation in sparse inverse problems using Bernoulli-Gaussian prior," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 5413–5417.
- [25] A. P. Dempster, N. M. Laird, and D. B. Rubin, "Maximum likelihood from incomplete data via the EM algorithm," *Journal of the Royal Statistical Society: Series B*, vol. 39, no. 1, pp. 1–22, 1977.
- [26] M. Protter, I. Yavneh, and M. Elad, "Closed-form MMSE estimation for signal denoising under sparse representation modeling over a unitary dictionary," *IEEE Transaction on Signal Processing*, vol. 58, no. 7, pp. 3471–3484, 2010.
- [27] G. Celeux and G. Govaert, "A classification EM algorithm for clustering and two stochastic versions," *Computational statistics & Data analysis*, vol. 14, no. 3, pp. 315–332, 1992.
- [28] S. M. Kay, *Fundamentals of statistical signal processing: estimation theory*. Prentice-Hall, Inc., 1993.
- [29] Y. Goussard, G. Demoment, and J. Idier, "A new algorithm for iterative deconvolution of sparse spike trains," in *Proc. IEEE ICASSP*, 1990, pp. 1547–1550.