

# SPARSITY VIA HYPERPRIORS: A THEORETICAL AND ALGORITHMIC STUDY UNDER EMPIRICAL BAYES FRAMEWORK\*

ZHITAO LI <sup>†</sup>, YIQU DONG <sup>‡</sup>, AND XUEYING ZENG <sup>§</sup>

**Abstract.** This paper presents a comprehensive analysis of hyperparameter estimation within the empirical Bayes framework (EBF) for sparse learning. By studying the influence of hyperpriors on the solution of EBF, we establish a theoretical connection between the choice of the hyperprior and the sparsity as well as the local optimality of the resulting solutions. We show that some strictly increasing hyperpriors, such as half-Laplace and half-generalized Gaussian with the power in  $(0, 1)$ , effectively promote sparsity and improve solution stability with respect to measurement noise. Based on this analysis, we adopt a proximal alternating linearized minimization (PALM) algorithm with convergence guarantees for both convex and concave hyperpriors. Extensive numerical tests on two-dimensional image deblurring problems demonstrate that introducing appropriate hyperpriors significantly promotes the sparsity of the solution and enhances restoration accuracy. Furthermore, we illustrate the influence of the noise level and the ill-posedness of inverse problems to EBF solutions.

**Key words.** empirical Bayes framework, sparse Bayesian learning, hyperparameter estimation, generalized Gamma distribution, proximal alternating linearized minimization

**AMS subject classifications.** 62F15, 65K10, 65F22

**1. Introduction.** Inferring signals or parameters from limited and noisy measurements is a central challenge in inverse problems [15], which commonly appear in many applications, e.g. medical imaging, geophysics, and engineering. The unknown often admits a sparse representation under an appropriate transform, motivating the search for its sparse representation over a possibly overcomplete dictionary or basis. The prior information on the sparsity can be incorporated through regularization technique. Some commonly used regularizers include the  $l_0$ -norm [23], the  $l_1$ -norm that is a convex relaxation of  $l_0$  [8, 27], and the  $l_p$  quasi-norm with  $0 < p < 1$  [13, 30, 36, 18, 11], etc. In these methods, the balance between the fitting of the data and regularization is controlled by a regularization parameter, and the selection of this parameter has a significant impact on the quality of the reconstruction results. Choosing an appropriate regularization parameter is always challenging.

Most of the above mentioned approaches can be interpreted within a Bayesian framework as maximum a posteriori (MAP) estimates with a sparsity-inducing prior. Hierarchical Bayesian models extend this idea by introducing additional layers of priors over hyperparameters, allowing the model to adaptively learn regularization strength and sparsity patterns from data [6]. In sparse dictionary learning, Gaussian-inverse Gamma hierarchical priors have been used to jointly infer sparse codes and noise variance [37]. Generalized Gamma hyperpriors further enhance the sparsity and

---

\*Submitted to the editors November 6, 2025.

**Funding:** ZL and XZ were supported by the National Key R&D Plan of China (No. 2024YFC2814403), and the Fundamental Research Funds for the Central Universities (No. 202264006). YD was supported by a Villum Investigator grant (No. 25893) from The Villum Foundation.

<sup>†</sup>School of Mathematical Sciences, Ocean University of China, Qingdao, China (lizhitao3378@stu.ouc.edu.cn)

<sup>‡</sup>Department of Applied Mathematics and Computer Science, Technical University of Denmark, Kgs. Lyngby, Denmark (yido@dtu.dk)

<sup>§</sup>Corresponding author. Laboratory of Marine Mathematics, Ocean University of China, Qingdao, China (zxying@ouc.edu.cn)

stability of solutions with respect to measurement noise, providing stronger control over the reconstructed solution [7, 5, 4]. Recently, horseshoe shrinkage prior, which encodes a half-Cauchy hyperprior in the conditional Gaussian prior, has been shown to be effective in promoting sparsity [29, 25]. In those methods, the unknown and the hyperparameters are usually solved in an alternate iteration scheme.

In contrast, sparse Bayesian learning (SBL) offers an alternative approach that estimates hyperparameters directly from data without requiring explicit hierarchical hyperpriors. SBL originated from the automatic relevance determination (ARD) framework [19, 28], which is a machine learning technique for regression and classification. Under SBL, the hyperparameters, which parameterize the variance of the Gaussian prior, are estimated through evidence maximization, also known as Type-II maximum likelihood [2]. By marginalizing the unknown, the hyperparameters and the data are linked directly. Furthermore, the marginalization process provides a regularization mechanism, which can shrink many hyperparameters to zero and further induce sparsity in the unknown. In [28] Tipping established the theoretical foundation of SBL, while subsequent analysis explored the geometry of the marginal likelihood and revealed both its sparsity-promoting nature and its nonconvexity [10]. In [33, 31] Wipf et al. analyzed the structure of SBL and showed that it is equivalent to an optimization problem with a nonconvex sparsity-inducing penalty, clarifying the underlying mechanism of sparsity formation. Later, they extended this framework to latent-variable Bayesian models, which further enhance the flexibility of sparse priors through hierarchical latent representations [35]. In [38] Yu et al. proposed a unified framework for hyperparameter estimation in SBL, which offers a new perspective to interpret, analyze and compare the algorithms used for hyperparameter estimation, both theoretically and numerically. Building on these insights, several extensions generalizing SBL have been developed to handle structured sparsity and large-scale imaging problems [1, 12].

Despite its success across inverse and regression tasks, SBL remains challenged by the nonconvexity of the corresponding optimization problem, the sensitivity to measurement noise and the computational burden of large matrix operations. It motivates continued research on more stable and interpretable empirical Bayes framework (EBF) [20]. In EBF instead of maximizing the marginal likelihood, it maximizes the marginal posterior by incorporating the hyperpriors. Recent advances extend EBF to several important directions. The work in [28, 14] discussed the properties of the EBF solution in the case of an inverse Gamma hyperprior. A hyperprior combining a concave function and a convex function was proposed in [26]. It was also pointed out that, except for certain cases, the generalized Gamma distribution does not impose sparsity, which is the reason why most existing EBF only rely on noninformative or weakly informative hyperpriors. In [34] Wipf et al. extended EBF to the simultaneous sparse approximation problem, providing a method for jointly reconstructing multiple signals with shared sparsity patterns. In addition, EBF has been applied in many applications or as a building block in other approaches, e.g. a unified Bayesian formulation covering EBF tailored for biomedical inverse problems such as MEG/EEG source imaging [32], adaptive schemes for deep models that integrate layer-wise sparsity [9], meta-prior learning that uses EBF to learn higher-level prior distributions [22], etc. Despite empirical successes, EBF inherits challenges of nonconvex marginal-likelihood optimization and sensitivity to hyperprior families, motivating research into robust hyperprior design and sparsity-promoting modeling strategies.

Among the work on EBF, it still lacks a solid study on how hyperpriors promote sparsity meanwhile stabilizing the problem. In this work, we provide a comprehensive

analysis of the hyperparameter estimation under EBF. By analyzing the properties of stationary points in this framework, we establish a theoretical connection between the structure of the hyperprior and both the sparsity and local optimality of the resulting solutions. We show that hyperpriors under certain conditions can effectively promote sparsity and enhance solution stability with respect to measurement noise. Based on these theoretical insights, we focus on the generalized Gamma family as hyperpriors and study their influence on sparsity under different settings. In particular, the hyperpriors following such as the half-Laplace distribution and the half-generalized Gaussian distribution with the power in  $(0, 1)$  can effectively promote the sparsity of hyperparameters and improve the stability of the solution to EBF. In addition, we introduce the proximal alternating linearized minimization (PALM) algorithm to solve our problem and provide a rigorous convergence analysis. The convergence result shows that regardless of whether the hyperprior is convex or concave, the convergence of PALM to a stationary point is guaranteed. Finally, we use two-dimensional image deblurring problems to illustrate the effects of different hyperpriors on the restoration quality and the influence of the ill-posedness and the noise level on the EBF results. The results confirm our theoretical findings and show that choosing an appropriate hyperprior, such as half-Laplace or half-generalized Gaussian with power in  $(0, 1)$ , can significantly promote sparsity and improve restoration accuracy.

The remainder of this paper is structured as follows. Section 2 reviews EBF briefly and formulates the hyperparameter estimation problem. Section 3 establishes the theoretical underpinnings of sparsity promotion and solution optimality in EBF. Section 4 presents the proposed numerical algorithm together with its convergence analysis. In Section 5 we test our method on the 2D image deblurring problems. Finally, we conclude our work with some perspectives for future research in Section 6.

**2. Empirical Bayes framework.** We consider the following discrete inverse problem:

$$(2.1) \quad \mathbf{y} = \mathbf{F}\mathbf{x} + \boldsymbol{\epsilon},$$

where  $\mathbf{y} \in \mathbb{R}^m$  denotes the vectorized observed data,  $\mathbf{x} \in \mathbb{R}^n$  is the unknown that we are interested in recovering,  $\mathbf{F} \in \mathbb{R}^{m \times n}$  is the forward operator modeling the physical relation between the parameters and the data, and  $\boldsymbol{\epsilon} \in \mathbb{R}^m$  is the vector of noises. We assume  $m \leq n$  and the noise vector  $\boldsymbol{\epsilon}$  follows an independent and identically distributed (i.i.d.) normal distribution with zero mean and inverse variance  $\sigma$ , i.e.,

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma^{-1} \mathbf{I}_m),$$

where  $\mathbf{I}_m$  is an identity matrix of size  $m$ , and  $\mathcal{N}$  denotes the multivariate Gaussian distribution. Then the likelihood is in the following form:

$$(2.2) \quad \pi(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{F}\mathbf{x}, \sigma^{-1} \mathbf{I}_m).$$

Furthermore, we assume that  $\mathbf{x}$  follows a Gaussian distribution due to its computational simplicity. However, Gaussian priors are not well suited to promote sparsity. One possible way to retain computational convenience while promoting sparsity is via a conditional Gaussian prior in the form of

$$(2.3) \quad \pi(\mathbf{x}|\boldsymbol{\gamma}) = \mathcal{N}(0, \boldsymbol{\Gamma}) \propto \prod_i \frac{1}{\sqrt{\gamma_i}} \exp\left(-\frac{x_i^2}{2\gamma_i}\right).$$

where  $\boldsymbol{\gamma} \in \mathbb{R}_+^n := \{\boldsymbol{\gamma} \in \mathbb{R}^n : \gamma_i \geq 0, 1 \leq i \leq n\}$  collects all the unknown variance  $\gamma_i$ , which is expected to be small when  $x_i$  vanishes, while it should be larger in correspondence of the few nonzero-entries of  $\mathbf{x}$ . In the extreme scenario where  $\gamma_i = 0$ ,  $x_i$  essentially reduces to a degenerate distribution, represented as a point mass at the mean, i.e. 0. In addition,  $\boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma})$  is the diagonal matrix with diagonal entries given by the hyperparameters  $\boldsymbol{\gamma}$ .

In the Bayesian paradigm, the hyperparameters  $\boldsymbol{\gamma}$  are modeled as random variables, and their estimation becomes part of the problem. A common approach to hyperparameter estimations is to use sparse Bayesian learning [10], also known as Type II maximum likelihood, which computes the marginal likelihood by integrating out the unknown and then maximizes this marginal likelihood with respect to the hyperparameters. Since both the likelihood  $\pi(\mathbf{y}|\mathbf{x})$  and the prior  $\pi(\mathbf{x}|\boldsymbol{\gamma})$  are Gaussian, it follows from the conjugate property of the Gaussian distribution that the marginal likelihood is given by

$$(2.4) \quad \pi(\mathbf{y}|\boldsymbol{\gamma}) = \int \pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}|\boldsymbol{\gamma})d\mathbf{x} = \mathcal{N}(\mathbf{y}|\mathbf{0}, \mathbf{S}(\boldsymbol{\gamma}))$$

with  $\mathbf{S}(\boldsymbol{\gamma}) = \sigma^{-1}\mathbf{I}_m + \mathbf{F}\mathbf{F}^\top$ , i.e.,

$$(2.5) \quad \pi(\mathbf{y}|\boldsymbol{\gamma}) \propto |\mathbf{S}(\boldsymbol{\gamma})|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{y}^\top(\mathbf{S}(\boldsymbol{\gamma}))^{-1}\mathbf{y}\right).$$

In this paper, we consider the case of limited and noisy data, i.e.,  $m \leq n$  and  $\epsilon \neq \mathbf{0}$ . Due to the existence of the noise and the under-determined problem, instead of marginal likelihood, we suggest using *empirical Bayes framework* [32] to maximize the marginal posterior given as

$$\pi(\boldsymbol{\gamma}|\mathbf{y}) \propto \pi(\mathbf{y}|\boldsymbol{\gamma})\pi(\boldsymbol{\gamma})$$

according to Bayes' theorem. Here,  $\pi(\boldsymbol{\gamma})$  represents the hyperprior and incorporates the prior information of  $\boldsymbol{\gamma}$ . In order to promote sparsity and keep the model simple, we expect that the elements in  $\boldsymbol{\gamma}$  are i.i.d. and that all but few should be close to zero. Thus, the hyperprior is in a product form:

$$(2.6) \quad \pi(\boldsymbol{\gamma}) \propto \prod_{i=1}^n \exp(-\mathcal{H}(\gamma_i)),$$

where  $\mathcal{H}(\cdot) : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  is a continuous function. Based on (2.5) and (2.6), the marginal posterior of  $\boldsymbol{\gamma}$  can be expressed as

$$(2.7) \quad \pi(\boldsymbol{\gamma}|\mathbf{y}) \propto |\mathbf{S}(\boldsymbol{\gamma})|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}\mathbf{y}^\top(\mathbf{S}(\boldsymbol{\gamma}))^{-1}\mathbf{y} - \sum_{i=1}^n \mathcal{H}(\gamma_i)\right).$$

A commonly used point estimate for (2.7) is the MAP estimate, where the mode of the posterior is set as the single point representative of the whole density function:

$$(2.8) \quad \boldsymbol{\gamma}^* = \arg \max_{\boldsymbol{\gamma} \in \mathbb{R}_+^n} \pi(\boldsymbol{\gamma}|\mathbf{y}) = \arg \min_{\boldsymbol{\gamma} \in \mathbb{R}_+^n} \mathcal{J}(\boldsymbol{\gamma}),$$

where the cost functional  $\mathcal{J}$  is defined as

$$(2.9) \quad \mathcal{J}(\boldsymbol{\gamma}) = \frac{1}{2}\mathbf{y}^\top(\mathbf{S}(\boldsymbol{\gamma}))^{-1}\mathbf{y} + \frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}) + \sum_{i=1}^n \mathcal{H}(\gamma_i).$$

Once we obtain the estimate  $\gamma^*$  of the hyperparameters, we can derive the conditional posterior distribution  $\pi(\mathbf{x}|\mathbf{y}, \gamma^*)$  of the unknown  $\mathbf{x}$  through

$$(2.10) \quad \pi(\mathbf{x}|\mathbf{y}, \gamma^*) \propto \pi(\mathbf{y}|\mathbf{x})\pi(\mathbf{x}|\gamma^*) = \mathcal{N}(\mathbf{x}|\mu(\gamma^*), \Sigma(\gamma^*)),$$

whose mean and covariance are

$$\mu(\gamma^*) = \mathbf{F}^* \mathbf{F}^T (\mathbf{S}(\gamma^*))^{-1} \mathbf{y} \quad \text{and} \quad \Sigma(\gamma^*) = \mathbf{F}^* - \mathbf{F}^* \mathbf{F}^T (\mathbf{S}(\gamma^*))^{-1} \mathbf{F} \mathbf{F}^*.$$

**3. Sparsity and local optimality.** In this section, we study the influence of the hyperprior on the sparsity and local optimality of the Karush–Kuhn–Tucker (KKT) points of the proposed optimization problem (2.8).

The Lagrangian function of (2.8) is given by

$$(3.1) \quad \mathcal{L}(\gamma, \mu) = \mathcal{J}(\gamma) - \mu^\top \gamma,$$

with the KKT conditions:

$$(3.2) \quad \begin{cases} \nabla \mathcal{J}(\gamma) - \mu = \mathbf{0}, \\ \mu \geq \mathbf{0}, \gamma \geq \mathbf{0}, \\ \mu^\top \gamma = 0, \end{cases}$$

where  $\mu \in \mathbb{R}_+^n$  is the vector of Lagrange multipliers.

**3.1. Derivatives of  $\mathcal{J}(\gamma)$ .** Let  $\mathbf{F} = [\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_n]$ . Then

$$(3.3) \quad \frac{\partial \mathbf{S}(\gamma)}{\partial \gamma_i} = \frac{\partial}{\partial \gamma_i} \left( \sigma^{-1} \mathbf{I}_m + \sum_{i=1}^n \gamma_i \mathbf{f}_i \mathbf{f}_i^\top \right) = \mathbf{f}_i \mathbf{f}_i^\top.$$

The derivatives of the inverse and log determinant of a matrix  $\mathbf{A}$  with respect to  $t$  are given in [21]

$$(3.4) \quad \frac{\partial \mathbf{A}^{-1}}{\partial t} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial t} \mathbf{A}^{-1}, \quad \frac{\partial \log \det \mathbf{A}}{\partial t} = \text{Tr} \left( \mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial t} \right),$$

where  $\text{Tr}(\mathbf{B})$  denotes the trace of the matrix  $\mathbf{B}$ . Then, we can obtain

$$(3.5) \quad \frac{\partial \mathbf{y}^\top \mathbf{S}(\gamma)^{-1} \mathbf{y}}{\partial \gamma_i} = -\mathbf{y}^\top \mathbf{S}(\gamma)^{-1} \mathbf{f}_i \mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{y},$$

and

$$\frac{\partial \log \det \mathbf{S}(\gamma)}{\partial \gamma_i} = \text{Tr} (\mathbf{S}(\gamma)^{-1} \mathbf{f}_i \mathbf{f}_i^\top) = \mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{f}_i.$$

Let  $\tilde{p}_i = \mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{y}$  and  $\tilde{q}_{ij} = \mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{f}_j$ . The partial derivative of  $\mathcal{J}(\gamma)$  with respect to  $\gamma_i$  is

$$(3.6) \quad \frac{\partial \mathcal{J}(\gamma)}{\partial \gamma_i} = \frac{1}{2} \tilde{q}_{ii} - \frac{1}{2} \tilde{p}_i^2 + \mathcal{H}'(\gamma_i).$$

Following the similar procedure in (3.5), we have

$$\frac{\partial \tilde{q}_{ii}}{\partial \gamma_j} = -(\mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{f}_j)^2 = -\tilde{q}_{ij}^2,$$

and

$$\frac{\partial \tilde{p}_i}{\partial \gamma_j} = -(\mathbf{f}_i^\top \mathbf{S}(\boldsymbol{\gamma})^{-1} \mathbf{f}_j) (\mathbf{f}_j^\top \mathbf{S}(\boldsymbol{\gamma})^{-1} \mathbf{y}) = -\tilde{q}_{ij} \tilde{p}_j.$$

Hence,

$$(3.7) \quad \frac{\partial^2 \mathcal{J}(\boldsymbol{\gamma})}{\partial \gamma_i \partial \gamma_j} = -\frac{1}{2} \tilde{q}_{ij}^2 + \tilde{p}_i \tilde{q}_{ij} \tilde{p}_j + \delta_{ij} \mathcal{H}''(\gamma_i),$$

where  $\delta_{ij}$  is the Kronecker delta function allowing us to separate out the additional (diagonal) term that appears only when  $i = j$ .

**3.2. Sparsity and local optimality of KKT points.** We now study the sparsity of the KKT points of (2.8). Since the KKT conditions (3.2) are separable componentwise, we focus on the  $i$ -th component. To this aim, we first rewrite the covariance matrix

$$(3.8) \quad \mathbf{S}(\boldsymbol{\gamma}) = \mathbf{S}_{-i} + \gamma_i \mathbf{f}_i \mathbf{f}_i^\top$$

to isolate its dependence on a single hyperparameter  $\gamma_i$ , where  $\mathbf{S}_{-i} = \sigma^{-1} \mathbf{I}_m + \sum_{j \neq i} \gamma_j \mathbf{f}_j \mathbf{f}_j^\top$  is the covariance matrix excluding the contribution of  $\mathbf{f}_i$ .

Using the Sherman-Morrison formula [17], we obtain

$$(3.9) \quad \mathbf{S}(\boldsymbol{\gamma})^{-1} = \mathbf{S}_{-i}^{-1} - \frac{\gamma_i \mathbf{S}_{-i}^{-1} \mathbf{f}_i \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1}}{1 + \gamma_i \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{f}_i},$$

and

$$(3.10) \quad \det \mathbf{S}(\boldsymbol{\gamma}) = (1 + \gamma_i \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{f}_i) \det \mathbf{S}_{-i}.$$

Substituting (3.9) and (3.10) into the cost function (2.9) yields

$$(3.11) \quad \mathcal{J}(\boldsymbol{\gamma}) = \mathcal{J}_{-i}(\boldsymbol{\gamma}_{-i}) - \frac{1}{2} \frac{\gamma_i (\mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{y})^2}{1 + \gamma_i \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{f}_i} + \frac{1}{2} \log(1 + \gamma_i \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{f}_i) + \mathcal{H}(\gamma_i),$$

where

$$\mathcal{J}_{-i}(\boldsymbol{\gamma}_{-i}) = \frac{1}{2} \mathbf{y}^\top \mathbf{S}_{-i}^{-1} \mathbf{y} + \frac{1}{2} \log \det \mathbf{S}_{-i} + \sum_{j \neq i} \mathcal{H}(\gamma_j),$$

and  $\boldsymbol{\gamma}_{-i} = [\gamma_1, \dots, \gamma_{i-1}, \gamma_{i+1}, \dots, \gamma_n]^\top$  with  $\gamma_i$  removed from  $\boldsymbol{\gamma}$ . We have now isolated the terms with respect to  $\gamma_i$ .

Let  $p_i = \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{y}$  and  $q_i = \mathbf{f}_i^\top \mathbf{S}_{-i}^{-1} \mathbf{f}_i$ . Note that  $q_i > 0$ . Then we can define the scalar function

$$(3.12) \quad \mathcal{L}(\gamma_i) = -\frac{1}{2} \frac{p_i^2 \gamma_i}{1 + q_i \gamma_i} + \frac{1}{2} \log(1 + q_i \gamma_i) + \mathcal{H}(\gamma_i),$$

and we have  $\mathcal{J}(\boldsymbol{\gamma}) = \mathcal{J}_{-i}(\boldsymbol{\gamma}_{-i}) + \mathcal{L}(\gamma_i)$ . The derivative of  $\mathcal{J}(\boldsymbol{\gamma})$  with respect to  $\gamma_i$  is

$$(3.13) \quad \frac{\partial \mathcal{J}(\boldsymbol{\gamma})}{\partial \gamma_i} = \mathcal{L}'(\gamma_i) = \frac{1}{2} \frac{q_i^2 \gamma_i + (q_i - p_i^2)}{(1 + q_i \gamma_i)^2} + \mathcal{H}'(\gamma_i).$$

Note that due to (3.9), we have

$$\mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{f}_i = q_i - \frac{\gamma_i q_i^2}{1 + \gamma_i q_i} = \frac{q_i}{1 + q_i \gamma_i},$$

and

$$\mathbf{f}_i^\top \mathbf{S}(\gamma)^{-1} \mathbf{y} = p_i - \frac{\gamma_i p_i q_i}{1 + \gamma_i q_i} = \frac{p_i}{1 + q_i \gamma_i}.$$

Therefore, the two different expressions of  $\frac{\partial \mathcal{J}}{\partial \gamma_i}$  in (3.6) and (3.13) are indeed consistent.

Let us drop the subscript for clarity, and consider the following univariate optimization problem:

$$(3.14) \quad \min_{\gamma \geq 0} \mathcal{L}(\gamma) = -\frac{1}{2} \frac{p^2 \gamma}{1 + q\gamma} + \frac{1}{2} \log(1 + q\gamma) + \mathcal{H}(\gamma).$$

Since  $\mathcal{L}(\gamma)$  is proper and  $\mathcal{L}(\gamma) \rightarrow +\infty$  as  $\gamma \rightarrow +\infty$ , this problem admits at least one minimizer, satisfying the KKT conditions:

$$(3.15) \quad \begin{cases} \frac{1}{2} \frac{q^2 \gamma + (q-p^2)}{(1+q\gamma)^2} + \mathcal{H}'(\gamma) - \mu = 0, \\ \mu \geq 0, \gamma \geq 0, \\ \mu \gamma = 0. \end{cases}$$

**THEOREM 3.1.** *Assume  $\mathcal{H}$  is nonnegative in  $[0, +\infty)$ . Then the following statements hold:*

- (I) *If  $\gamma^* = 0$  satisfies the KKT conditions (3.15), then it must hold that  $q - p^2 \geq -2\mathcal{H}'(0^+)$  with  $\mathcal{H}'(0^+) = \lim_{\gamma \rightarrow 0^+} \mathcal{H}'(\gamma)$ .*
- (II) *If  $q - p^2 > -2\mathcal{H}'(0^+)$ , then  $\gamma^* = 0$  satisfies the KKT conditions (3.15).*
- (III) *Any  $\gamma^* > 0$  satisfies the KKT conditions (3.15) if and only if it is a positive root of*

$$(3.16) \quad \frac{1}{2} \frac{q^2 \gamma + (q - p^2)}{(1 + q\gamma)^2} + \mathcal{H}'(\gamma) = 0.$$

*Proof.* Let  $\mu^* = \mathcal{H}'(0^+) - \frac{p^2 - q}{2}$ . For (I), if  $\gamma^* = 0$  satisfies the KKT conditions, then  $\mu^*$  satisfies the stationarity of (3.15). By the dual feasibility, we must have  $\mu^* \geq 0$ , which implies  $q - p^2 \geq -2\mathcal{H}'(0^+)$ . For (II), if  $q - p^2 > -2\mathcal{H}'(0^+)$ , then  $\mu^* > 0$  and hence  $(\gamma^*, \mu^*) = (0, \mu^*)$  satisfies the KKT conditions.

For (III), if  $\gamma^* > 0$  is a solution to (3.16), then it is straightforward to verify that  $(\gamma^*, 0)$  is a KKT point. Conversely, if  $\gamma^* > 0$  satisfies the KKT conditions (3.15), by the complementary slackness, we have  $\mu^* = 0$ , which together with the stationarity implies that  $\gamma^*$  is a root of (3.16).  $\square$

In the sparse Bayesian learning, which does not use any hyperprior, i.e.,  $\mathcal{H}'(\gamma) = 0$  on  $[0, +\infty)$ , the unique solution to (3.15) is given by

$$(3.17) \quad \gamma^* = \begin{cases} 0, & \text{if } q - p^2 \geq 0, \\ \frac{p^2 - q}{q^2}, & \text{if } q - p^2 < 0. \end{cases}$$

According to Theorem 3.1, it is reasonable to choose  $\mathcal{H}$  as a strictly increasing function to promote the sparsity. This leads to the following corollary.

**COROLLARY 3.2.** *Assume  $\mathcal{H}$  is nonnegative and strictly increasing in  $[0, +\infty)$ . Then any KKT point  $\gamma^*$  of (3.15) satisfies either  $\gamma^* = 0$  or  $\gamma^* \in \left(0, \frac{p^2 - q}{q^2}\right)$ .*

*Proof.* Since  $\mathcal{L}'(\gamma) > 0$  for  $\gamma \geq \frac{p^2-q}{q^2}$ , any positive root of (3.16) must lie in  $(0, \frac{p^2-q}{q^2})$ . Based on Theorem 3.1, if  $q - p^2 > -2\mathcal{H}'(0^+)$ , then  $\gamma^* = 0$  satisfies the KKT conditions. If  $q - p^2 < -2\mathcal{H}'(0^+)$ , we have  $\mathcal{L}'(0^+) < 0$ . Together with  $\mathcal{L}'(\frac{p^2-q}{q^2}) = \mathcal{H}'(\frac{p^2-q}{q^2}) > 0$  and the continuity of  $\mathcal{L}$ , we can conclude that there exists at least one  $\gamma^* \in (0, \frac{p^2-q}{q^2})$  that is a root of (3.16), hence a KKT point. For  $q - p^2 = -2\mathcal{H}'(0^+)$ , the pair  $(\gamma^*, \mu^*) = (0, 0)$  also satisfies the KKT conditions.  $\square$

*Remark 3.3.* Compared to the case without any hyperprior, a strictly increasing  $\mathcal{H}$  promotes stronger sparsity in  $\gamma^*$  by relaxing the threshold from  $q - p^2 \geq 0$  to  $q - p^2 \geq -2\mathcal{H}'(0^+)$ . Furthermore, such a choice of  $\mathcal{H}$  also shrinks the nonzero solution  $\gamma^*$  from  $\frac{p^2-q}{q^2}$  toward values closer to zero.

If  $\mathcal{H}$  is further convex, we have the following theorem.

**THEOREM 3.4.** *Assume  $\mathcal{H}$  is nonnegative, strictly increasing and convex in  $[0, +\infty)$ . Then any  $\gamma^*$  satisfying (3.15) is a local minimizer of (3.14).*

*Proof.* The second derivative of  $\mathcal{L}$  is calculated as

$$(3.18) \quad \mathcal{L}''(\gamma) = \frac{q}{(1+q\gamma)^2} \left( \frac{p^2}{1+q\gamma} - \frac{q}{2} \right) + \mathcal{H}''(\gamma)$$

Since  $\mathcal{H}$  is nonnegative and strictly increasing, according to Corollary 3.2 we know that either  $\gamma^* = 0$  or  $0 < \gamma^* < \frac{p^2-q}{q^2}$ .

If  $\gamma^* = 0$ , then by item (I) of Theorem 3.1, we have

$$\mathcal{L}'(0^+) = \frac{q-p^2}{2} + \mathcal{H}'(0^+) \geq 0,$$

which guarantees  $\gamma^* = 0$  is a local minimizer.

If  $\gamma^* \neq 0$ , then

$$\frac{p^2}{1+q\gamma^*} - \frac{q}{2} \geq \frac{q}{2} > 0.$$

Since  $\mathcal{H}'' \geq 0$ , we get that  $\mathcal{L}''(\gamma^*) > 0$ . Therefore,  $\gamma^*$  is a local minimizer.  $\square$

*Remark 3.5.* According to Theorem 3.4, any KKT point  $\gamma^*$  of problem (2.8) is a coordinate-wise minimizer. However,  $\gamma^*$  is not necessarily a local minimizer, as ascent directions may exist outside the coordinate axes. This is due to the structure of the Hessian matrix  $\nabla^2 \mathcal{J}(\gamma^*)$ . As seen in (3.7), while the diagonal entries are positive at  $\gamma^*$ , the off-diagonal entries usually are nonzero and the positive definiteness of the Hessian at  $\gamma^*$  cannot be guaranteed. However, if  $\mathcal{H}$  is selected to be strongly convex with sufficiently large curvature such that the Hessian of  $\mathcal{J}$  becomes positive semidefinite, then  $\gamma^*$  is a global minimizer.

**3.3. General Gamma prior.** In this subsection, we consider a specific type of distribution as the hyperprior, namely the generalized Gamma distribution, which encompasses several notable probability density functions (PDFs). Our aim is to analyze the effect of incorporating this prior into the empirical Bayes framework.

The PDF of the generalized Gamma distribution is given by

$$(3.19) \quad \pi(\gamma) = \frac{|\zeta|}{\Gamma(\alpha)\beta^{\zeta\alpha}} \gamma^{\zeta\alpha-1} \exp\left(-\frac{\gamma^\zeta}{\beta^\zeta}\right), \quad \gamma \geq 0,$$

where  $\Gamma(\cdot)$  denotes the Gamma function. We require  $\zeta \neq 0$ , and  $\alpha, \beta > 0$  are the shape and scale parameters, respectively. Here, we restrict the domain to  $\gamma \in [0, +\infty)$  and, for notational convenience, define  $1/0 = \infty$ .

Under the assumption of independence, the joint hyperprior over  $\gamma$  is

$$\pi(\gamma) \propto \prod_{i=1}^n \gamma_i^{\zeta\alpha-1} \exp\left(-\frac{\gamma_i^\zeta}{\beta^\zeta}\right) = (\gamma_1 \cdots \gamma_n)^{\zeta\alpha-1} \exp\left(-\sum_{i=1}^n \frac{\gamma_i^\zeta}{\beta^\zeta}\right).$$

The corresponding cost function (2.9) becomes

$$(3.20) \quad \mathcal{J}(\gamma) = \frac{1}{2} \mathbf{y}^\top \mathbf{S}(\gamma)^{-1} \mathbf{y} + \frac{1}{2} \log \det \mathbf{S}(\gamma) - (\zeta\alpha - 1) \sum_{i=1}^n \log \gamma_i + \sum_{i=1}^n \frac{\gamma_i^\zeta}{\beta^\zeta}.$$

As in the previous subsection, we isolate a single component and define the univariate objective

$$(3.21) \quad \mathcal{L}(\gamma) = -\frac{1}{2} \frac{p^2 \gamma}{1 + q\gamma} + \frac{1}{2} \log(1 + q\gamma) + \mathcal{H}(\gamma),$$

where

$$(3.22) \quad \mathcal{H}(\gamma) = -(\zeta\alpha - 1) \log \gamma + \frac{\gamma^\zeta}{\beta^\zeta}.$$

According to the definition of  $\mathcal{H}$ , we have

$$(3.23) \quad \mathcal{H}'(\gamma) = \frac{1 - \zeta\alpha}{\gamma} + \frac{\zeta}{\beta^\zeta} \gamma^{\zeta-1}.$$

We now analyze the KKT points under different combinations of hyperparameters  $\alpha, \beta$  and  $\zeta$ , using Theorem 3.1 and Corollary 3.2.

(I)  $\zeta\alpha > 1$ : Since  $\alpha > 0$ , we must have  $\zeta > 0$ , and thus  $\mathcal{H}'(0^+) = -\infty$ . Therefore, any KKT point  $\gamma^*$  must be positive.

(II)  $\zeta\alpha < 1$ : In this case,  $\zeta$  can be either positive or negative.

- (1) If  $\zeta > 0$ , then  $\mathcal{L}(0^+) = -\infty$ , and hence  $\gamma^* = 0$  is the global minimizer.  
When  $q - p^2 < 0$ , the stationarity equation  $\mathcal{L}'(\gamma) = 0$  takes the form

$$(3.24) \quad \frac{1}{2} \frac{q^2 \gamma + (q - p^2)}{(1 + q\gamma)^2} + \frac{1 - \zeta\alpha}{\gamma} + \frac{\zeta}{\beta^\zeta} \gamma^{\zeta-1} = 0,$$

which may admit positive roots in the interval  $(0, \frac{p^2 - q}{q^2})$ . Each such root is a KKT point.

(2) If  $\zeta < 0$ , then  $\mathcal{H}'(0^+) = -\infty$ , and thus any KKT point  $\gamma^*$  must be positive.

(III)  $\zeta\alpha = 1$ : Then  $\zeta > 0$ .

(1) If  $\zeta > 1$ , then  $\mathcal{H}'(0^+) = 0$ . Hence  $\gamma^* = 0$  is a KKT point if  $q - p^2 \geq 0$ ; otherwise, there exists a KKT point  $\gamma^* \in (0, \frac{p^2 - q}{q^2})$ .

(2) If  $\zeta = 1$ , then  $\mathcal{H}'(0^+) = \frac{1}{\beta}$ . Hence  $\gamma^* = 0$  is a KKT point if  $q - p^2 \geq -\frac{2}{\beta}$ ; otherwise, a KKT point lies in  $(0, \frac{p^2 - q}{q^2})$ .

(3) If  $0 < \zeta < 1$ , then  $\mathcal{H}'(0^+) = +\infty$ , and thus  $\gamma^* = 0$  must be a KKT point.  
When  $q - p^2 < 0$ , then  $\mathcal{L}'(\gamma) = 0$  becomes

$$(3.25) \quad \frac{1}{2} \frac{q^2 \gamma + (q - p^2)}{(1 + q\gamma)^2} + \frac{\zeta}{\beta^\zeta} \gamma^{\zeta-1} = 0,$$

TABLE 1  
*Characterization of the KKT points under different hyperparameter regimes.*

$\alpha$ -condition	$\zeta$ -condition	$\mathcal{H}'(0^+)$	KKT points
$\zeta\alpha > 1$	$\zeta > 0$	$\mathcal{H}'(0^+) = -\infty$	$\gamma^* > 0$ .
$\zeta\alpha < 1$	$\zeta > 0$	$\mathcal{H}'(0^+) = +\infty$	$\gamma^* = 0$ (global minimizer); may exist other KKT points in $(0, \frac{p^2-q}{q^2})$ if $q - p^2 < 0$ .
	$\zeta < 0$	$\mathcal{H}'(0^+) = -\infty$	$\gamma^* > 0$ .
$\zeta\alpha = 1$	$\zeta > 1$	$\mathcal{H}'(0^+) = 0$	$\begin{cases} \gamma^* = 0, & \text{if } q - p^2 \geq 0, \\ \gamma^* \in \left(0, \frac{p^2-q}{q^2}\right), & \text{if } q - p^2 < 0. \end{cases}$
	$\zeta = 1$	$\mathcal{H}'(0^+) = \frac{1}{\beta}$	$\begin{cases} \gamma^* = 0, & \text{if } q - p^2 \geq -\frac{2}{\beta}, \\ \gamma^* \in \left(0, \frac{p^2-q}{q^2}\right), & \text{if } q - p^2 < -\frac{2}{\beta}. \end{cases}$
	$0 < \zeta < 1$	$\mathcal{H}'(0^+) = +\infty$	$\gamma^* = 0$ ; may exist other KKT points in $\left(0, \frac{p^2-q}{q^2}\right)$ if $q - p^2 < 0$ .

which may have positive roots in  $\left(0, \frac{p^2-q}{q^2}\right)$ , each corresponding to a KKT point.

We summarize the above results in Table 1 and provide some commonly used examples.

EXAMPLE 3.1. **Gamma hyperprior.** With a Gamma hyperprior,  $G(\alpha, \beta)$ , the PDF is

$$\pi(\gamma) = \frac{1}{\Gamma(\alpha)\beta^\alpha} \gamma^{\alpha-1} \exp\left(-\frac{\gamma}{\beta}\right), \quad \gamma \geq 0.$$

It corresponds to the cases with “ $\zeta = 1$ ” in Table 1. We arrive at the following conclusion:

- If  $\alpha > 1$ , then  $\gamma^* > 0$  necessarily holds, and thus no sparse solution can be obtained.
- If  $0 < \alpha < 1$ , we have  $\gamma^* = 0$  as the global minimizer, hence  $\gamma^* = 0$  solves (2.8) and the corresponding hyperprior promotes the sparsity.
- If  $\alpha = 1$ , the Gamma distribution reduces to the half Laplace case, see Example 3.4.

EXAMPLE 3.2. **Inverse Gamma hyperprior.** With an inverse Gamma hyperprior,  $\text{InvG}(\alpha, \beta)$ , the PDF is

$$\pi(\gamma) = \frac{\beta^\alpha}{\Gamma(\alpha)} \gamma^{-(\alpha+1)} \exp\left(-\frac{\beta}{\gamma}\right), \quad \gamma \geq 0.$$

It corresponds to the case “ $\zeta < 0$ ” in Table 1 with  $\zeta = -1$ . Hence,  $\gamma^*$  must be positive, and we can conclude that adding an inverse Gamma hyperprior would reduce the sparsity of  $\gamma$ .

EXAMPLE 3.3. **Half-Gaussian hyperprior.** The PDF of half-Gaussian prior  $\mathcal{N}^+(0, \theta^2)$  is

$$(3.26) \quad \pi(\gamma) = \frac{\sqrt{2}}{\sqrt{\pi}\theta} \exp\left(-\frac{\gamma^2}{2\theta^2}\right), \quad \gamma \geq 0.$$

It corresponds to the case “ $\zeta\alpha = 1$  and  $\zeta > 1$ ” in Table 1 with  $\zeta = 2$ ,  $\alpha = \frac{1}{2}$ ,  $\beta = \sqrt{2}\theta$ . If  $q - p^2 \geq 0$ , then  $\gamma^* = 0$  is the global minimizer. If  $q - p^2 < 0$ , then the stationarity condition  $\mathcal{L}'(\gamma) = 0$  is equivalent to the cubic equation

$$(3.27) \quad 4q^2\gamma^3 + 8q\gamma^2 + (4 + \beta^2 q^2)\gamma + \beta^2(q - p^2) = 0.$$

By Vieta’s relations one can show that (3.27) has exactly one positive real root in  $(0, \frac{p^2 - q}{q^2})$ , which is the global minimizer in this case. A closed form is available via the Cardano formula. Comparing with the result without the hyperprior given in (3.17), we can see that half-Gaussian hyperprior won’t change the sparsity but can push  $\gamma^*$  closer to zero.

**EXAMPLE 3.4. Half-Laplace hyperprior.** With a half-Laplace hyperprior, the PDF is

$$\pi(\gamma) = \frac{1}{\beta} \exp\left(-\frac{\gamma}{\beta}\right), \quad \gamma \geq 0,$$

which corresponds to the case “ $\zeta\alpha = 1$  and  $\zeta = 1$ ” in Table 1 with  $\mathcal{H}(\gamma) = \frac{\gamma}{\beta}$ . Since  $\beta > 0$ , the condition promoting  $\gamma = 0$  is weaker than the no-hyperprior condition  $q - p^2 \geq 0$ . Thus adding a half-Laplace prior can promote sparsity. When  $q - p^2 < -\frac{2}{\beta}$ , the equation  $\mathcal{L}'(\gamma) = 0$  reduces to

$$(3.28) \quad 2q^2\gamma^2 + (4q + \beta q^2)\gamma + 2 + \beta(q - p^2) = 0,$$

which has a unique positive root that can be explicitly given by the quadratic formula. Therefore, the global minimizer is

$$\gamma^* = \begin{cases} 0, & \text{if } q - p^2 \geq -\frac{2}{\beta}, \\ \frac{-(4 + \beta q) + \sqrt{\beta^2 q^2 + 8\beta p^2}}{4q}, & \text{if } q - p^2 < -\frac{2}{\beta}. \end{cases}$$

**EXAMPLE 3.5. Half-Generalized Gaussian hyperprior with  $0 < \zeta < 1$ .** For the half-Generalized Gaussian hyperprior with  $0 < \zeta < 1$ , the PDF is

$$(3.29) \quad \pi(\gamma) = \frac{\zeta}{\Gamma(\alpha)\beta} \exp\left(-\left(\frac{\gamma}{\beta}\right)^\zeta\right), \quad \gamma \geq 0.$$

It corresponds to the case “ $\zeta\alpha = 1$  and  $0 < \zeta < 1$ ” in Table 1 with a concave  $\mathcal{H}(\gamma) = (\frac{\gamma}{\beta})^\zeta$ . Since  $\mathcal{H}'(0^+) = +\infty$ , one always has  $\gamma^* = 0$  as a local minimizer of  $\mathcal{L}(\gamma)$  in  $[0, +\infty)$ . And any positive root of the equation

$$(3.30) \quad 2\zeta q^2\gamma^{\zeta+1} + 4\zeta q\gamma^\zeta + 2\zeta\gamma^{\zeta-1} + \beta^\zeta q^2\gamma + \beta^\zeta(q - p^2) = 0,$$

is a KKT point. In particular, when  $\zeta = \frac{1}{2}$ , any positive solution  $s$  of the quartic equation

$$(3.31) \quad q^2 s^4 + \sqrt{\beta} q^2 s^3 + 2q s^2 + \sqrt{\beta}(q - p^2)s + 1 = 0,$$

yields a KKT point with  $\gamma = s^2$ . Therefore, the global minimizer depends on the specific values of  $\beta$ ,  $p$ , and  $q$ , illustrating the added complexity induced by a nonconvex  $\mathcal{H}$ .

**4. Numerical algorithm.** In this section, referring to the approach in [38], we present a numerical algorithm for minimizing (2.9) within the empirical Bayes framework and analyze its convergence.

**4.1. Proximal alternating linearized minimization.** We begin by reformulating the problem using an auxiliary variable, which enables an alternating minimization scheme between the original variables and the auxiliary variable. This reformulation allows us to apply a proximal alternating linearized minimization algorithm [3] to efficiently compute a stationary point, even in the presence of nonconvex terms.

Following the idea in [38], we construct a strict upper-bounding auxiliary function  $\mathcal{F}$  for the first term in  $\mathcal{J}(\gamma)$ :

$$(4.1) \quad \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma}) = \sigma \|\mathbf{F}\mathbf{x} - \mathbf{y}\|^2 + \mathbf{x}^\top \boldsymbol{\Gamma}^\dagger \mathbf{x} + \sum_{i \in I_\gamma} \iota_{\{0\}}(x_i), \quad \mathbf{x} \in \mathbb{R}^n, \boldsymbol{\gamma} \in \mathbb{R}_+^n,$$

where  $\boldsymbol{\Gamma}^\dagger$  is the Moore–Penrose inverse of  $\boldsymbol{\Gamma}$ , i.e.,

$$(4.2) \quad [\boldsymbol{\Gamma}^\dagger]_{ii} = \begin{cases} \gamma_i^{-1}, & \text{if } \gamma_i \neq 0, \\ 0, & \text{if } \gamma_i = 0, \end{cases}$$

$I_\gamma = \{1 \leq i \leq n : \gamma_i = 0\}$ , and  $\iota_A(x)$  is the indicator function defined as

$$(4.3) \quad \iota_A(x) = \begin{cases} +\infty, & \text{if } x \notin A, \\ 0, & \text{if } x \in A. \end{cases}$$

For any  $\boldsymbol{\gamma} \in \mathbb{R}_+^n$  and  $\mathbf{y} \in \mathbb{R}^m$ , the optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma})$$

admits a unique minimizer  $\mathbf{x}^*(\boldsymbol{\gamma})$  given by

$$(4.4) \quad \mathbf{x}^*(\boldsymbol{\gamma}) = \boldsymbol{\Gamma} \mathbf{F}^\top (\mathbf{S}(\boldsymbol{\gamma}))^{-1} \mathbf{y} \quad \text{and} \quad \mathcal{F}(\mathbf{x}^*(\boldsymbol{\gamma}), \boldsymbol{\gamma}) = \mathbf{y}^\top (\mathbf{S}(\boldsymbol{\gamma}))^{-1} \mathbf{y}.$$

Therefore,  $\mathcal{J}$  in (2.9) can be equivalently expressed as

$$(4.5) \quad \mathcal{J}(\boldsymbol{\gamma}) = \min_{\mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma}) + g(\boldsymbol{\gamma}),$$

where

$$(4.6) \quad g(\boldsymbol{\gamma}) = \frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}) + \sum_{i=1}^n \mathcal{H}(\gamma_i).$$

Minimizing  $\mathcal{J}$  on  $\boldsymbol{\gamma} \in \mathbb{R}_+^n$  leads to the joint minimization problem

$$(4.7) \quad \min_{\boldsymbol{\gamma} \in \mathbb{R}_+^n, \mathbf{x} \in \mathbb{R}^n} \frac{1}{2} \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma}) + g(\boldsymbol{\gamma}).$$

A natural approach for such two-block minimization problems is the alternating minimization algorithm (AMA). Specifically, starting with an initial guess  $\boldsymbol{\gamma}^{(0)}$ , AMA for the above problem alternates between the following updates

$$(4.8) \quad \mathbf{x} - \text{subproblem :} \quad \mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma}^{(k)}),$$

$$(4.9) \quad \boldsymbol{\gamma} - \text{subproblem :} \quad \boldsymbol{\gamma}^{(k+1)} = \arg \min_{\boldsymbol{\gamma} \in \mathbb{R}_+^n} \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}) + g(\boldsymbol{\gamma}).$$

According to (4.4), the  $\mathbf{x}$ -subproblem has a closed-form solution

$$(4.10) \quad \mathbf{x}^{(k+1)} = \mathbf{\Gamma}^{(k)} \mathbf{F}^\top (\mathbf{S}(\boldsymbol{\gamma}^{(k)}))^{-1} \mathbf{y}.$$

However, the  $\boldsymbol{\gamma}$ -subproblem is substantially more challenging: the log-determinant term is nonconvex, and  $\mathcal{H}$  may also be nonconvex, which prevents a closed-form solution and making the global minimizer computationally demanding.

To address this difficulty, we replace  $g(\boldsymbol{\gamma})$  with a surrogate function  $\tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)})$ , leading to the PALM update

(4.11)

$$\mathbf{x} - \text{subproblem} : \quad \mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \mathcal{F}(\mathbf{x}, \boldsymbol{\gamma}^{(k)}),$$

(4.12)

$$\boldsymbol{\gamma} - \text{subproblem} : \quad \boldsymbol{\gamma}^{(k+1)} = \arg \min_{\boldsymbol{\gamma} \in \mathbb{R}_+^n} \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}) + \tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)}) + \frac{\tau}{2} \|\boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)}\|^2,$$

where  $\tau > 0$  is the proximal parameter. The choice of the surrogate function  $\tilde{g}$  is critical for both computational efficiency and convergence. Based on the property of  $\mathcal{H}$ , we present two cases:

- for a convex  $\mathcal{H}$ :

(4.13)

$$\tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)}) = \frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}) + \frac{1}{2} \langle \nabla \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}), \boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)} \rangle + \sum_{i=1}^n \mathcal{H}(\gamma_i),$$

- for a nonconvex  $\mathcal{H}$ :

$$(4.14) \quad \tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)}) = g(\boldsymbol{\gamma}^{(k)}) + \langle \nabla g(\boldsymbol{\gamma}^{(k)}), \boldsymbol{\gamma} - \boldsymbol{\gamma}^{(k)} \rangle.$$

**4.2. Solving the  $\boldsymbol{\gamma}$ -subproblem.** Before discussing the solver of the  $\boldsymbol{\gamma}$ -subproblem, we can show the sequence  $\{(\mathbf{x}^{(k)}, \boldsymbol{\gamma}^{(k)})\}$  generated by PALM satisfies the following proposition.

**PROPOSITION 4.1.** *Let  $\mathcal{H}$  be nonnegative, and let  $\{(\mathbf{x}^{(k)}, \boldsymbol{\gamma}^{(k)})\}$  be generated by PALM (4.11)-(4.12) with an initial  $\boldsymbol{\gamma}^{(0)}$ . If we have  $\gamma_i^{(K)} = 0$  for some  $K > 0$  and  $1 \leq i \leq n$ . Then  $x_i^{(k)} = \gamma_i^{(k)} = 0$  for all  $k > K$ .*

*Proof.* It suffices to prove that  $\gamma_i^{(K)} = 0$  implies  $x_i^{(K+1)} = 0$  and  $\gamma_i^{(K+1)} = 0$ . From (4.10),  $\gamma_i^{(K)} = 0$  directly yields  $x_i^{(K+1)} = 0$ .

For the  $\boldsymbol{\gamma}$ -subproblem given in (4.12), note that when  $\gamma_i^{(K)} = 0$  and  $x_i^{(K+1)} = 0$ , replacing a positive  $\gamma_i$  by zero decreases the terms  $(\mathbf{x}^{(K+1)})^\top \mathbf{\Gamma}^\dagger \mathbf{x}^{(K+1)}$ ,  $\log \det \mathbf{S}(\boldsymbol{\gamma})$  and  $\|\boldsymbol{\gamma} - \boldsymbol{\gamma}^{(K)}\|^2$ , while leaving  $\sum_{i \in I_\gamma} \iota_{\{0\}}(x_i^{(K+1)})$  unchanged. Hence,  $\gamma_i^{(K+1)} = 0$  for such  $i$ .  $\square$

Proposition 4.1 shows that once  $\gamma_i^{(k)} = 0$  for some  $i$ , both  $x_i$  and  $\gamma_i$  will remain zero in all subsequent iterations. Consequently, these components can be removed from the optimization, and only the remaining indices need to be updated. Moreover, if  $x_i^{(k)} \neq 0$ , then necessarily  $\gamma_i^{(k+1)} > 0$ ; otherwise, the term  $\sum_{i \in I_\gamma} \iota_{\{0\}}(x_i^{(k)})$  would become infinite. Therefore, in the following we only focus on those  $\gamma_i$ 's with  $\gamma_i^{(k)} > 0$  and  $x_i^{(k)} \neq 0$ .

---

**Algorithm 4.1** The PALM algorithm for minimizing  $\mathcal{J}(\boldsymbol{\gamma})$  defined in (2.9).

---

**Input:**  $\mathbf{y}$ ,  $\mathbf{F}$ ,  $\sigma$ ,  $\tau$  and  $\omega$ ;

**Output:**  $\mathbf{x}^{(k+1)}$  and  $\boldsymbol{\gamma}^{(k+1)}$ ;

```

1: Initialize  $\boldsymbol{\gamma}^{(0)}$  and set  $k = 0$ 
2: while  $\mathbf{x}^{(k)}$  not converged do
3:    $\mathbf{x}$ -update:  $\mathbf{x}^{(k+1)} = \mathbf{\Gamma}^{(k)} \mathbf{F}^\top (\mathbf{S}(\boldsymbol{\gamma}^{(k)}))^{-1} \mathbf{y}$ 
4:    $\boldsymbol{\gamma}$ -update:
5:   for  $i = 1, \dots, n$  do
6:     if  $\gamma_i^{(k)} = 0$  then
7:       set  $\gamma_i^{(k+1)} = 0$ 
8:     else
9:        $\gamma_i^{(k+1)} = \arg \min_{\gamma_i \in \mathbb{R}_{++}} \frac{(x_i^{(k+1)})^2}{2\gamma_i} + [\tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)})]_i + \frac{\tau}{2}(\gamma_i - \gamma_i^{(k)})^2$ 
10:      if  $\gamma_i^{(k+1)} < \omega$  then
11:        set  $\gamma_i^{(k+1)} = 0$ 
12:      end if
13:    end if
14:   end for
15:    $k = k + 1$ 
16: end while
17: return  $\mathbf{x}^{(k+1)}$  and  $\boldsymbol{\gamma}^{(k+1)}$ 

```

---

Due to the separability of the objective in (4.12) with respect to  $\gamma_i$ , the update of  $\boldsymbol{\gamma}^{(k+1)}$  reduces to solving a set of independent univariate problems:

$$(4.15) \quad \gamma_i^{(k+1)} = \arg \min_{\gamma_i > 0} \frac{(x_i^{(k+1)})^2}{2\gamma_i} + \mathcal{L}(\gamma_i) + \frac{\tau}{2}(\gamma_i - \gamma_i^{(k)})^2,$$

where

$$\mathcal{L}(\gamma_i) = \begin{cases} \frac{1}{2}\tilde{q}_{ii}^{(k)}\gamma_i + \mathcal{H}(\gamma_i), & \text{if } \mathcal{H} \text{ is convex,} \\ \left(\frac{1}{2}\tilde{q}_{ii}^{(k)} + \mathcal{H}'(\gamma_i^{(k)})\right)\gamma_i, & \text{if } \mathcal{H} \text{ is nonconvex,} \end{cases}$$

and

$$\tilde{q}_{ii}^{(k)} = \mathbf{f}_i^\top \mathbf{S}(\boldsymbol{\gamma}^{(k)})^{-1} \mathbf{f}_i.$$

For the hyperpriors considered in Example 3.1–3.5, the update formula (4.15) reduces to finding the positive root of a cubic polynomial. In this case,  $\gamma_i^{(k+1)}$  can be computed exactly via Cardano's formula, enabling an efficient closed-form update of  $\boldsymbol{\gamma}$  without requiring any inner iterations.

**4.3. Convergence analysis.** In this subsection, we study the convergence of the proposed PALM algorithm. First, we summarize the PALM algorithm to minimize  $\mathcal{J}(\boldsymbol{\gamma})$  in Algorithm 4.1. For the convergence analysis we restrict  $\mathcal{H}$  to be either convex or concave. This condition is satisfied by the hyperpriors in Example 3.1, and Example 3.3–3.5, which are well-suited for promoting sparsity. The convergence results of PALM are given in the following theorem.

**THEOREM 4.2.** *Let  $\tau > 0$ ,  $\mathcal{H}$  be nonnegative and either convex or concave, and  $\{(\mathbf{x}^{(k)}, \boldsymbol{\gamma}^{(k)})\}$  be generated by PALM given in Algorithm 4.1 with an initial positive  $\boldsymbol{\gamma}^{(0)}$ , i.e.,  $\boldsymbol{\gamma}^{(0)} \in \mathbb{R}_{++}^n$ . Then the following convergence results hold:*

- (I) The sequence  $\{\mathcal{J}(\boldsymbol{\gamma}^{(k)})\}$  is strictly decreasing and convergent.
- (II) There exists a subsequence  $\{s_k\} \subset \mathbb{N}$  such that  $\{\mathbf{x}^{(s_k)}\}$  and  $\{\boldsymbol{\gamma}^{(s_k)}\}$  converge.
- (III) The limit  $\boldsymbol{\gamma}^*$  of  $\{\boldsymbol{\gamma}^{(s_k)}\}$  is a KKT point of (2.8), that is, for each  $1 \leq i \leq n$ , either  $\gamma_i^* = 0$  or  $\frac{\partial \mathcal{J}(\boldsymbol{\gamma}^*)}{\partial \gamma_i} = 0$  holds.

*Proof.* Proof of (I). From (4.5) and the  $\mathbf{x}$ -subproblem in (4.11), we obtain

$$(4.16) \quad \mathcal{J}(\boldsymbol{\gamma}^{(k)}) = \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k)}) + g(\boldsymbol{\gamma}^{(k)})$$

and

$$(4.17) \quad \mathcal{J}(\boldsymbol{\gamma}^{(k+1)}) = \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+2)}, \boldsymbol{\gamma}^{(k+1)}) + g(\boldsymbol{\gamma}^{(k+1)}) \leq \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)}) + g(\boldsymbol{\gamma}^{(k+1)}).$$

According to (4.6) and the concavity of  $\log \det \mathbf{S}(\boldsymbol{\gamma})$ , we have

$$(4.18) \quad g(\boldsymbol{\gamma}^{(k+1)}) \leq \frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}) + \frac{1}{2} \langle \nabla \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}), \boldsymbol{\gamma}^{(k+1)} - \boldsymbol{\gamma}^{(k)} \rangle + \sum_{i=1}^n \mathcal{H}(\gamma_i^{(k+1)}).$$

According to the definition of  $\tilde{g}(\boldsymbol{\gamma}, \boldsymbol{\gamma}^{(k)})$  in the case of a convex  $\mathcal{H}$ , we obtain

$$(4.19) \quad g(\boldsymbol{\gamma}^{(k+1)}) \leq \tilde{g}(\boldsymbol{\gamma}^{(k+1)}, \boldsymbol{\gamma}^{(k)}).$$

For a concave  $\mathcal{H}$ , based on (4.14) and the concavity of  $g(\boldsymbol{\gamma})$  we can obtain (4.19) in the similar way.

From the  $\boldsymbol{\gamma}$ -subproblem in (4.12), it follows that

$$\begin{aligned} & \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)}) + \tilde{g}(\boldsymbol{\gamma}^{(k+1)}, \boldsymbol{\gamma}^{(k)}) + \frac{\tau}{2} \|\boldsymbol{\gamma}^{(k+1)} - \boldsymbol{\gamma}^{(k)}\|^2 \\ & \leq \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k)}) + \tilde{g}(\boldsymbol{\gamma}^{(k)}, \boldsymbol{\gamma}^{(k)}) \\ & = \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k)}) + \frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}) + \sum_{i=1}^n \mathcal{H}(\gamma_i^{(k)}) \\ & = \frac{1}{2} \mathcal{F}(\mathbf{x}^{(k+1)}, \boldsymbol{\gamma}^{(k)}) + g(\boldsymbol{\gamma}^{(k)}). \end{aligned}$$

Combining the above inequalities yields

$$(4.20) \quad \mathcal{J}(\boldsymbol{\gamma}^{(k+1)}) \leq \mathcal{J}(\boldsymbol{\gamma}^{(k)}) - \frac{\tau}{2} \|\boldsymbol{\gamma}^{(k+1)} - \boldsymbol{\gamma}^{(k)}\|^2.$$

Hence,  $\{\mathcal{J}(\boldsymbol{\gamma}^{(k)})\}$  is strictly decreasing. Moreover, by the definition of  $\mathcal{J}(\boldsymbol{\gamma})$  in (2.9), we know it is always bounded below by  $-\frac{1}{2} \log \sigma$ . Therefore, the sequence  $\{\mathcal{J}(\boldsymbol{\gamma}^{(k)})\}$  converges.

Proof of (II). We now show that  $\{\boldsymbol{\gamma}^{(k)}\}$  is bounded. Observe that

$$\log \det \mathbf{S}(\boldsymbol{\gamma}) = \log \det \left( \sigma^{-1} \mathbf{I}_m + \sum_{i=1}^n \gamma_i \mathbf{f}_i \mathbf{f}_i^\top \right)$$

is coercive, i.e.,  $\lim_{\|\boldsymbol{\gamma}\| \rightarrow +\infty} \log \det \mathbf{S}(\boldsymbol{\gamma}) = +\infty$ . Moreover,

$$\frac{1}{2} \log \det \mathbf{S}(\boldsymbol{\gamma}^{(k)}) \leq \mathcal{J}(\boldsymbol{\gamma}^{(k)}) \leq \mathcal{J}(\boldsymbol{\gamma}^{(0)}), \quad \forall k \geq 0.$$

which implies boundedness of  $\{\boldsymbol{\gamma}^{(k)}\}$ . Thus, there exists a subsequence  $s_k$  such that  $\{\boldsymbol{\gamma}^{(s_k)}\}$  converges. Due to the closed-form expression of  $\mathbf{x}^{(k)}$  in (4.10), we conclude that  $\{\mathbf{x}^{(s_k)}\}$  also converges.

Proof of (III). Taking the limit in the sufficiently descent inequality (4.20), we find

$$\lim_{k \rightarrow +\infty} \|\boldsymbol{\gamma}^{(k+1)} - \boldsymbol{\gamma}^{(k)}\|^2 = 0.$$

Due to the convergence of  $\{\boldsymbol{\gamma}^{(s_k)}\}$ , we define the limit point  $\boldsymbol{\gamma}^*$ :

$$\lim_{s_k \rightarrow +\infty} \boldsymbol{\gamma}^{(s_k)} = \boldsymbol{\gamma}^*.$$

For any  $1 \leq i \leq n$ , if  $\gamma_i^{(s_k)} = 0$  for some  $s_k \in \mathbb{N}$ , based on Proposition 4.1 we have  $\gamma_i^* = 0$ . Otherwise, the stationarity condition for (4.15) gives

$$(4.21) \quad -\frac{1}{2} \left( \frac{x_i^{(s_k)}}{\gamma_i^{(s_k+1)}} \right)^2 + \frac{1}{2} \tilde{q}_{ii}^{(s_k)} + \mathcal{H}'(\gamma_i^{(s_k+1)}) + \tau(\gamma_i^{(s_k+1)} - \gamma_i^{(s_k)}) = 0, \text{ if } \mathcal{H} \text{ is convex},$$

$$(4.22) \quad -\frac{1}{2} \left( \frac{x_i^{(s_k)}}{\gamma_i^{(s_k+1)}} \right)^2 + \frac{1}{2} \tilde{q}_{ii}^{(s_k)} + \mathcal{H}'(\gamma_i^{(s_k)}) + \tau(\gamma_i^{(s_k+1)} - \gamma_i^{(s_k)}) = 0, \text{ if } \mathcal{H} \text{ is concave}.$$

Define  $\tilde{p}_i^{(k)} = \mathbf{f}_i^\top \mathbf{S}(\boldsymbol{\gamma}^{(k)})^{-1} \mathbf{y}$ . By (4.10), we have  $x_i^{(s_k)} = \gamma_i^{(s_k)} \tilde{p}_i^{(s_k)}$ . Substituting it into (4.21) and (4.22), after taking the limit on both sides we obtain

$$\frac{1}{2} \tilde{q}_{ii}^* - \frac{1}{2} (\tilde{p}_i^*)^2 + \mathcal{H}'(\gamma_i^*) = 0.$$

From the expression of  $\frac{\partial \mathcal{J}(\boldsymbol{\gamma})}{\partial \gamma_i}$  in (3.6), we conclude that  $\frac{\partial \mathcal{J}(\boldsymbol{\gamma}^*)}{\partial \gamma_i} = 0$ .  $\square$

**5. Numerical experiments.** In this section, we present some numerical results that illustrate how the choice of the hyperprior from the generalized Gamma family affects the sparsity of the solution to (2.8).

We consider a 2D image deblurring problem with a Gaussian blurring operator  $\mathbf{K}$  as the test problem. We assume that the inverse variance  $\sigma$  of the noise and the standard deviation  $\sigma_{ker}$  of the Gaussian blurring kernel are known. In practice, images are often not sparse in the spatial domain but can exhibit sparsity under suitable bases or dictionaries. In our tests, we assume that the test images have a sparse representation under the Discrete Cosine Transform (DCT), i.e., we have

$$(5.1) \quad \mathbf{x} = \mathbf{R}\mathbf{z},$$

where  $\mathbf{z}$  denotes the test image,  $\mathbf{x}$  is the sparse DCT coefficients, and  $\mathbf{R}$  represents DCT that satisfies

$$\mathbf{R}^\top \mathbf{R} = \mathbf{R} \mathbf{R}^\top = \mathbf{I}.$$

Then, we rewrite our forward problem according to (2.1) as

$$(5.2) \quad \mathbf{y} = \mathbf{K} \mathbf{R}^\top \mathbf{x} + \epsilon$$

with  $\mathbf{F} = \mathbf{K} \mathbf{R}^\top$ .

In Algorithm 4.1, most of the computational cost lies in the calculation of the main diagonal of  $\mathbf{F}^\top \mathbf{S}(\boldsymbol{\gamma})^{-1} \mathbf{F}$  in the  $\boldsymbol{\gamma}$ -update. Applying the Woodbury matrix identity [16], we have

$$\begin{aligned} \mathbf{F}^\top \mathbf{S}(\boldsymbol{\gamma})^{-1} \mathbf{F} &= \mathbf{F}^\top (\sigma^{-1} \mathbf{I} + \mathbf{F} \boldsymbol{\Gamma} \mathbf{F}^\top)^{-1} \mathbf{F} \\ (5.3) \quad &= \sigma \mathbf{F}^\top \mathbf{F} - \sigma \mathbf{F}^\top \mathbf{F} (\boldsymbol{\Gamma}^{-1} + \sigma \mathbf{F}^\top \mathbf{F})^{-1} \sigma \mathbf{F}^\top \mathbf{F}. \end{aligned}$$

Under the symmetric boundary condition, the blurring matrix  $\mathbf{K}$  can be diagonalized by DCT [24], i.e.,  $\mathbf{R} \mathbf{K} \mathbf{R}^\top = \boldsymbol{\Lambda}$ , where  $\boldsymbol{\Lambda}$  is a diagonal matrix. Further, we obtain

$$\mathbf{F}^\top \mathbf{F} = \mathbf{R} (\mathbf{R}^\top \boldsymbol{\Lambda}^2 \mathbf{R}) \mathbf{R}^\top = \boldsymbol{\Lambda}^2.$$

Then, according to (5.3)  $\mathbf{F}^\top \mathbf{S}(\boldsymbol{\gamma})^{-1} \mathbf{F}$  becomes a diagonal matrix and can be easily calculated.

In order to prevent unnecessary calculations caused by too small  $\gamma_i^{(k)}$ , we introduce a threshold  $\omega$  in Algorithm 4.1. When  $\gamma_i^{(k)} < \omega$ , we set  $\gamma_i^{(k)} = 0$ . According to Proposition 4.1, we can neglect these zero components from the optimization problem. In all tests, we set  $\omega = 10^{-16}$ . In addition, the initialization of  $\boldsymbol{\gamma}$  is set as the absolute value of the DCT coefficients of the degraded image and the stopping criteria of Algorithm 4.1 are

$$\frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|}{\|\mathbf{x}^{(k)}\|} < 10^{-8} \quad \text{and} \quad k \leq 200.$$

To quantitatively evaluate the performance of different hyperpriors, we use the relative error and the sparsity rate as the main evaluation metrics. The relative error between the image  $\hat{\mathbf{z}}$  and the ground truth  $\mathbf{z}$  is defined as  $\|\hat{\mathbf{z}} - \mathbf{z}\|/\|\mathbf{z}\|$ , and the sparsity rate of  $\mathbf{x}$  is defined as the ratio of zero elements in  $\mathbf{x}$  to its total number of elements.

**5.1. Impact of different hyperpriors.** In this subsection, we test the performance and characteristics of EBF with different hyperpriors compared with SBL (without any hyperprior) to illustrate the influence of the hyperpriors on the sparsity. We use the 256-by-256 gray image *Cameraman* as the test image, and the intensity range is  $[0, 1]$ . Since *Cameraman* is not sufficiently sparse under DCT, we truncate its DCT coefficients at 0.025, which gives the sparsity rate as 54.94%, and the relative error to the original image is 0.0172. In addition, we set  $\sigma_{ker} = 1$  and the noise level as 10% for all tests. Figure 1 shows the original image, the compressed image that is used as the ground truth, the degraded image and their DCT coefficient maps. In each map, frequencies increase from top to bottom and left to right. The absolute value of each DCT coefficient indicates the contribution of the corresponding frequency component.

**5.1.1. Comparison of the half-Laplace and half-Gaussian hyperpriors.** In this numerical test, we illustrate the influence of using the half-Laplace and half-Gaussian hyperpriors on the solutions of (4.7), and further compare the effect of these two hyperpriors with the one without hyperprior, i.e., SBL. To ensure that the strengths of two hyperpriors are comparable, we set  $\beta = \theta = 0.1$ .

In Figure 2, it is obvious that the half-Laplace hyperprior significantly promotes the sparsity of the DCT coefficients, making the sector-shaped region of non-zero coefficients more compact, which leads to fewer artifacts in the restoration. The half-Gaussian hyperprior produces a smoother decay of the coefficients from low frequency to high frequency, but keeps many small non-zero components, resulting in more

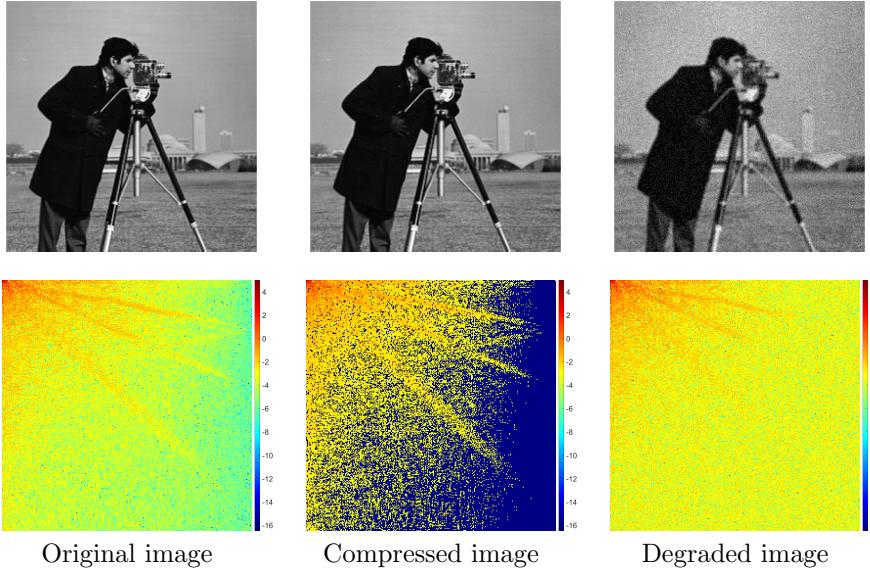


FIG. 1. *The Cameraman image. Top row: test images; Bottom row: absolute values of DCT coefficients.*

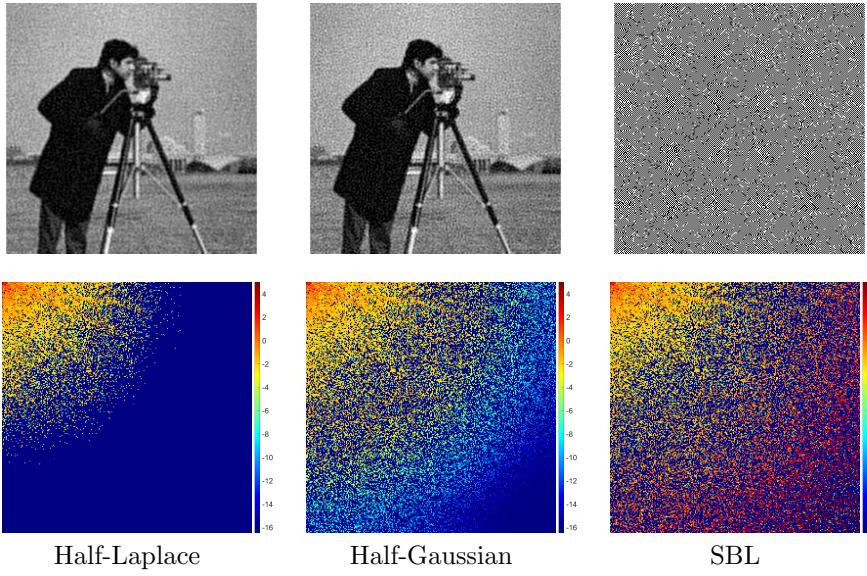


FIG. 2. *Comparisons of the results using the half-Laplace, half-Gaussian and no hyperprior. Top row: restored images; Bottom row: absolute values of DCT coefficients.*

artifacts in the restoration compared with the half-Laplace result. In contrast, the SBL method without any hyperprior exhibits numerous scattered abnormal peaks in the coefficient map, leading to high-frequency noise and resulting in obvious speckles and artifacts in the restoration. According to the relative errors listed in Table 2, we can see that the half-Laplace hyperprior provides the best results quantitatively. Comparing the sparsity level, i.e. the percentage of the zero elements in  $\mathbf{x}$ , in Table 2

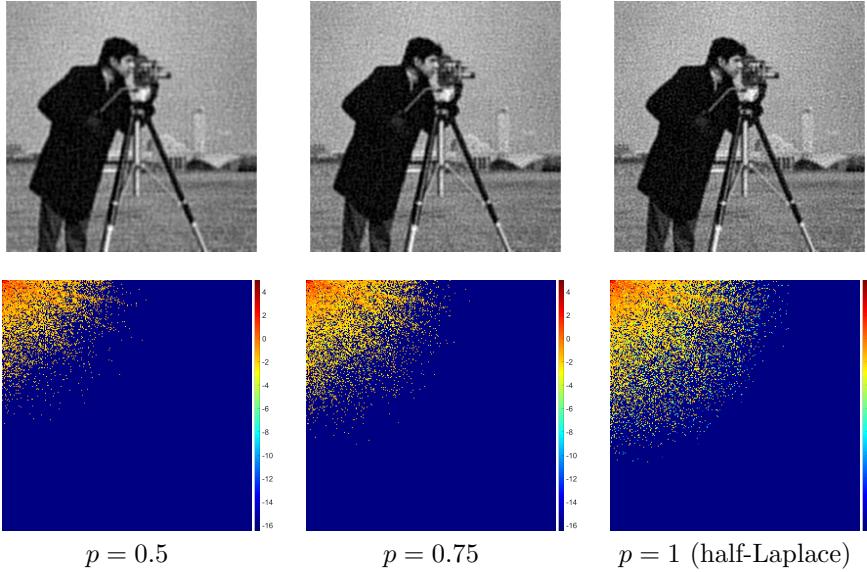


FIG. 3. Comparisons of the results using the half-generalized Gaussian hyperpriors with different  $p$ . Top row: restored images; Bottom row: absolute values of DCT coefficients.

we conclude that the half-Laplace hyperprior can promote the sparsity. The half-Gaussian hyperprior can stabilize the problem by reducing the contributions of high-frequency components, but cannot promote sparsity.

**5.1.2. Influence of the half-generalized Gaussian hyperpriors.** To investigate the influence of the half-generalized Gaussian hyperprior on the sparsity of the DCT coefficients, we test it with different choice of  $p$ , which controls the sharpness and tail thickness of the hyperprior. When  $p = 1$ , the hyperprior reduces to the half-Laplace hyperprior. When  $p = 2$ , it becomes the half-Gaussian hyperprior. To obtain a prior that strongly promotes sparsity, we focus on the case  $0 < p \leq 1$ . In this test, we set  $\beta = 0.1$ . Figure 3 shows the restored images and the DCT coefficients with  $p = 0.5, 0.75$  and  $1$ .

In the half-generalized Gaussian distribution, the smaller the  $p$ -value, the sharper the peak of the distribution. Therefore, the half-generalized Gaussian hyperprior with a smaller  $p$  has a stronger capability to promote sparsity. This property can be easily observed by comparing the DCT coefficients shown in Figure 3, where more and more high-frequency coefficients are driven to zero as  $p$  decreasing. The sparsity levels reported in Table 2 also confirm this trend. However, setting  $p$  too small (e.g.,  $p = 0.5$ ) can lead to oversmoothing and loss of image detail. As shown in the figure, the restoration for  $p = 0.5$  appears overly smooth compared to the results with larger  $p$ .

Although the half-generalized Gaussian hyperpriors with  $0 < p < 1$  can offer excellent sparse solutions, they correspond to a nonconvex  $\mathcal{H}$ , making analysis and computation more challenging. In contrast, with the half-Laplace hyperprior, i.e.,  $p = 1$ , we can obtain the similar restoration quality and sparsity level, and the corresponding  $\mathcal{H}$  is convex. Therefore, in subsection 5.2 we choose the half-Laplace hyperprior to illustrate the performance of EBF with respect to the ill-posedness and noise.

TABLE 2  
*Restoration performance of half-generalized Gaussian hyperpriors with different  $p$ .*

		$\ell_p$ Hyperprior			
	None (SBL)	$p = 2$ (Gaussian)	$p = 1$ (Laplace)	$p = 0.75$	$p = 0.5$
Relative error	0.5246	0.1279	0.1055	0.1017	0.1041
Sparsity (%)	55.95	56.36	84.73	90.32	93.25

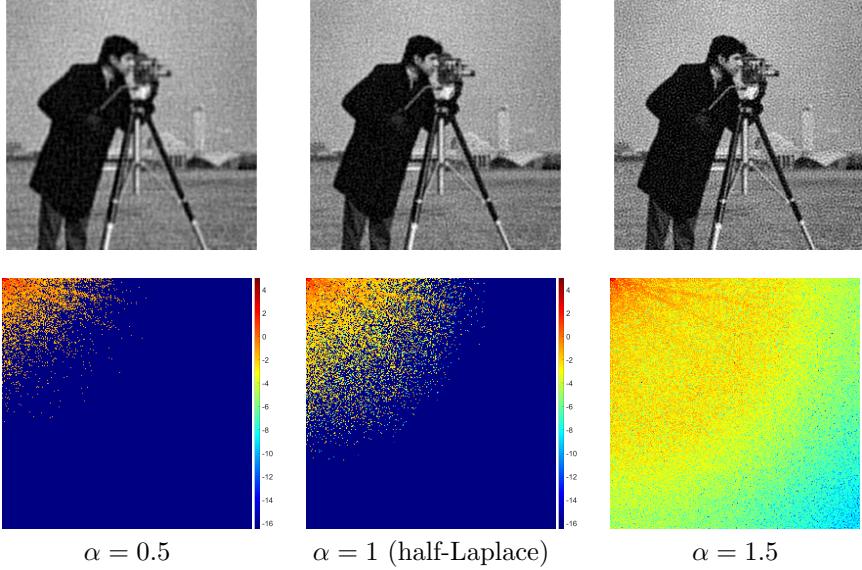


FIG. 4. Comparisons of the results using the Gamma hyperprior with different  $\alpha$ . Top row: restored images with relative errors of 0.1084, 0.1055 and 0.1423, respectively; Bottom row: absolute values of DCT coefficients with sparsity rates of 94, 74%, 84.73% and 0%, respectively.

**5.1.3. Influence of the Gamma hyperprior.** Before discussing the performance of EBF equipped with half-Laplace, we test EBFs with the Gamma hyperprior. According to Example 3.1, we know that the relationship between  $\alpha$  and 1 determines the sparsity-promoting effect of the Gamma hyperprior. In this test, we fix  $\beta = 0.1$  and show how the hyperparameter  $\alpha$  affects the sparsity of the solution. Figure 4 shows the restored images and the corresponding DCT coefficient maps for  $\alpha = 0.5$ , 1 and 1.5, respectively.

In the Gamma hyperprior, when  $0 < \alpha < 1$ ,  $\gamma^* = 0$  is always a global minimizer. Therefore, the Gamma hyperprior is a strongly sparsity-promoting prior in this case, thereby pushing all  $\gamma_i$  values towards 0. However, there are still quite a few non-zeros  $\gamma_i$  in the low-frequency region in Figure 4 because the algorithm converges to other KKT points instead of the global minimizer. When  $\alpha > 1$ , many mid- and high-frequency coefficients that should have been suppressed remain relatively large magnitudes. It can be observed from both Figure 4 and its 0% sparsity rate that the Gamma hyperprior in this case completely destroys the sparse structure with a positive  $\gamma$ , but it can still stabilize the problem. The case  $\alpha = 1$ , corresponding to the half-Laplace hyperprior, has been analyzed in previous tests.

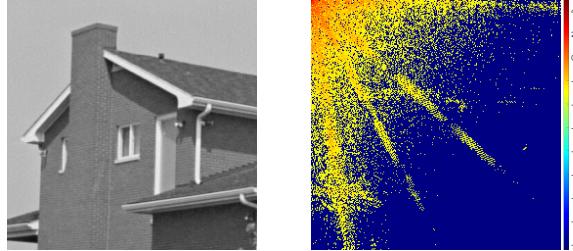


FIG. 5. The compressed *House* image and its DCT coefficients.

**5.2. Influence of ill-posedness and noise.** In subsection 5.1, we have compared the influence of different hyperpriors applied in EBF on sparsity and restoration. In this subsection, we will use the half-Laplace hyperprior as an example to study the performance of EBF under different levels of ill-posedness and noise. In our tests, besides *Cameraman* we include 256-by-256 gray image *House* as an extra test image. To ensure the sparsity in its DCT coefficients, we pre-process the image similarly by truncating its DCT coefficients with the threshold 0.025. It gives the sparsity rate as 75.18%, and the relative error between the compressed and original images is 0.0157. In Figure 5 we show the compressed *House* image with its DCT coefficient map.

**5.2.1. Ill-posedness.** In all tests, we consider the Gaussian blur. The standard deviation  $\sigma_{ker}$  of the Gaussian blurring kernel controls the ill-posedness of the inverse problem. A larger  $\sigma_{ker}$  results in stronger smoothing and faster decay of the singular values of the forward operator, thereby making the inverse problem more ill-posed. In this test, we vary the ill-posedness of the problem by adopting different values of  $\sigma_{ker}$ . We set  $\sigma_{ker} = 0.5, 1$  and  $1.5$ . In the half-Laplace hyperprior we still set  $\beta$  as 0.1. To reduce the interference of noise and ensure that the differences in performance in the tests are mainly caused by the variation in ill-posedness, we set a rather low noise level as 5%.

The restored results are shown in Figure 6. We can see that in a rather well-posed case with  $\sigma_{ker} = 0.5$ , EBF provides accurate restoration with coefficients concentrated in the low-frequency region forming a distinct fan-shaped pattern to preserve most details. As the level of ill-posedness increases, the sector-shaped areas of the DCT coefficient maps shrink. That is because the stronger the ill-posedness of the problem, the more severely the high-frequency information is corrupted by noise. Therefore, the sparsity of the restored coefficients gradually increases. Furthermore, because more and more high-frequency coefficients are suppressed to zero, the restored image loses more details as  $\sigma_{ker}$  increases. In Table 3, we list the relative errors of the restored images and the sparsity rates of the DCT coefficients. It is clear that the sparsity rate of the restored result increases as the problem becomes more ill-posed due to amplified noise on the high-frequency. Moreover, more ill-posed problem is more challenging to solve, leading to larger relative error. We observe that for the *House* test image the relative error in the case of  $\sigma_{ker} = 1.5$  is slightly smaller than the case of  $\sigma_{ker} = 1.0$ . This is due to the trade-off between suppressing noise and preserving details.

**5.2.2. Noise.** In this test, we fix  $\sigma_{ker} = 1$  and vary the noise level to assess the performance of EBF with a half-Laplace hyperprior ( $\beta = 0.1$ ). Specifically, we test three levels of Gaussian noise: 5%, 10%, and 20%.

The restored results are given in Figure 7. It shows that as the noise level in-

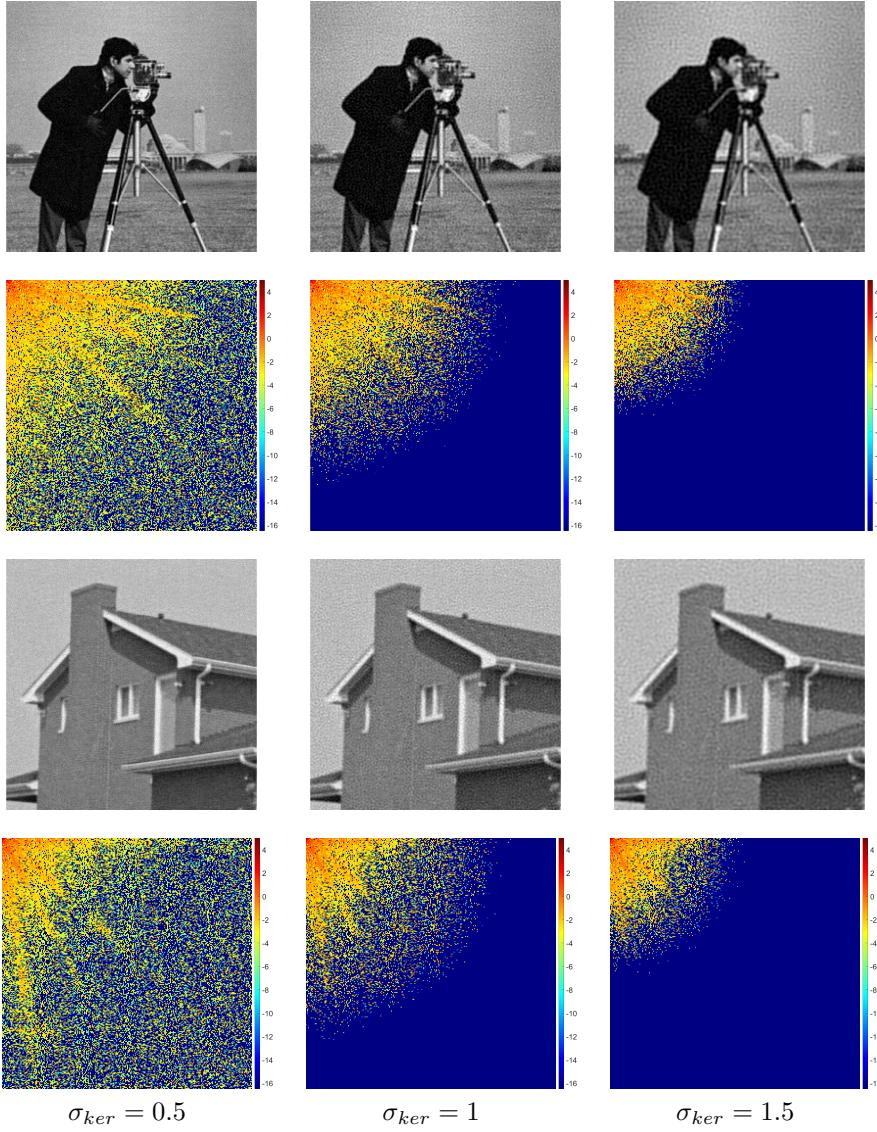


FIG. 6. The results of EBF with the Laplace hyperprior for restoring images degraded by different Gaussian blurring kernels. Top and third rows: restored images; Second and fourth rows: absolute values of DCT coefficients.

creases, the sparsity of the restored DCT coefficients gradually increases and more high-frequency coefficients are turned to zero. In low noise level cases, more the high-frequency DCT coefficients contribute to the restoration, therefore the restored images preserve most of the details. However, in high noise level cases, in order to suppress the influence of the noise most of high-frequency coefficients become zero, so some details in the restored image are lost and some compression artifacts appear. In Table 3 we list the relative errors of the restored images and the sparsity rates of the DCT coefficients. It is obvious that the relative error and the sparsity rate increase as the noise level increases. It confirms that EBF with a sparsity-promotting hyperprior

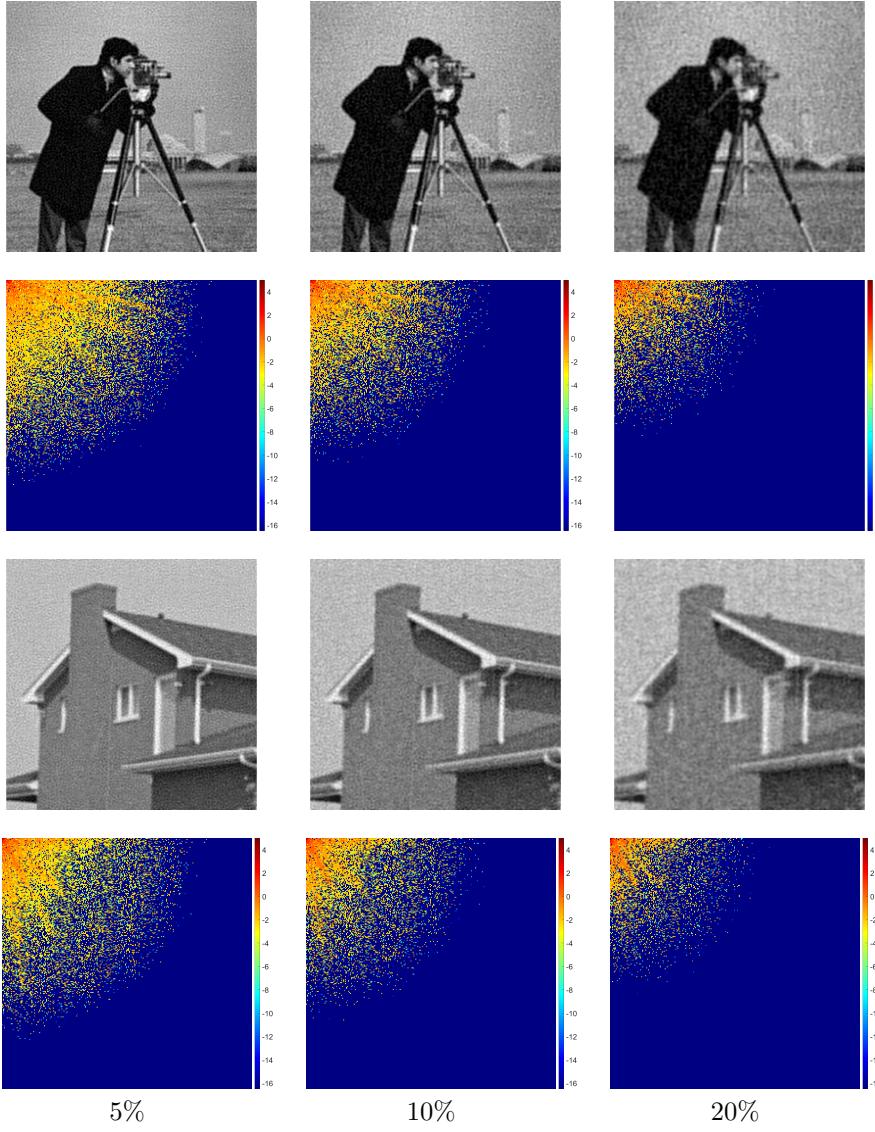


FIG. 7. The results of EBF with the half-Laplace hyperprior for restoring images degraded by different noise levels. Top and third rows: restored images; Second and fourth rows: absolute values of DCT coefficients.

can automatically adjust the sparsity with respect to the noise level to obtain a good balance between suppressing noise and preserving details.

**5.3. Convergence of PALM.** In the last test, we study the convergence of Algorithm 4.1 numerically. Here, we use EBF equipped with a half-Laplace hyperprior to solve the image deblurring problem with 10% noise and  $\sigma_{ker} = 1$ . We test our PALM algorithm with different choices of  $\beta$  in the half-Laplace hyperprior. The convergence plots are shown in Figure 8, where we plot the values of the objective function defined in (4.5) in terms of iterations. It is obvious that the objective function values decrease monotonically, and the algorithm converges. Furthermore, the algo-

TABLE 3

Restoration performance under different levels of ill-posedness (with 5% noise) and different noise levels (with  $\sigma_{ker} = 0.1$ ).

Test type	Image	Parameter	Relative error	Sparsity (%)
Ill-posedness	<i>Cameraman</i>	$\sigma_{ker}$	0.5	0.0627
			1.0	0.0909
			1.5	0.1006
	<i>House</i>	$\sigma_{ker}$	0.5	0.0550
			1.0	0.0684
			1.5	0.0627
Noise level	<i>Cameraman</i>	Noise level (%)	5	0.0909
			10	0.1055
			20	0.1273
	<i>House</i>	Noise level (%)	5	0.0684
			10	0.0751
			20	0.0870

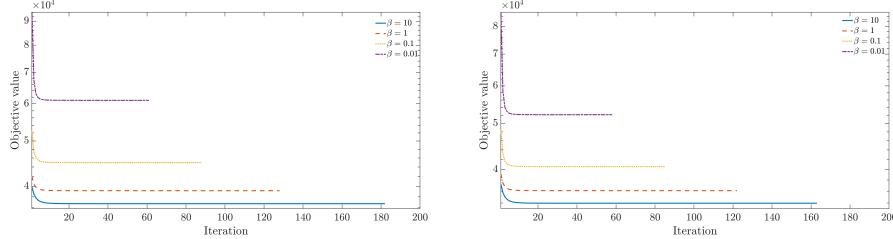


FIG. 8. The convergence results of PALM in Algorithm 4.1. Left: *Cameraman*; Right: *House*.

rithm converges faster with smaller  $\beta$ , since in this case the half-Laplace hyperprior is closer to the  $\delta$  distribution and dominates the posterior.

**6. Conclusions.** This work investigated the fundamental question of how different hyperpriors influence sparsity and stability within the empirical Bayes framework. We provided a comprehensive theoretical analysis that study the choice of the hyperprior to both the sparsity and the local optimality of the resulting solutions. In particular, we demonstrated that certain hyperpriors, such as the half-Laplace prior and half-generalized Gaussian prior with shape parameter  $0 < p < 1$ , can not only promote sparsity in the estimated hyperparameters, but also stabilize the solution with respect to measurement noise. To address the resulting nonconvex optimization problem, we proposed a PALM algorithm and established its convergence under both convex and concave hyperpriors. Extensive numerical experiments on 2D image deblurring tasks were conducted to validate our theoretical findings. The results confirm that the choice of hyperprior plays a decisive role in balancing sparsity, stability, and restoration accuracy, especially under increasing levels of ill-posedness and noise.

Several directions remain open for future work. First, it is important to extend the convergence theory of PALM beyond the convex or concave setting, to accommodate more general nonconvex hyperpriors. Second, we plan to investigate the sparsity-promoting properties of broader classes of group or heavy-tailed hyperpriors, such as the Student's t prior. Finally, to further improve image restoration quality, we will explore alternative sparsifying transforms beyond DCT, and develop fast solvers tailored to these bases within EBF.

## REFERENCES

- [1] A. AL HILLI, L. NAJAFIZADEH, AND A. PETROPULU, *Weighted sparse Bayesian learning (WSBL) for basis selection in linear underdetermined systems*, IEEE Trans. Veh. Technol., 68 (2019), pp. 7353–7367.
- [2] C. M. BISHOP AND N. M. NASRABADI, *Pattern Recognition and Machine Learning*, vol. 4, Springer, Berlin, 2006.
- [3] J. BOLTE, S. SABACH, AND M. TEBOULLE, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Math. Program., 146 (2014), pp. 459–494.
- [4] D. CALVETTI, M. PRAGLIOLA, AND E. SOMERSALO, *Sparsity promoting hybrid solvers for hierarchical Bayesian inverse problems*, SIAM J. Sci. Comput., 42 (2020), pp. A3761–A3784.
- [5] D. CALVETTI, M. PRAGLIOLA, E. SOMERSALO, AND A. STRANG, *Sparse reconstructions from few noisy data: analysis of hierarchical Bayesian models with generalized gamma hyperpriors*, Inverse Problems, 36 (2020), p. 025010.
- [6] D. CALVETTI AND E. SOMERSALO, *Hypermodels in the Bayesian imaging framework*, Inverse Problems, 24 (2008), p. 034013.
- [7] D. CALVETTI, E. SOMERSALO, AND A. STRANG, *Hierarchical Bayesian models and sparsity:  $\ell_2$ -magic*, Inverse Problems, 35 (2019), p. 035003.
- [8] S. S. CHEN, D. L. DONOHO, AND M. A. SAUNDERS, *Atomic decomposition by basis pursuit*, SIAM Rev., 43 (2001), pp. 129–159.
- [9] W. DENG, X. ZHANG, F. LIANG, AND G. LIN, *An adaptive empirical Bayesian method for sparse deep learning*, in Advances in Neural Information Processing Systems, vol. 32, Curran Associates, Inc., 2019.
- [10] A. FAUL AND M. TIPPING, *Analysis of sparse bayesian learning*, in Advances in Neural Information Processing Systems, vol. 14, MIT Press, 2001.
- [11] D. GHILLI, D. A. LORENZ, AND E. RESMERITA, *Nonconvex flexible sparsity regularization: theory and monotone numerical schemes*, Optimization, 71 (2022), pp. 1117–1149.
- [12] J. GLAUBITZ, A. GELB, AND G. SONG, *Generalized sparse Bayesian learning and application to image reconstruction*, SIAM/ASA J. Uncertain. Quantif., 11 (2023), pp. 262–284.
- [13] M. GRASMAIR, M. HALTMEIER, AND O. SCHERZER, *Sparse regularization with  $\ell_q$  penalty term*, Inverse Problems, 24 (2008), p. 055020.
- [14] X. GU, H. LEUNG, AND X. GU, *Bayesian sparse estimation using double lomax priors*, Math. Probl. Eng., 2013 (2013), pp. 1–17.
- [15] P. C. HANSEN, *Discrete Inverse Problems: Insight and Algorithms*, SIAM, Philadelphia, 2010.
- [16] N. J. HIGHAM, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 2002.
- [17] L. HOGBEN, *Handbook of Linear Algebra*, CRC press, Boca Raton, 2006.
- [18] D. A. LORENZ AND E. RESMERITA, *Flexible sparse regularization*, Inverse Problems, 33 (2016), p. 014002.
- [19] D. J. MACKAY, *Bayesian interpolation*, Neural Comput., 4 (1992), pp. 415–447.
- [20] D. J. MACKAY, *The evidence framework applied to classification networks*, Neural Comput., 4 (1992), pp. 720–736.
- [21] J. R. MAGNUS AND H. NEUDECKER, *Matrix Differential Calculus with Applications in Statistics and Econometrics*, John Wiley & Sons, Chichester, 2019.
- [22] S. NABI, H. NASSIF, J. HONG, H. MAMANI, AND G. IMBENS, *Bayesian meta-prior learning using empirical Bayes*, Manage. Sci., 68 (2022), pp. 1737–1755.
- [23] B. K. NATARAJAN, *Sparse approximate solutions to linear systems*, SIAM J. Comput., 24 (1995), pp. 227–234.
- [24] M. K. NG, R. H. CHAN, AND W.-C. TANG, *A fast algorithm for deblurring models with Neumann boundary conditions*, SIAM J. Sci. Comput., 21 (1999), pp. 851–866.
- [25] M. PRAGLIOLA AND Y. DONG, *Inducing sparsity via the horseshoe prior in imaging problems*, Inverse Problems, 39 (2023), p. 074001.
- [26] S. H. RUDY AND T. P. SAPSIS, *Sparse methods for automatic relevance determination*, Phys.

- D, 418 (2021), p. 132843.
- [27] R. TIBSHIRANI, *Regression shrinkage and selection via the Lasso: a retrospective*, J. R. Stat. Soc. Ser. B. Stat. Methodol., 73 (2011), pp. 273–282.
  - [28] M. E. TIPPING, *Sparse Bayesian learning and the relevance vector machine*, J. Mach. Learn. Res., 1 (2001), pp. 211–244.
  - [29] F. URIBE, Y. DONG, AND P. C. HANSEN, *Horseshoe priors for edge-preserving linear Bayesian inversion*, SIAM J. Sci. Comput., 45 (2023), pp. 337–365.
  - [30] M. WANG, W. XU, AND A. TANG, *On the performance of sparse recovery via  $\ell_p$ -minimization ( $0 \leq p \leq 1$ )*, IEEE Trans. Inform. Theory, 57 (2011), pp. 7255–7278.
  - [31] D. WIPF AND S. NAGARAJAN, *A new view of automatic relevance determination*, in Advances in Neural Information Processing Systems, vol. 20, Curran Associates, Inc., 2007.
  - [32] D. WIPF AND S. NAGARAJAN, *A unified Bayesian framework for MEG/EEG source imaging*, NeuroImage, 44 (2009), pp. 947–966.
  - [33] D. P. WIPF AND B. D. RAO, *Sparse Bayesian learning for basis selection*, IEEE Trans. Signal Process., 52 (2004), pp. 2153–2164.
  - [34] D. P. WIPF AND B. D. RAO, *An empirical Bayesian strategy for solving the simultaneous sparse approximation problem*, IEEE Trans. Signal Process., 55 (2007), pp. 3704–3716.
  - [35] D. P. WIPF, B. D. RAO, AND S. NAGARAJAN, *Latent variable Bayesian models for promoting sparsity*, IEEE Trans. Inform. Theory, 57 (2011), pp. 6236–6255.
  - [36] Z. XIE AND J. HU, *Reweighted  $\ell_1$ -minimization for sparse solutions to underdetermined linear systems*, in 2013 6th International Congress on Image and Signal Processing (CISP), vol. 3, IEEE, 2013, pp. 1660–1664.
  - [37] L. YANG, J. FANG, H. CHENG, AND H. LI, *Sparse Bayesian dictionary learning with a Gaussian hierarchical model*, Signal Process., 130 (2017), pp. 93–104.
  - [38] F. YU, L. SHEN, AND G. SONG, *Hyperparameter estimation for sparse Bayesian learning models*, SIAM/ASA J. Uncertain. Quantif., 12 (2024), pp. 759–787.